

STATISTICS

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

Data: "facts or pieces of information"

Ex: Height of students in a classroom = {175 cm, 150 cm, 160 cm, 155 cm}

↳ Intelligence Quotient (IQ) of 5 randomly selected individuals.
= (109, 89, 129, 101, 105)

Two types of statistics:

Descriptive stats

Inferential stats

Descriptive Statistics

- It consists of organizing and summarizing of data
- Ex: Probability Density Function (PDF), Histogram, Box plot, Bar chart, Pie chart

Inferential Statistics

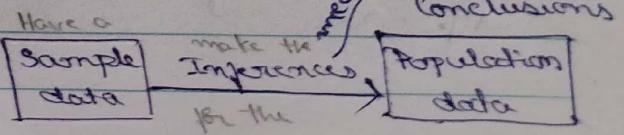
- It consists of using data that you've measured to form conclusions.
- Ex: Hypothesis Testing, P-value, Z-test, t-test, Anova test, chi-square test

Example: Let's say there are 20 maths classes at your university and you've collected the ages of students in one class

Ages = {21, 20, 18, 34, 17, 22, 24, 25, 26, 23, 22, 20}

- What is the average (or mode, median or min) age of student in this maths class similar to what you would expect in a general maths class at this university?

Are the ages of the students in this maths class similar to what you would expect in a general maths class at this university?

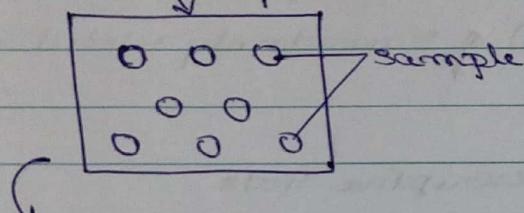


Population and Sample data

Elections:

working with
 • Population and Sample data
 is very much important
 for statistics
inferential

Population of State



- There are different ways of sampling too.

News channels can not ask each and every person, as to whom they have cast their votes to... for prediction

So, news channel pick up some sample of people to know to whom they have voted

Let total no. of ppl selected for sample is 1000.

Parties for {BJP, Congress, Independent election 60%. 30%. 20%.

⇒ Exit Poll Results

Example: Have to take up the task of ordering the jacket for the employees of a company (Considering 4 branches at different locations) 60k - 70k employees approx.

- It is not possible to ask each and every employee the size.
- Let's take sample data & make some assumptions.

40% - L

30% - XL

20% - S

5% - XS

4% - XXL

1% - XXXL

} with this only 1-2% of jackets were wasted (like not of correct size for ppl)

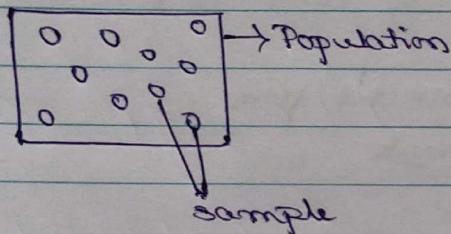
Population (N)

Sample (n)

Sampling Techniques

1) Simple Random Sampling

Every member of the population (N) has an equal chance of being selected for your sample (n)



2) Stratified Sampling

Strata (means) → Layers
also

Non overlapping groups
↓
clusters

- Ex: • Gender
 - Male
 - Female
- Age groups
 - 10-19
 - 19-35
 - 35-60

- Blood groups
- Courses
- Education Qualification

3) Systematic Sampling

In a population
of (N) → Select every n^{th} individual

Ex: Survey → Mall (SBI credit card)

Airport → Customs (Checking luggage)

4) Convenience Sampling

Only those people who are interested will only be participating

Ex: Data Science Survey $\xrightarrow[\text{about 1/2}]{\text{Sampling only to ppl who know anyone interested in it}}$ Artificial Intelligence

- Youtube survey

- Healthcare Disease

Blind people Survey \rightarrow Only blind ppl participate in it

- Economics - RBI

House Hold Survey \rightarrow Info collected from Female

Exit Poll: Stratified + Random Sampling Sampling.

Variable:

A variable is a property that can take on any value.

Ex: Height = 182

$$a = 150$$

$$b = 145$$

$$xy = 160$$

$$A = \{7, 9, 10\}$$

There are 2 kinds of variable

Quantitative variable
Qualitative variable

1) Quantitative variable:

Can be measured numerically $\left\{ \begin{array}{l} \text{e.g., can also perform major operation like} \\ \text{Add, Sub, } \times, \div \end{array} \right\}$

→ Qualitative Variable: Based on some characteristics we can derive categorical variables

Ex: Gender ↗ Male

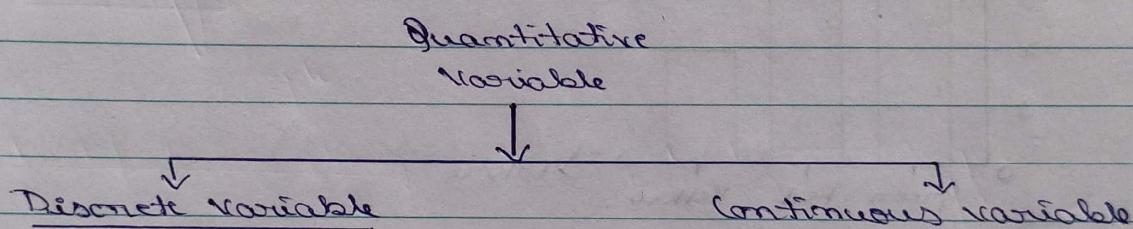
↘ Female

- Can we convert a Quantitative variable → Qualitative variable?
Yes it is possible.

Ex: IQ

0 - 10	10 - 50	50 - 100
Low IQ	Medium IQ	High IQ

→ 3 categories
(Qualitative)



Ex: • Whole numbers

• Number of Bank Accounts
= {2, 3, 4, 7} ↗

$$\begin{matrix} 3.5x \\ 1.5x \end{matrix}$$

• Total no. of children in a family = {2, 3, 7, 10}

• Total no. of employees in a company = {10k, 30k, 50k}

Ex: • Height = {172.5, 165.5, 130.78} cm

• Rainfall = 1.35, 1.25, 1.75 cm

• Weight of a person

• Temperature

• Measurement of ingredients

Questions:

What kind of variable does it belong to;

[a] Marital Status ↗ Categorical (Qualitative)

[b] Nile River length } Continuous quantitative

[c] Movie duration } variable

[d] IQ

Frequency Distribution

Consider

Sample dataset : Green, Red, Yellow, Green, Red, Yellow, Green, Red
WRT colour

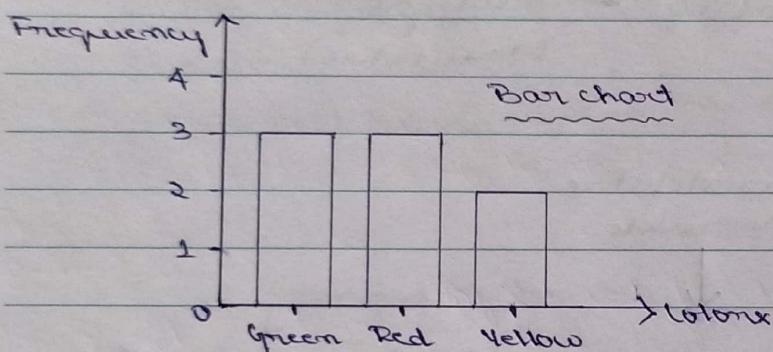
Color	Frequency
Green	3
Red	3
Yellow	2

- Bar graph - categorical variable

- Histogram - continuous variable

↳ Categorical variable.

BAR GRAPH



4: Variable Measurement Scales

- 4 types of measured variables.

→ Ranking is not that important

- 1) Nominal data : {Categorical data}

Ex: Colors, Gender, Flower types

- 2) Ordinal data : {Categorical data}

→ Ranks can be assigned

to it

Ex: [] Student (Marks)

Rank

100	1	{ } ordinal data
96	2	
85	3	
57	4	
45	5	

Tenatiles

Ranking will definitely help us to understand things in a better way.

M T W T F S S

--	--	--

[B]	Degree	Salary
	PHD	(1)
	B.E	(2)
	Masters	(3)
	BCA	(4)
	12 th	(5)

3) Interval data

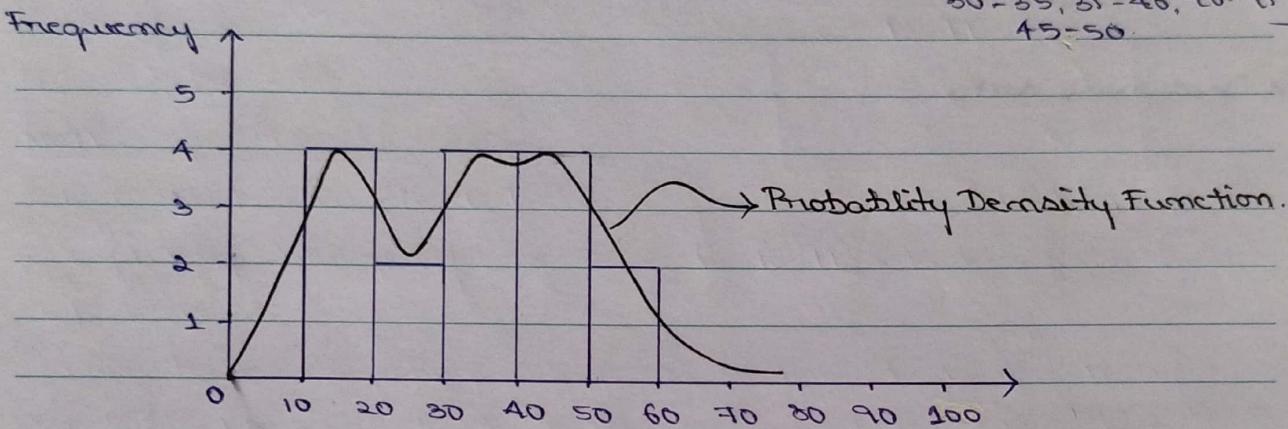
Bins: The number of parts an interval/range is divided into. F S S

Histograms [continuous variables]

Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51}

Bins = 10 (Default) \Rightarrow x-axis
 ↳ group/range size

Ex: 0-50 $\xrightarrow{\text{bins}}$ 0-5, 5-10, 10-15,
 15-20, 20-25, 25-30,
 30-35, 35-40, 40-45
 45-50.



Ages	Frequency
0-10	0
10-20	4
20-30	2
30-40	4
40-50	4
50-60	2

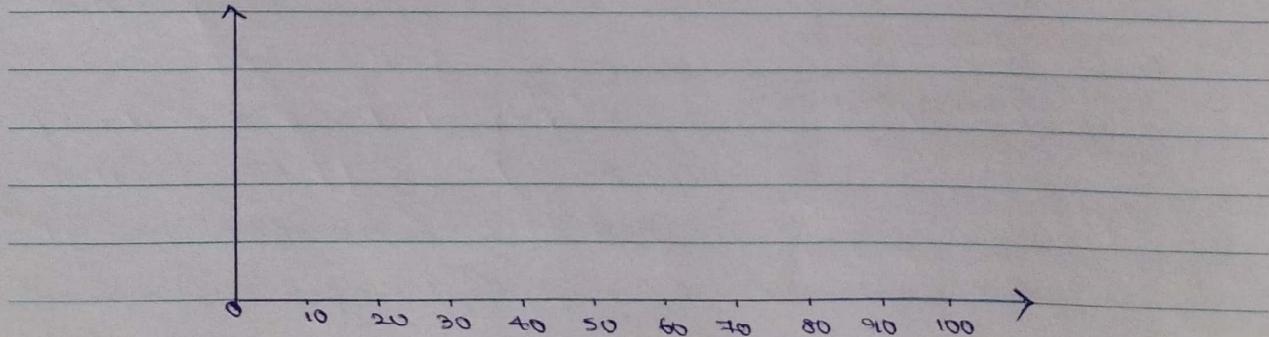
Ages \Rightarrow a-b (excluding b)

- Probability Density Function (PDF)
 It is the smoothed version
 of Histogram
- kernel Density Estimator
 used to smoothen histogram.

Example: 10, 13, 18, 22, 27, 32, 38, 40, 45, 51, 56, 57,
 88, 90, 92, 94, 99

Create a histogram with total bins as 10.

- 0-10, 10-20, 20-30, 30-40, 40-50, 50-60, 60-70, 70-80, 80-90, 90-100



Intermediate Stats

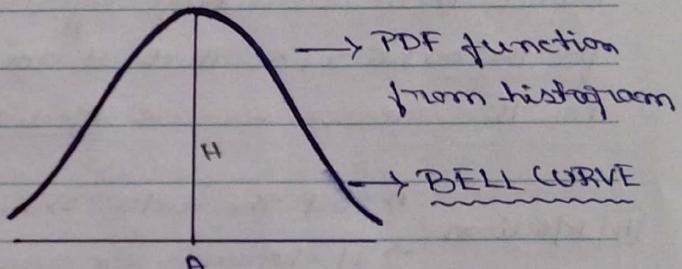
- 1) Measure of Central Tendency
- 2) Measure of Dispersion
- 3) Gaussian Distribution
- 4) Z-Score
- 5) Standard Normal Distribution
- 6) Central Limit Theorem

1) Measure of Central Tendency

↳ Basically identifies the central position of the dataset

- [a] Mean
- [b] Median
- [c] Mode

• Most of the elements are present in the central regions only near A.



• A can be mean/median/mode

H determines in what region the curve is at the top

• Let's consider:

Population (N)

$$\mathcal{X} = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$$

Sample (n)

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

$$\text{Population mean } (\mu) = \frac{\sum_{i=1}^N x_i}{N}$$

where, N = size of dataset

x_i = elements in dataset.

$$\mu = \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$\mu = \frac{32}{10} = 3.2$$

Outlier: M T W T F S S

It is an element that looks completely different than that of the entire distribution.

[a] Mean

Let $x = \{1, 2, 2, 3, 4, 5\}$
↳ sample

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

$$\bar{x} = \frac{1+2+2+3+4+5}{6} = \frac{17}{6}$$

$$\bar{x} = 2.83$$

- adding outlier to x

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

$$\bar{x} = \frac{1+2+2+3+4+5+100}{7} = \frac{117}{7}$$

$$\bar{x} = 16.71$$

NOTE:

- We can observe that because of an outlier addition, mean completely changes from 2.83 to 16.71, which gives a very huge difference.
- Since there is a huge difference, so have to do something for the outlier, so that it does not affect the entire distribution.
For this reason we use Median

[b] Median

1) Sort the data. 2) Pickup the central element from data odd.
3) If elements are even taken 2 elements $\frac{a+b}{2}$

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

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

(outlier)

$$\text{Median} = \frac{2+3}{2} = \frac{5}{2} = 2.5 \approx 2.83$$

(Mean)

$$\text{Median} = 3 \approx 2.83$$

(Mean)

- Median definitely works in the case of outliers

[c] Mode - Element with highest frequency is chosen

$$x = \{1, 2, 2, 3, 3, 3, 4, 5, 6, 6, 7\}$$

$$\text{Mode} = 3$$

Example:

A girl went to a market and picked up some flowers.

→ Lilly, Sunflower, Rose, ..., Lilly
Lotus, NAN, NAN

• NAN is replaced by frequently picked up flower. Lilly here.

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

$$\text{Mode} = [2, 3, 4]$$
 is soln in python

If such scenario, use Median occurs

Usage of Mode:

In EDA → Feature Engineering

If it has [in the dataset]

NAN values
(NULL)

Categorical variable

then NAN is
replaced
using

Mode

If the data set has

continuous values

then we can replace
NAN values with

Mean

continuous + outliers
values

then the technique
used to replace NAN
values

Medium

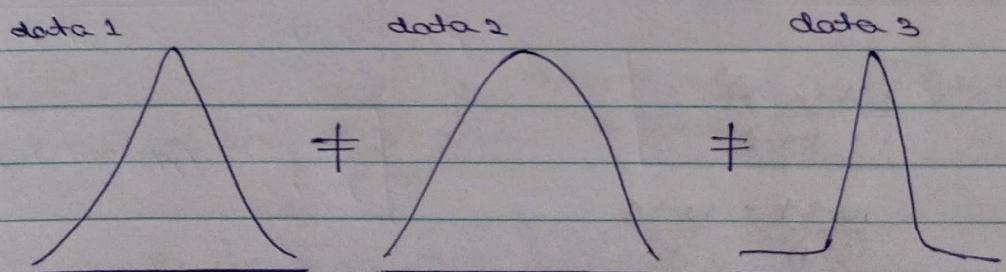
2) Measure of Dispersion

↳ Spread {How the data is spread}

a) Variance

b) Standard Deviation

- Consider a specific data (here data 1, data 2, data 3) giving different spread for the data respectively, which is measured using a) & b)



With the help of these different spreadness, we can understand something about data

a) Variance

It basically says how well the data is spread

We defined Variance based on 2 things

Population Variance (σ^2) ^{sigma square}

^{Population Variance}
Sample Variance
Sample Variance (s^2)

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

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

where,

μ = population mean

where

\bar{x} = sample mean

Example: $x = \{1, 2, 2, 3, 4, 5\} \Rightarrow \bar{x} = 2.83$

Q) Why sample variance
is divided by $(n-1)$?

$n = 6$

x	\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
10.84			

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

$$= \frac{10.84}{(6-1)} = \frac{10.84}{5}$$

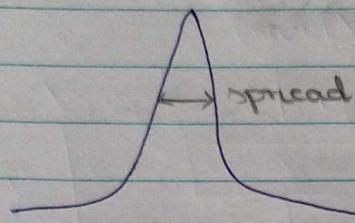
$$\therefore [s^2 = 2.168]$$

Now, let's consider, variance (defines the spread)

Variance spread

- From previous example

$$\text{Variance} = 2.168$$

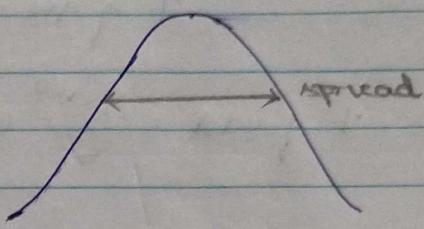


Mumbai (less space)

Vertical growth (Tall buildings)

- Random consideration

$$\text{Variance} = 6.45$$



Bangalore - villages to cities

Horizontal growth

(Outskirts of Bangalore expanded to now comprising entire Bangalore)

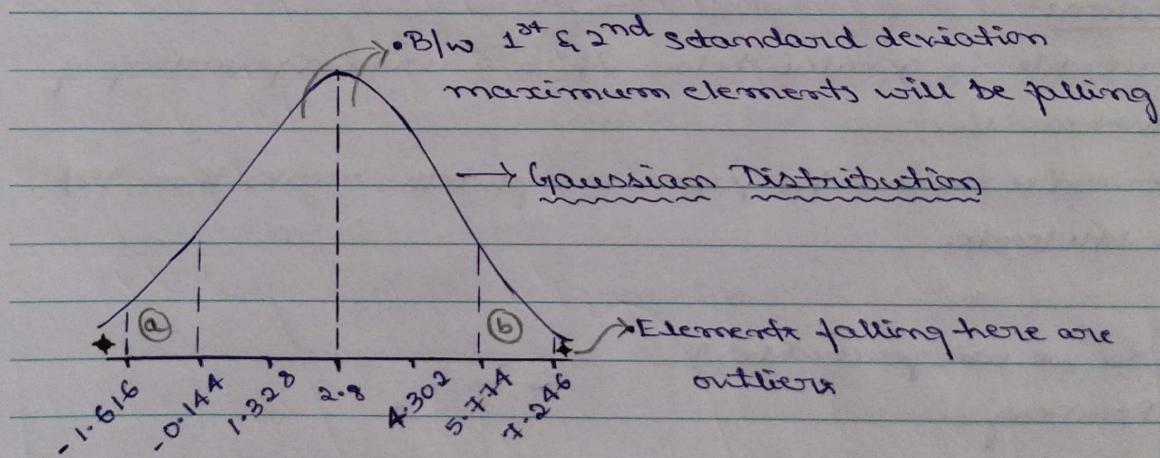
b) Standard Deviation (-)

$$\sigma = \sqrt{\text{Variance}}$$

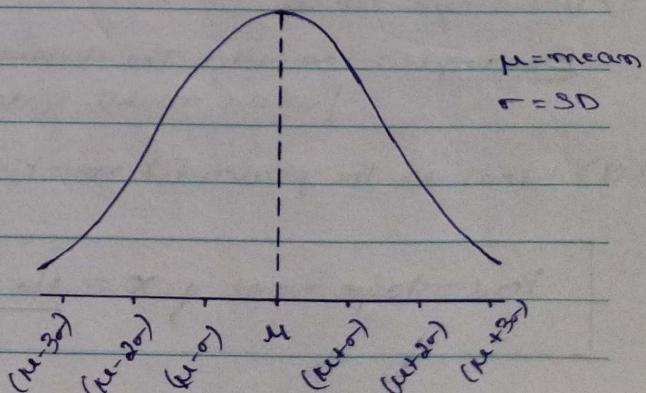
↳ Granular knowledge $\rightarrow \frac{1}{n}$

Example: $x = \{1, 2, 2, 3, 4, 5\}$, $s^2 = 2.168$

$$\sigma = \sqrt{2.168} = 1.472$$



Regions (a) & (b) minimum elements will be falling



Concepts:

PERCENTILES AND QUARTILES

Percentages:

Example: 1, 2, 3, 4, 5.

∴ % of the numbers that are odd?

$$\% \text{ of odd} = \frac{3}{5} \times 100 = 60\%$$

Percentiles:

Exams like CAT, GATE, SAT here the setting happens in terms of percentile

99 percentile ≠ 99%.

↓ it means

Out of all the people that have written that exam, his/her rank is somewhere near this particular (99) value.

Definition:

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

99 percentile means the person has got better marks than 99% of the students.

Why do we use percentiles?

1) For Ranking

2) used in Stats

Example: Consider the dataset {2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 10, 11, 12}
[where n=20, & data already sorted]

Q) What is the percentile ranking of 10

$$\text{Percentile rank of } x = \frac{\text{No. of values below } x}{n} \times 100$$

$$\text{Percentile Rank} = \frac{16}{20} \times 100 = 80 \text{ percentile}$$

of 10

Indicates that 10 is greater than 80% of entire dataset distribution

28) What is the percentile rank of 11?

$$\text{Percentile Rank} = \frac{17}{20} \times 100 = 85 \text{ percentile}$$

of 11

30) What value exists at percentile ranking of 25%?

$$\text{Value} = \frac{\text{Percentile} \times (\text{mti})}{100}$$

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

A

= 5.25 → index not value.

∴ dataset: $\frac{1}{2}, \frac{2}{2}, \frac{3}{2}, \frac{4}{2}, \frac{5+6}{2}, 5, 6, \dots, 12$

↳ combine them & find the average
as index 5.25 is not a whole no.

$$\text{Average} = \frac{5+5}{2} = \frac{10}{2} = 5$$

Value = 5 for 25 percentile

- These all things are for understanding purpose only no need of calculation as code will automatically do it.
- Formulas derivation are not necessary.

Quartiles

Quartiles basically mean 25 percentile

M T W T F S S

FIVE NUMBER SUMMARY

- 1) Minimum Value
- 2) First Quartile Value (25 percentile) Q_1
- 3) Median (50 percentile) (Measure of central tendency)
- 4) Third Quartile Value (75 percentile) Q_3
- 5) Maximum Value

Removing the outliers

Dataset: {1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27}

- How do we get to know some element is an outlier?

We have to define the fences so that i.e.,

Lower Fence \leftrightarrow Higher Fence.

Any element ^{that} does not belong to the fence is considered as outliers

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

$$\text{Higher Fence} = Q_3 + 1.5(\text{IQR})$$

IQR = Inter Quartile Range

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

$$Q_1 = \frac{25}{100} \times (n+1)$$

$$Q_3 = \frac{75}{100} \times (n+1)$$

For given dataset: $n=19$

$$Q_1 = \frac{25}{100} \times (20) = 5^{\text{th}} \text{ index}$$

$$Q_1 = 3$$

$$Q_3 = \frac{75}{100} \times (20) = 15^{\text{th}} \text{ index}$$

$$Q_3 = 13$$

$$\text{IQR} = 13 - 3$$

$$\boxed{\text{IQR} = 10}$$

$$\text{Lower Fence} = 3 - 1.5(10) = -3$$

$$\text{Higher Fence} = 13 + 1.5(10) = 23$$

$$-3 \longleftrightarrow 23$$

$\therefore 27$ does not come into this fence, hence it is a outlier.

→ Found after excluding outliers
Median = $\frac{5+5}{2} = 5$

M	T	W	T	F	S	S

- Five number summary
(For the dataset)

Box plot

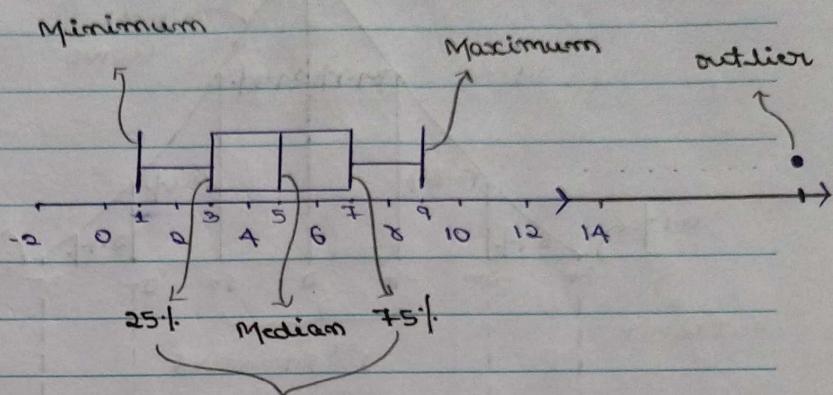
Minimum = 1

$Q_1 = 3$

Median = 5

$Q_3 = 7$

Maximum = 9



Use of Box Plot

IQR

- It is a visualization technique to see the outlier.

3) Distributions

[a] Normal / Gaussian Distribution

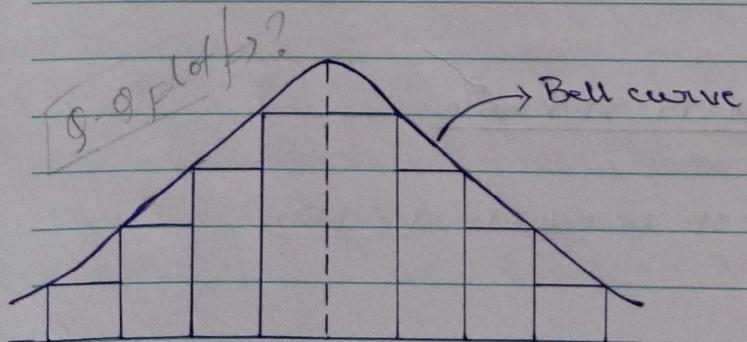
[b] Standard Normal Distribution

[c] I-Score

[d] Log Normal Distribution

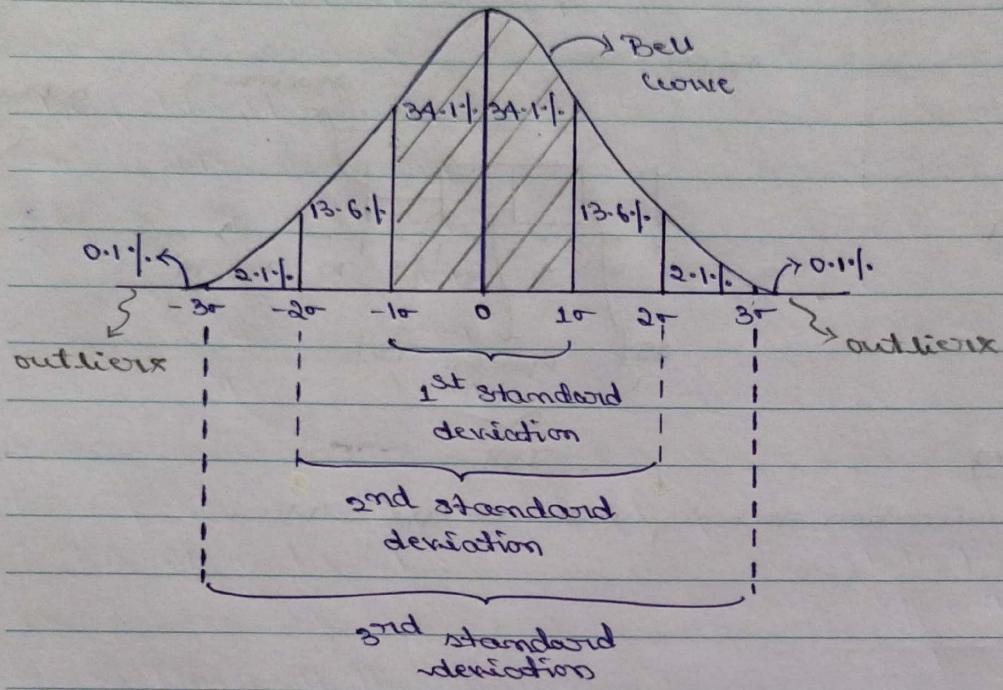
[e] Bernoulli's Distribution } No need of
[f] Binomial Distribution } formulas to
[g] Power law } remember.
[h] Normal / Gaussian Distribution

It is nothing but the bell curve (as seen before) constructed based on mean / median / mode and the standard deviation.



Properties:

- Empirical Rule of Gaussian Distribution



→ Area b/w 1st standard deviation (to the right & left)

$$34.1\% + 34.1\% = 68.2\%$$

Most of the elements are present over here i.e., around 68.2%.

→ Area b/w 2nd standard deviation (to the right & left)

$$68.2\% + 13.6\% + 13.6\% = 95.4\%$$

Here around 95.4% of the entire dataset will be falling in 2nd SD

→ Area b/w 3rd standard deviation (to the right & left)

$$95.4\% + 2.1\% + 2.1\% = 99.7\%$$

Here around 99.7% of data will be found

Empirical Rule:

This Rule is known as 68-95-99.7 rule

→ Anything falling after 3rd SD is considered outlier.

- Which all kind of dataset will fall into Gaussian distribution (bell curve)?
- 1) IRIS Dataset
 2) Height of human being
 3) Height } Domain expertise have confirmed them

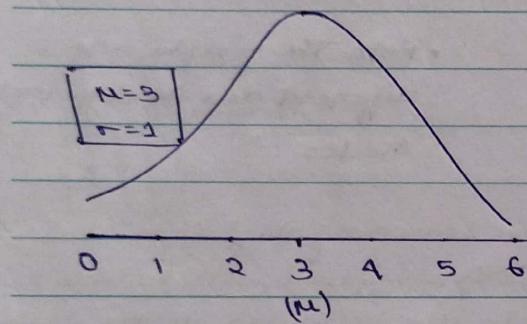
b) Standard Normal Distribution

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

$$(M) \text{Mean} = 3 \left(\frac{1+2+3+4+5}{5} \right)$$

$$(SD) \text{Standard deviation} = 1.414$$

↳ Let's roughly consider $SD \approx 1$ (to make the calculating easier)

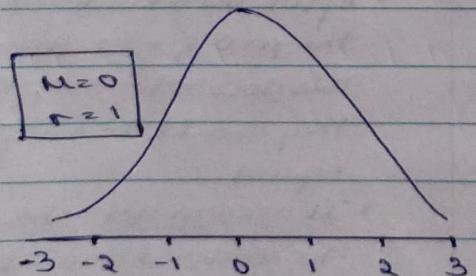


Why?

- Can we convert this entire dataset with $\mu = 3$ & $SD = 1/\sqrt{10}$ in such a way that $\mu = 0$ & $SD = 1$?
- Yes, it is possible using Z-score.

[c] Z-score {It helps us understand how far/how many SD (-) that particular no. falls away from the mean}

$$\text{Z-score} = \frac{x_i - \mu}{\sigma}$$



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

⇒ Z-score

$$\cdot \frac{1-0}{1} = 1$$

$$\cdot \frac{2-0}{1} = 2$$

$$\cdot \frac{3-0}{1} = 3$$

$$\cdot \frac{4-0}{1} = 4$$

$$\cdot \frac{5-0}{1} = 5$$

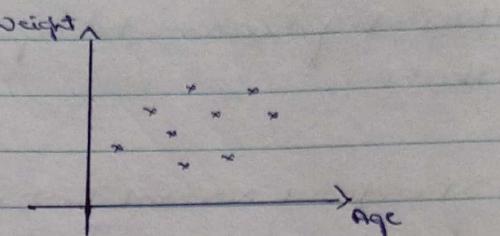
ConceptStandardization vs NormalizationStandardization:

Consider the dataset, with

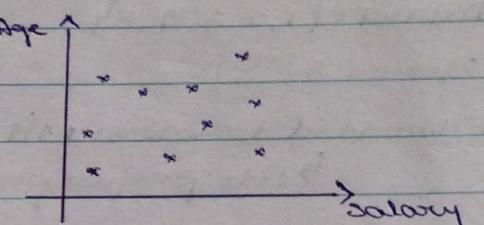
Age(yr)	Weight(kg)	Salary(INR)
25	75	35k
26	80	30k
28	85	40k
30	60	20k
32	72	35k

$$\mu = 28.2$$

weight

Scatter plot

Age



salary

- Each feature here (Age, Weight, Salary) has different unit.

Both the graphs are different due to different scales

- Can we bring all these features

into the same unit scale? [For Machine Learning Algorithm whatever Maths gets applied, will be able to calculate it quickly as all the values are in the same scale]

Yes, using Standardization applying Z-Score w.r.t each & every feature.

using z-score

Consider,

Age	μ	$x_i - \mu$	$(x_i - \mu)^2$	New Age (Standardized)
25	28.2	-3.2	10.24	-1.25
26	28.2	-2.2	4.84	-0.977
28	28.2	-0.2	0.04	-0.078
30	28.2	1.8	3.24	0.703
32	28.2	3.8	14.44	1.484

$$\mu = 28.2$$

$$\text{total} = 32.8$$

$$\left(\frac{141}{5}\right)$$

- After standardization, the $\mu = 0$ & $\sigma = 1$. However, this does not mean that they should be exactly equal.

- It means, data should be centered around zero & the spread of the data should be similar for all features.

$$\text{Variance} = \frac{\text{total}}{n} = \frac{32.8}{5} = 6.56$$

$$\text{SD}(\sigma) = \sqrt{\text{Variance}} = \sqrt{6.56} = 2.56$$

Normalization:

Here we specify the min & max value, so that values get entirely converted b/w them.

[Min Max Scalar]

$$\downarrow \downarrow$$

$$0 \rightarrow 1$$

$$0 \rightarrow 1$$

- It is done using the formula:

Min Max Scalar

$$x_{\text{Nor}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

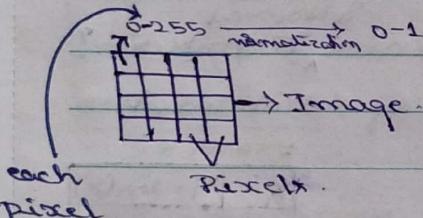
- When should we use standardization & normalization?

→ Most of the ML use case - Standardization

→ Convolution Neural Network - Normalization (CNN)

In images we have pixels

Pixels are normalized b/w 0-1



consider:

f_i	f'_i
2	0.14
5	0.54
6	0.71
8	1
1	0

calculation

$$x_{\text{Nor}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

using

After Min Max Scalar

Now there values are b/w 0 to 1

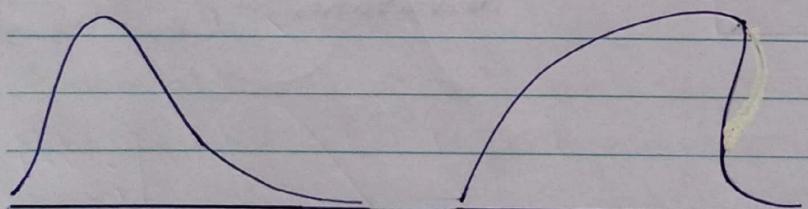
$$x_{\text{min}} = 1$$

$$x_{\text{max}} = 8$$

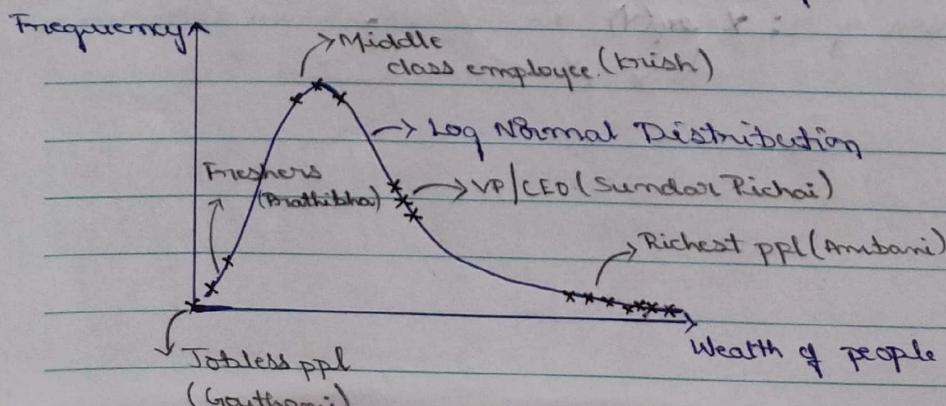
[d] Log Normal Distribution

skewed curves

They have lot of outliers



Log Normal Distribution is a type of skewed curves distribution.



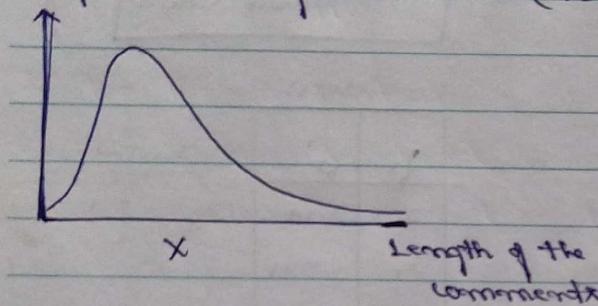
data \rightarrow histogram \rightarrow smoothen \rightarrow  \rightarrow log normal distribution M T W T F S S

Basically log Normal Distribution is saying that..

- Very less no. of people will be present in one sufficient area (richest)
- There will be max. no. of people falling in one area (middle class)
- There will be v. less no. of people in another area (Jobless, Freshers)

Another Example:

- People writing comments (length) in YouTube channel (krush)



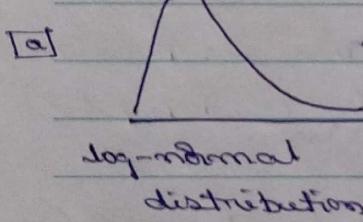
Q) Can we convert Lognormal Distribution to Gaussian Distribution?

\rightarrow Yes, if the random variable X is lognormally distributed, then

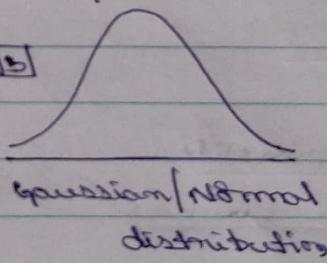
$Y = \ln(X)$ has a normal distribution.

Example.

X	$Y = \ln(x)$
20	-
25	-
30	-
40	-
45	-



when we apply
'ln' to [a] it
gives [b]



NOTE:

$\log_{10} \Rightarrow$ Base 10

$\ln \Rightarrow$ Base e

\rightarrow Equivalently, if Y has a normal distribution, then the exponential function of Y ; $X = \exp(Y)$, has a log-normal distribution.
 $\hookrightarrow X = e^Y$

Look at wikipedia for more examples.

e) Bernoulli's Distribution

- Categorical Variable - with binary outcome.

Example: Tossing a coin

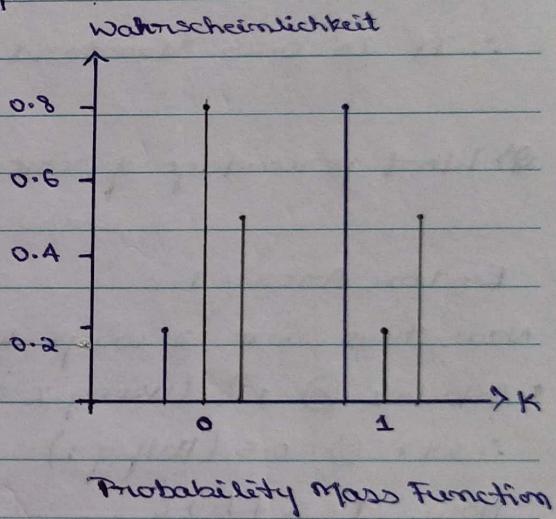
Probability of passing exam

- Bernoulli's Distribution is the discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability $q = 1 - p$.

- Probability Mass Function (PMF)

Whenever we are trying to plot for a categorical variable, at that point of time we say it as PMF

$$\text{PMF} = \begin{cases} q = 1-p & \text{if } k=0 \\ p & \text{if } k=1 \end{cases}$$



Examples of Bernoulli distribution:

- $P(x=0) = 0.2$ and $P(x=1) = 0.8$
- $P(x=0) = 0.8$ and $P(x=1) = 0.2$
- $P(x=0) = 0.5$ and $P(x=1) = 0.5$

f) Binomial Distribution (Extension of Bernoulli's Distribution)

Here we combine multiple Bernoulli Distribution

Denoted as $B(n, p)$

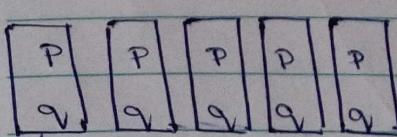
n = no. of experiments/trials, $n \in \{0, 1, 2, \dots\}$

p = probability (success probability for each trial) $p \in [0, 1]$

$q = 1 - p$; $k \in \{0, 1, \dots, n\}$ - no. of successes

Example: Tossing a coin 5 times

$$\text{PMF} = \binom{n}{k} p^k q^{n-k}$$



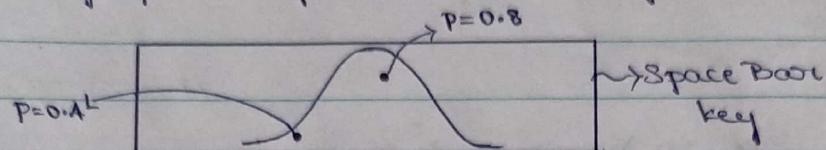
combining
multiple Bernoulli
is Binomial

Bernoulli Distribution

(Power law
80-20 Rule)

P value (Probability value)

Pressing/operating a space bar key, maximum clicks at the middle region.



- $p=0.8$, means, Out of 100 touches, the no. of touches at $p=0.8$ is 80
- $p=0.1$, means, Out of all 100 touches, the no. of times we are touching at $p=0.1$ is 10

Hypothesis Testing

Consider, Coin experiment \rightarrow Test whether the coin is fair coin or not by tossing the coin 100 times

Probability of tossing $\begin{cases} \text{Head (50 times)} \Rightarrow P(H) = 0.5 \\ \text{Tail (50 times)} \Rightarrow P(T) = 0.5 \end{cases}$ \Rightarrow Fair coin

[It is a coin that has an equal chance of landing heads or tails when it is flipped]

- Now how do we test a coin is fair or not using Hypothesis Testing?

Steps in Hypothesis Testing:

1) Null Hypothesis - (Coin is fair) $[H_0]$

By default whichever statement is true, that is taken as Null hypothesis.

Example: Criminal (Has made some crime) $\xrightarrow{\text{goes to the}} \text{Court}$

By default the court will say that criminal is not guilty until proved.

2) Alternative Hypothesis - (Coin is not fair). $[H_1]$

It should be the opposite of Null Hypothesis

3) Experiment - (Z-test, t-test, chi-square test, ANOVA test)

It is done to prove the coin is fair or not

M T W T F S S

4) Reject or Accept the Null Hypothesis.

Based on the experiments we perform this.

Now, lets test whether the coin is fair or not, 100 tosses.

Step 1:

Possibilities of tossing 100 times \Rightarrow

- 1) 50 times head $\left. \begin{array}{l} \\ \end{array} \right\}$ coin is fair
- 50 times tail

- 2) 60 times head $\left. \begin{array}{l} \\ \end{array} \right\}$ coin is fair
- 40 times tail

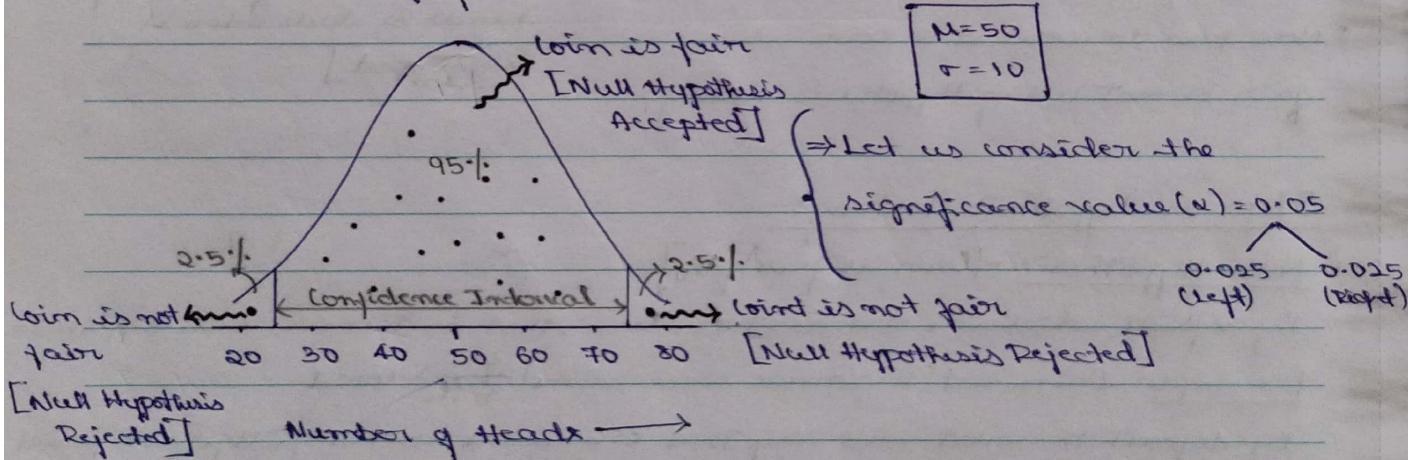
- 3) 30 times head $\left. \begin{array}{l} \\ \end{array} \right\}$ can't decide fair or not
- 70 times tail

\rightarrow So to determine this we use

[Confidence interval,
Significance value]

Confidence Interval, Significance value.

Based on the significance value we derive our confidence Interval (CI)



Confidence Interval : (CI)

It basically specifies a specific range based on significance interval, if the value falls within that range. Then, it is considered the coin is fair out of that range coin is not fair.

$$CI = 1 - \alpha = 1 - 0.05 = 0.95$$

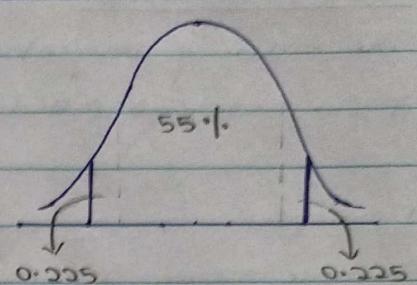
$$\therefore CI = 95\%.$$

- Consider, Health care domain where Covid Vaccine Test is conducted. Error here should be very less.
- Let significance value (α) = 0.15.

$$CI = 1 - 0.15 = 0.85$$

0.225
(Left) 0.225
(Right)

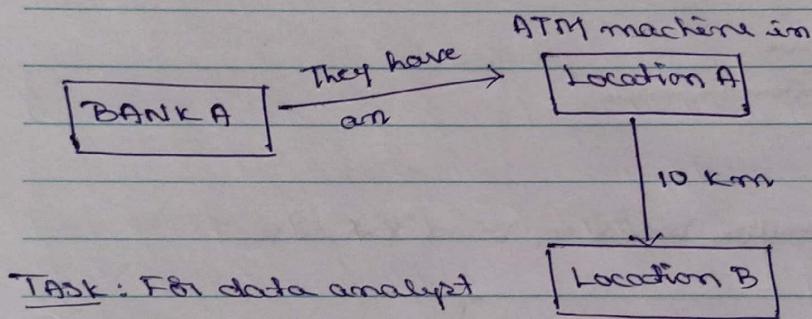
$$CI = 55\% \text{ (Restricted as } \alpha\text{-value is high)}$$



- Can we derive p-value from significance value?

Yes, possible

- Real world project: (Question Interview).



TASK: FBI data analyst

Based on the data that they have wrt Location A.

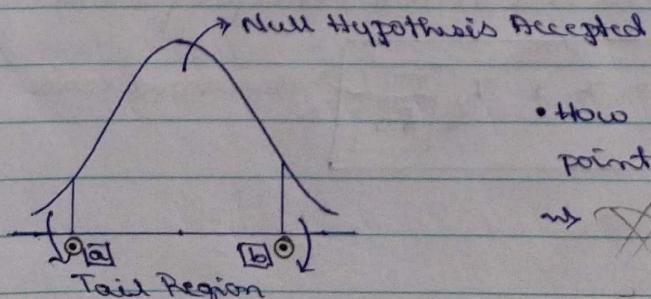
They have to find out whether an ATM should be opened in Location B or not.

This is done with the help of hypothesis testing

Confidence Interval:

There will be some range, if our values are within this range Null hypothesis will be accepted.

If value comes in tail region null hypothesis will be rejected.



- How do we find the values at point a & b for marking range?

\Rightarrow Z

Null Hypothesis Rejected.

Point Estimate :

The value of any statistics that estimates the value of a parameter is called point estimate.

Example:

Sample mean \bar{x}
Population mean μ] A perfect example of point estimate.

$\boxed{\bar{x}}$ is a point estimate of μ

which means, if $\mu = 3$ then \bar{x} can be 2.8 / 2.9 / 3.0 / 3.1 / 3.2 / 3.3 / 3.4 / 3.5
↓

Keeping sample mean value
we can estimate population mean value.
(concrete)

Confidence Interval formula based on Point Estimate:

$$\boxed{CI = \text{Point Estimate} \pm \text{Margin of Error}} \rightarrow \text{This in total gives the estimate of population.}$$

\downarrow
(sample data)

- Q) On the quantitative test of CAT Exam, the Standard Deviation is known to be 100. A sample of 25 test takers has a mean of 520. Construct a 95% CI about the mean.

Given: \therefore

SD = 100 = σ
(population)

n = 25

$$\bar{x} = 520$$

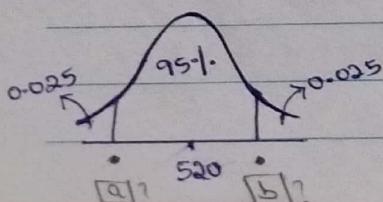
$$CI = 95\%$$

$$\alpha = 0.05$$

• Population SD is given $\frac{30}{\text{use}} \rightarrow Z\text{-score} \rightarrow Z\text{-table}$

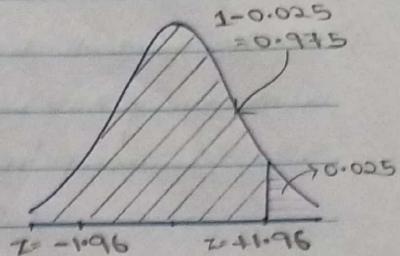
• CI = Point Estimate \pm Margin of Error

$$\therefore CI = \bar{x} \pm Z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \rightarrow \text{standard error}$$



M T W T F S S

- [a] • Lower Fence $C.I = \bar{x} - I_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$
- [b] • Higher Fence $C.I = \bar{x} + I_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$



• Lower Fence $C.I = 520 - (1.96) \frac{100}{\sqrt{25}}$

$$I_{\frac{\alpha}{2}} = I_{\frac{0.05}{2}} = I_{0.025} = 1.96$$

$$= 520 - (1.96) 20$$

$$= 480.8$$

Z-table

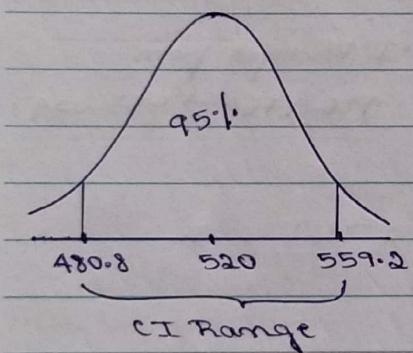
For value

0.9750

$$\downarrow \\ Z = 1.96$$

• Higher Fence $C.I = 520 + (1.96) \frac{100}{\sqrt{25}}$

$$= 559.2$$



→ For the above problem compute C.I Range values for $\sigma = 50$

• Lower Fence $= 520 - (1.96) \frac{50}{5} = 488.64$

• Higher Fence $= 520 + (1.96) \frac{50}{5} = 551.36$

- Q) On the recent test of CAT Exam, a sample of 25 test takers had a mean of 520 with a sample standard deviation of 20. Construct 95% CI about the mean.

Given

$$n = 25$$

$$\bar{x} = 520$$

$$CI = 95\%$$

$$\alpha = 0.05$$

$$S.D(\text{sample}) = 80 = s$$

NOTE:

When ever population Standard Deviation

is not given, we need to apply t-test (t-table)

$$C.I = \bar{x} \pm t_{\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}}$$

Standard Error.

For t-test we need to calculate

Degree of freedom = $n - 1$

$$25 - 1 = 24$$

$$C.I = 520 \pm 2.064 \left(\frac{80}{\sqrt{24}} \right)$$

C.I Ranges from

$$486.9776 \rightarrow 553.024$$

$$\frac{t_{\frac{\alpha}{2}}}{2} = \frac{t_{0.05}}{2} = t_{0.025} = 2.064$$

From t-table

Type 1 and Type 2 Error.

Consider the coin tossing as Example,

$$H_0 = \text{coin is fair}$$

$$H_1 = \text{coin is not fair}$$

Reality check:

⇒ The Null Hypothesis is True or Null Hypothesis is False

Decision: [Based on Experiment]

⇒ Null Hypothesis is True or False

H_0 - Criminal is not guilty

H_1 - Criminal is guilty

M T W T F S S

Outcome 1:

We reject the Null Hypothesis in reality if it is False - Yes

Example: Criminal is guilty, making them a prisoner-right decision

Outcome 2:

We reject the Null Hypothesis when in reality if it is True - No (Error)

Example: Criminal not guilty, but making them a prisoner-wrong decision

This kind of Error is Type 1 Error

Outcome 3:

We accept the Null Hypothesis, when in reality it is False - No (Error)

Example: Criminal is guilty, but not prisoned / set free - wrong decision

This kind of Error is Type 2 Error

Outcome 4:

We accept the Null Hypothesis, when in reality it is True - Yes

Example: Criminal is not guilty, so not prisoned / set free - right decision

Confusion matrix:

		H_0 (Not guilty)	H_1 (Guilty)
Not Prisoned (or set free)	Correct	Type II Error	
	Type I Error	Correct	

Example for Error:

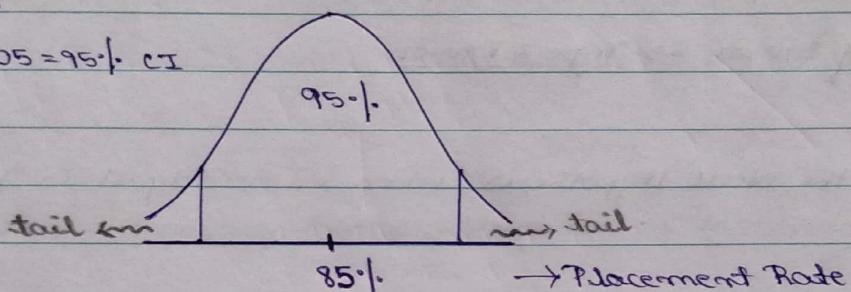
- 1) Person having cancer \rightarrow ML model predicts not having cancer.
- 2) Predicting stock market is going to crash \rightarrow but it did not crash next day

1 Tail and 2 Tail TestExample:

College in Karnataka has an 85% placement rate. A new college was recently opened and it was found that a sample of 150 students had a placement rate of 88%, with a standard deviation of 4%. Does this college have a different placement rate?

Let $\alpha = 0.05$

then $1 - 0.05 = 95\%$; CI



[a] Does this college have different placement rate?

(placement rate)

It can either increase or decrease - 2 Tail Test

(i.e., it can either be greater than 85% or lesser than 85%)

[b] Does this college have placement rate greater than 85%?

Placement rate focused only on one region - 1 Tail Test
(Right side)

[c] Does this college have placement rate lesser than 85%?

Placement rate focused only on one region - 1 Tail Test
(Left side)

[g] Power Law (Distribution) not fit

20%
• Pareto Distribution

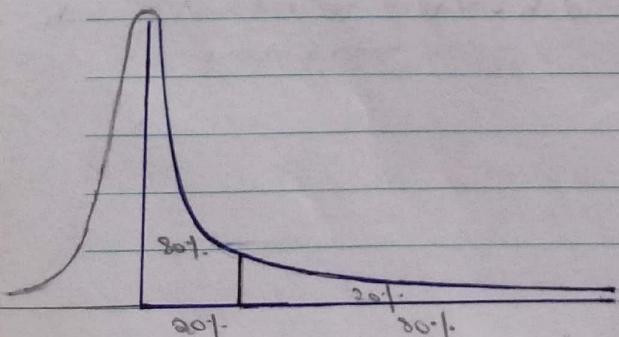
[Non-Gaussian Distribution]

• 80-20% rule

80% of data acquired with 20% of area

& remaining 20% of data acquired with 80% of area

• This graph is extension of Log Normal Distribution.



Examples:

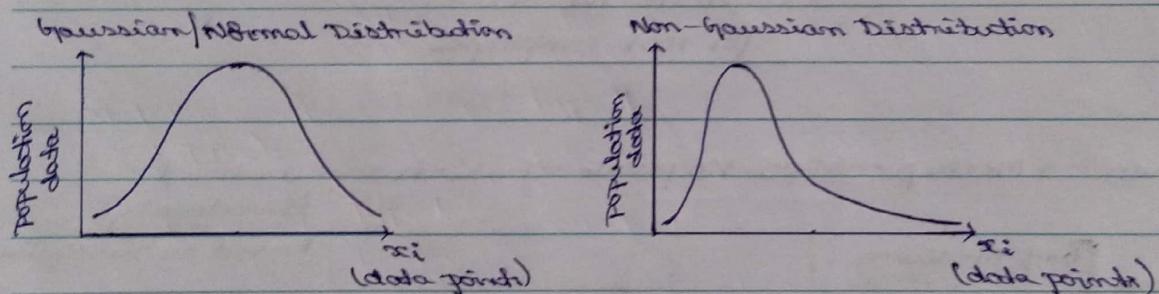
- 1) 20% of team is responsible for winning 80% of match.
- 2) 20% of score was scored by 20% of the team.
- 3) 20% of entire nation hold 80% of the oil.
- 4) 20% of wealth is held by 20% of the people.
- 5) 20% of Amazon revenue comes from 20% of the products.

Day 3

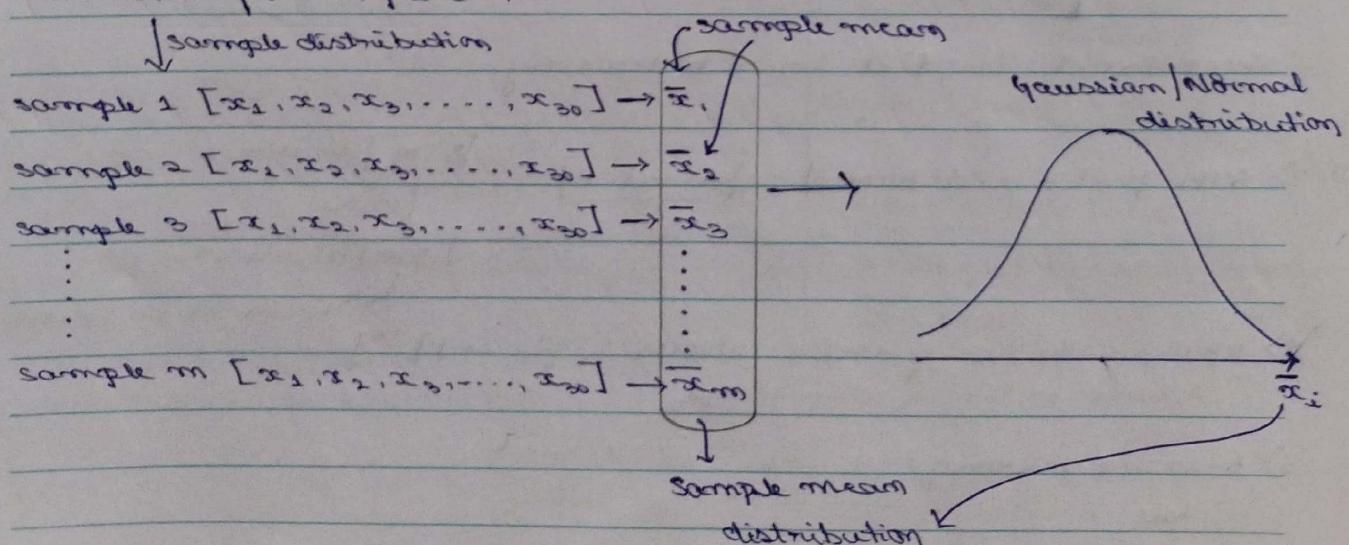
Central Limit Theorem

Let's say we have a distribution, it might be a Gaussian / Normal Distribution or Non-Gaussian Distribution.

If we keep on taking multiple samples from distribution where sample size is greater than or equal to 30, i.e., $n \geq 30$, then sample mean distribution that is obtained is plotted with the help of histogram we get Gaussian Distribution.



- From distribution pick up n elements for sample ($n \geq 30$)



M T W T F S S

• Is it compulsory to take sample size ($n \geq 30$)?

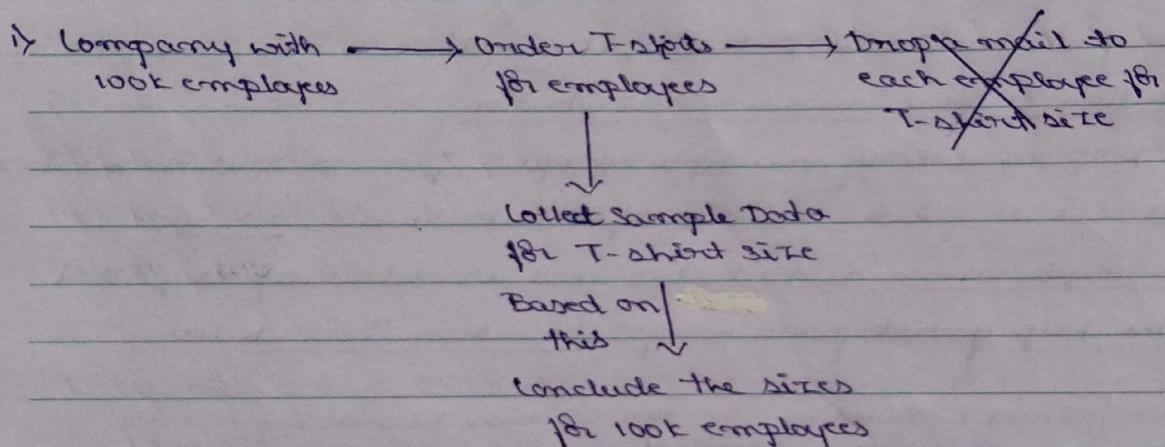
[a] Yes, it's compulsory for Non-Gaussian Distribution. If $n \geq 30$ is not followed we can not convert it into Gaussian Distribution

[b] If we consider $n \geq 30$ for Gaussian Distribution, we will get much better distribution from the sample.

Inferential Statistics (Important for Data Analyst & Data Scientist)

- Z-test (Z-table)
- T-test (T-table)
- Z-test proportion population.

Example:

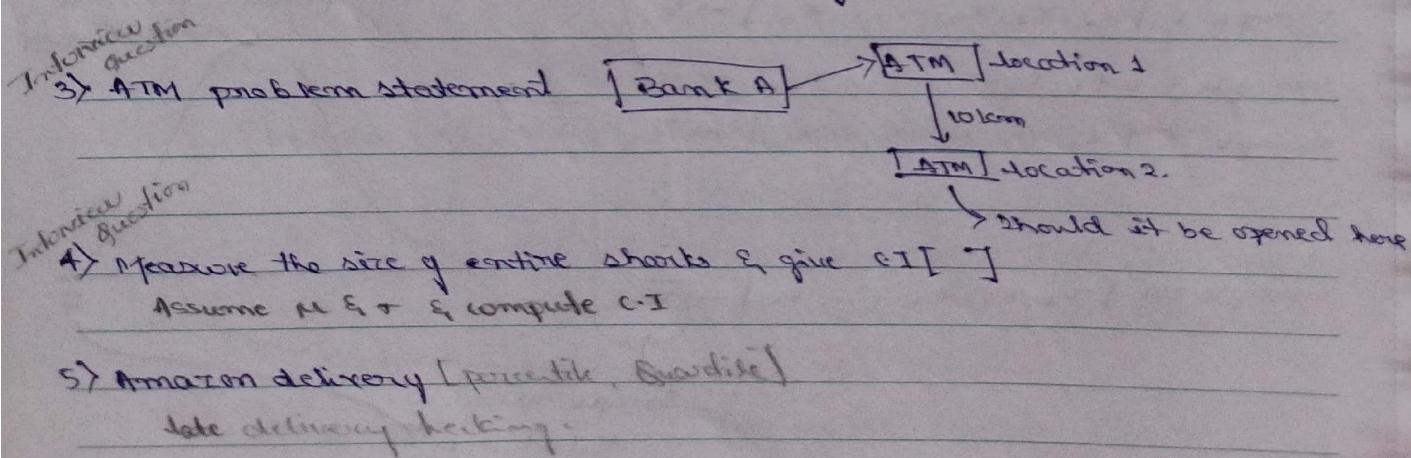


2) iNeuron → Meetup → 300-400 people → T-shirts → ordered (Randomly)
Based on Sampling

500 T-shirts

60%	Medium
10%	XL
10%	XXL
20%	L

Next Event: Plan it in a better way



Hypothesis Testing

a) Z-test

Q) A factory has a machine that fills 80ml of baby medicine in a bottle. An employee believes the average amount of baby medicine is not 80 ml. Using 40 samples, he measures the average amount dispensed by the machine to be 78ml with a standard deviation of 2.5.

[a] State Null and Alternate Hypothesis

[b] At a 95% CI, is there enough evidence to support the machine is not working properly? 2-Tail Test (can be < 80 or > 80)
(2 rejection regions)

Step 1:

Null Hypothesis (H_0) $\Rightarrow \mu = 80$

Alternate Hypothesis (H_1) $\Rightarrow \mu \neq 80$

$n = 40$

$\bar{x} = 78$ } sample

$s = 2.5$

Step 4: Calculate Test Statistic

$$\text{For sample standard deviation} \quad \text{WKT, } \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

σ / \sqrt{n} \rightarrow standard error

\downarrow SD for z-test can be ± 1.96 or depending on the question

$$Z = \frac{78 - 80}{2.5 / \sqrt{40}} = -2 = -5.05$$

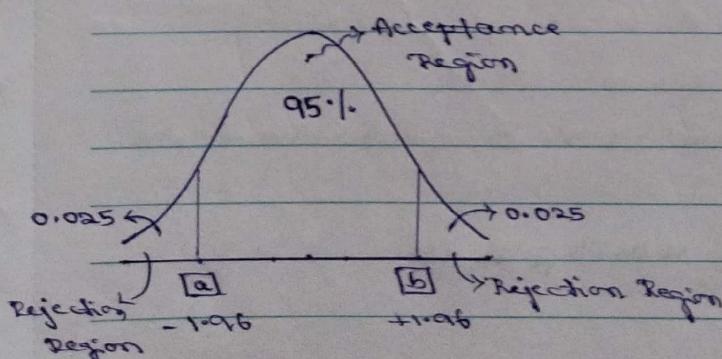
$Z = -5.05$

Step 2:

C.I. = 95%.

$\alpha = 0.05$

Step 3: Create Decision Boundary



NOTE:

- Why Z-test?

if $n > 30$ or $n < 30$

standard deviation

2) Population deviation &
sample Standard deviation

- Why t-test?

1) sample Standard deviation

2) $n < 30$

$Z\text{-score} \Rightarrow (1 - 0.025) = 0.9750$

\downarrow look for this

value in Z-table $\Rightarrow 1.96$

Step 5: State the results

Decision Rule: If $|z| = |z_0|$ is less than -1.96 or greater than 1.96 ,
then reject the null hypothesis with 95% C.I.

Result: Reject Null Hypothesis (H_0)

{ There is some fault in the machine with 95% C.I}

- Q) In the population the average IQ is 100 with a standard deviation of 15. A team of scientists wants to test a new medication to see if it has a +ve or -ve effect, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect Intelligence? 95%. C.I (2 Tail Test)

Step 1:

$$H_0: \mu = 100$$

$$H_1: \mu \neq 100$$

$$n = 30$$

$$\sigma = 15$$

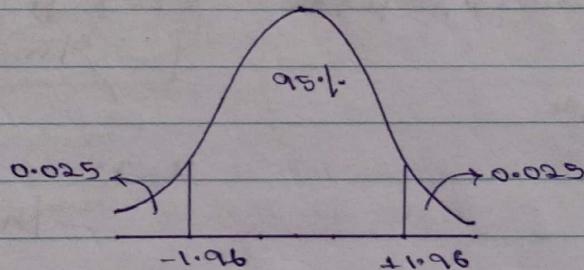
$$\bar{x} = 140$$

Step 3:

• Use Z-Score table

$$(1 - 0.025) = 0.9750$$

$$\downarrow 1.96$$

Step 4: Z-testStep 2:

$$C.I = 95\%$$

$$\alpha = 0.05$$

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{140 - 100}{15 / \sqrt{30}} = \frac{40 \times \sqrt{30}}{15} = 14.60$$

$$Z = 14.60$$

Step 5:

Decision Rule: If $|z|$ is less than -1.96 or greater than 1.96
Reject the H_0

Result: $14.60 > 1.96$ Reject the Null Hypothesis

+ve impact
increasing IQ

Q) A complaint was registered, the boys in the Municipal Primary School are underfed. Average weight of boys of age 10 is 32 kg with standard deviation is 9 kgs. A sample of 25 boys was selected from the municipal school and the average weight was found to be 29.5 kgs? with $(\alpha = 95\%)$, check whether it is True or False? (1 Tail Test)

Step 1:

$$H_0: \mu = 32$$

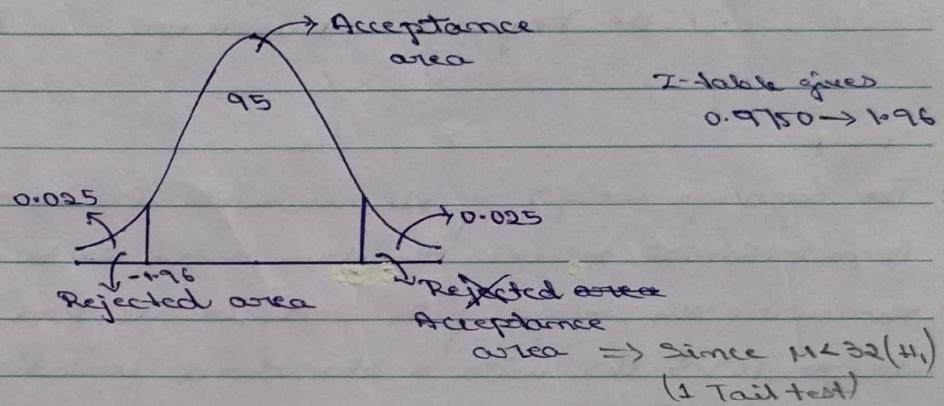
$$H_1: \mu < 32$$

$$\sigma = 9 \text{ kg}$$

$$n = 25$$

$$\bar{x} = 29.5$$

Step 3:



Step 2:

$$\alpha = 95\%$$

$$\alpha = 0.05$$

Step 4:

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{29.5 - 32}{9 / \sqrt{25}} = \frac{-2.5 \times 5}{9} = -1.38$$

Step 5:

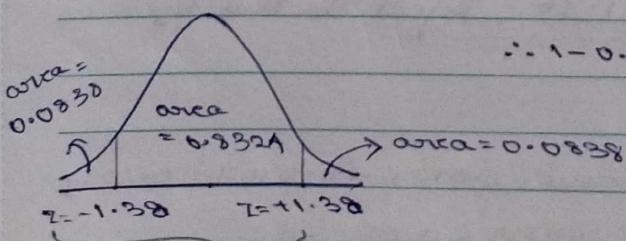
Decision Rule: If Z is less than -1.96 or greater than 1.96 , Reject Null Hypothesis.

Result: $-1.38 > -1.96$, Therefore we accept the Null Hypothesis.

\therefore The boys are not underfed.

\Rightarrow We all know significance value α , but what about p-value?

$$\boxed{\alpha = 0.05} \Rightarrow p\text{-value} = ?$$



$$\therefore 1 - 0.9162 = 0.0838$$

$$p\text{-value} = 0.0838 + 0.0838$$

$$\boxed{p\text{-value} = 0.1676}$$

For $Z = 1.38$, from Z-table

gives value 0.9162 (It is the area to the left of +1.38)

the process of finding

Is p-value same for both 1 tail & 2 tail test?

M T W T F S S

For all the problems, we found out H_0 is accepted or not using Z value or value obtained from ^{new} t-test value

We can also find H_0 is accepted or not

using p-value & significance value.

p-value calculated by
This is

wrt
2 tail

$$p = 0.1676 \quad] \text{ Here } 0.1676 > 0.05$$

$$\alpha = 0.05 \quad] \text{ i.e., if } p > \alpha \text{ we accept the } H_0$$

\therefore The Boys are not underfed.

NOTE:

p-value is calculated based on new Z-value

Q) The average weight of all residents in town XYZ is 168 lbs. A nutritionist believes the true mean to be different. She measured the weight of 36 individuals and found the mean to be 169.5 lbs with a standard deviation of 3.9.

a) At 95% C.I. is there enough evidence to discard the Null Hypothesis? (2 Tail Test)

Step 1:

$$H_0: \mu = 168$$

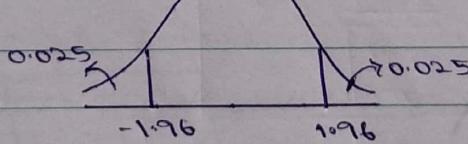
$$H_1: \mu \neq 168$$

$$n = 36$$

$$\bar{x} = 169.5 \quad \left. \begin{array}{l} \text{sample} \\ \text{mean} \end{array} \right.$$

$$s = 3.9$$

Step 3:



Step 4:

$$Z = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{169.5 - 168}{3.9/\sqrt{36}} = \frac{1.5}{3.9} = 0.39$$

Step 2:

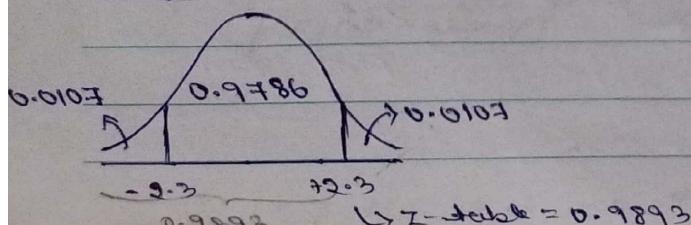
$$\alpha = 0.05$$

$$C.I. = 95\%$$

Step 5:

As $0.3 > 1.96$, Reject the Null Hypothesis

Step 6: wrt p-value.



$$p\text{-value} = 0.0107 + 0.0107 = 0.0214$$

$$0.0214 < 0.05$$

\therefore Reject the Null Hypothesis

b) t-test

M T W T F S S

Q) A company manufactures bike batteries with an average life span of 2.81 miles years. An engineer believes this value to be less. Using 10 samples, he measures the average life span to be 1.8 years, with a standard deviation of 0.15.

a) State the Null & Alternative hypothesis

b) At a 99% C.I., is there enough evidence to discard the H_0 ? (1 Tail Test)

Step 1:

$$H_0: \mu \geq 2$$

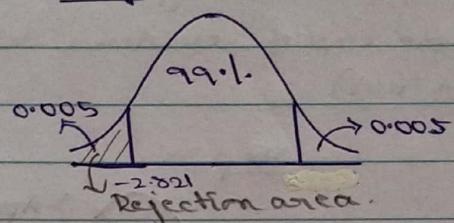
$$H_1: \mu < 2$$

$$n = 10$$

$$\bar{x} = 1.8$$
 } sample

$$s = 0.15$$

Step 3:



$$1 - 0.0005 = 0.9995$$

[T-test as
 $n < 30$ & 3 sig figs]

T-test

• Degree of freedom $\Rightarrow n - 1 = 10 - 1 = 9$.

$$\alpha = 0.01$$

Step 2:

$$C.I. = 99\%$$

$$\alpha = 0.01$$

$$\begin{array}{c} \diagdown \\ 0.005 \quad 0.005 \end{array}$$

Step 4: Calculate t-test statistics

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{1.8 - 2}{0.15/\sqrt{10}} = \frac{-0.2 \times \sqrt{10}}{0.15} = -4.216$$

Step 5:

$-4.216 < -2.821$, Reject the Null Hypothesis.

P-value can be calculated for this too.

This value is high & not within 3rd standard deviation.

so will not be able to see in the table.

C) Z-test with proportions (Proportions mean Ratio)
 (No standard deviation required here)

Q) A tech company believes that the percentage of residents in town XYZ, that owns a cell phone is 70%. A marketing manager believes that this value to be different. He conducts a survey of 200 individuals and found that 130 responded yes to owning a cell phone.

Q) State the Null and Alternate Hypotheses.

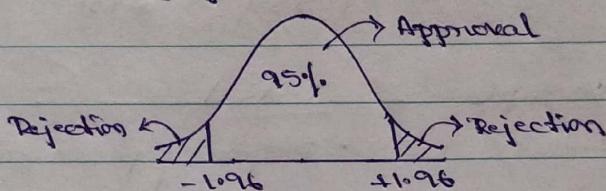
Q) At a 95% CI, is there enough evidence to reject the H_0 ?
 (2 tail test)

Step 1: \rightarrow Percentage

$$H_0: p = 0.70$$

$$H_1: p \neq 0.70$$

Step 3:



$$n = 200$$

$$x = 130 \quad \left[\begin{array}{l} \text{calculate} \\ \text{proportion } (\hat{p}) \\ \text{WRT this.} \end{array} \right]$$

$$\hat{p} = \frac{x}{n} = \frac{130}{200} = 0.65 \Rightarrow 65\%$$

Step 4: \rightarrow calculate Z -test, with proportion value

$$Z_{\text{test}} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0 q_0}{n}}}$$

$$Z = \frac{0.65 - 0.70}{\sqrt{\frac{0.70 \times 0.30}{200}}}$$

$$= \frac{-0.05}{0.0324} = -1.54$$

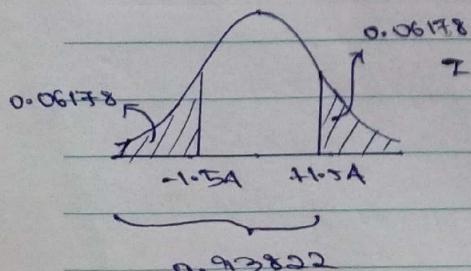
$$\therefore Z = -1.54$$

Step 5:

At 95% c.i there is $-1.54 > -1.96$,

so we Accept the Null Hypothesis.

Step 6: p-value



$$Z \text{ table} \Rightarrow 1.54 \rightarrow 0.93822$$

$$1 - 0.93822 = 0.06178$$

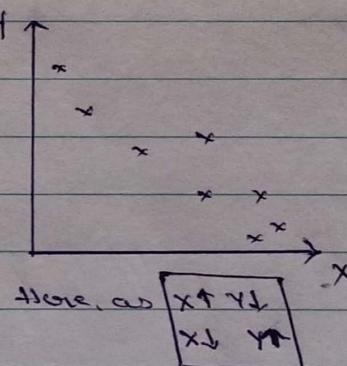
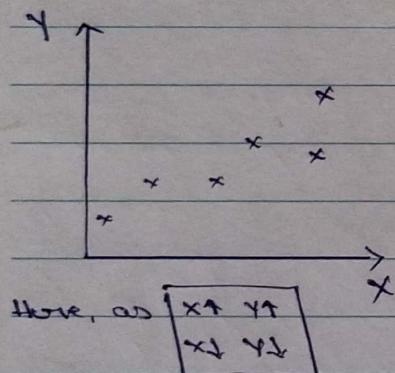
$$\Rightarrow P\text{-value} = 0.06178 + 0.06178 = 0.12356$$

$$\therefore P > \alpha \Rightarrow \text{Accept } H_0$$

Day 4

Covariance (Numerical variables)Consider 2 features $X \& Y$

- $X \quad Y$ • We need to quantify the relationship b/w $X \& Y$
- --- (Meaning what is the exact relationship b/w $X \& Y$)
- --- • When $X \uparrow$ is $Y \uparrow$
- --- $X \uparrow$ is $Y \uparrow$
- --- $X \uparrow$ is $Y \downarrow$
- --- $X \downarrow$ is $Y \uparrow$
- --- & vice versa



$$\text{cov}_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

where \bar{x} = mean of x
 \bar{y} = mean of y

$$\text{var}_x = \frac{\sum (x_i - \bar{x})^2}{N-1}$$

Relation b/w $\text{cov}(x,y)$ and variance(s) is:

$$\therefore \text{cov}(x,x) = \frac{\sum (x_i - \bar{x})^2}{N-1} = \text{var}(x)$$

- gives
 • Covariance - relationship b/w 2 variables

- If $x \uparrow y \uparrow \rightarrow$ Positively correlated

Example:

x	y
2	3
4	5
6	7
$\bar{x} = 4$	$\bar{y} = 5$
$(\frac{12}{3})$	$(\frac{15}{3})$

$$\text{cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

$$= \frac{(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)}{(3-1)}$$

$$= \frac{(-2)(-2) + 0 + (2)(2)}{2} = \frac{8}{2} = 4$$

$$\text{cov}(x, y) = 4 \Rightarrow \text{positive value}$$

\therefore Positively correlated

Why $x \uparrow y \downarrow \Rightarrow$ positive value

Case 1: $\text{cov}(x, y) = +ve$

Positive correlation

Case 2: $\text{cov}(x, y) = -ve$

Negative correlation

AND

$x \uparrow y \downarrow \left\{ \begin{array}{l} \text{cov}(x, y) \\ \text{Negative value} \end{array} \right.$

$x \downarrow y \uparrow \left\{ \begin{array}{l} \text{cov}(x, y) \\ \downarrow \\ \text{Negative correlation} \end{array} \right.$

Disadvantages of covariance:

- $\text{cov}(x, y) \Rightarrow$ Value will be either

+ ve value - ve value $\left[\begin{array}{l} \text{or} \\ \text{or} \end{array} \right] \Rightarrow$ there is no limit to the value like,
 $(500 / -400 / -300 / 1000 \dots)$ so it may go to infinity

problem is to restrict wrt the covariance, with some specific value

lets say -1 to +1



This is done by another correlation

called as Pearson correlation coefficient

Pearson Correlation Coefficient (For Numerical Values)

Each symbol called rho

$$\rho_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y}$$

because of this the value gets restricted b/w [-1 to 1]

σ_x = standard deviation of x

- More the value towards 1 more +ve it is correlated.

σ_y = standard deviation of y

- Why do we use Pearson Correlation?

Let say we have a dataset with 1000 features
features.

x y z A B c o/p

which all feature should be selected to predict the output (o/p) feature.

x, y, z, A, B, c

↓
(Dependent) feature

↓
(Independent) features

→ we will not use all the features in ML models

↓
It is not possible due to

“curse of dimensionality”
(As features ↑ cause problems, detailed in ML).

→ We are obviously supposed to select the most important feature that will definitely help us to predict the output

→ To find the relationship b/w x & o/p, y & o/p, ...

↓

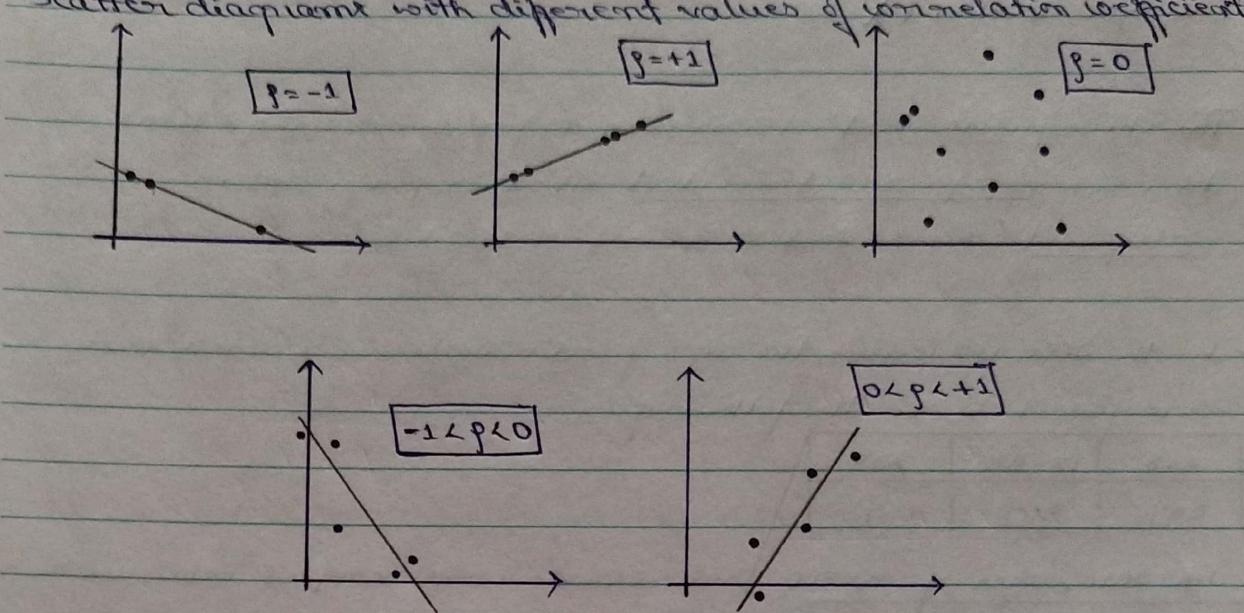
If they are highly correlated (either +ve or -ve)
keep those features

Now correlation b/w x & y is 99%.

then definitely delete one of the feature x, y as they are similar.

→ If $\rho=0$ (not at all correlated to output) remove that feature

Scatter diagrams with different values of correlation coefficient (ρ)



Spearman Rank Correlation (For Numerical Values)

$$r_s = \frac{\text{cov}(R(x), R(y))}{\sqrt{R(x) \cdot R(y)}}$$

(graph. disadvantages of Pearson correlation
Rank Correlation)

$R(x)$ = Rank of x

$R(y)$ = Rank of y

Consider,

x	y	$R(x)$	$R(y)$
1	2	4	4
3	4	3	3
7	5	2	2
0	7	5	1
8	1	1	5
7	2	2	4

Rank = 1, for highest values

Rank = same, for repeated values

BPT class 2 pp 15

d) Chi square Test (χ^2)

- The chi square test claims about population proportions.
- It is a non parametric test that is performed on categorical (numerical or ordinal) data.

Q) In the 2000 U.S census, the ages of individuals in a small town were found to be the following.

<18	18-35	>35
20%	30%	50%

In 2010, ages of $n=500$ individuals were sampled. Below are the results.

<18	18-35	>35
121	288	91

Using $\alpha = 0.05$, would you conclude the population distribution of ages has changed in the last 10 years?

$$n = 500$$

Expected	<18	18-35	>35	$\frac{20}{100} \times 500 = 100$
	20%	30%	50%	
Observed	100	150	250	$\frac{30}{100} \times 500 = 150$
	121	288	91	$\frac{50}{100} \times 500 = 250$

Step 2:

$$C.I = 95\%$$

Step 1:

H_0 : The data meet the expected distribution

$$\alpha = 0.05$$

H_1 : The data does not meet the expected distribution

Step 3: calculate the degree of freedom = (3 categories)

$$df = n - 1 = 3 - 1 = 2$$

In Chi-square table look for value using df and α .

$$\text{Now, } df = 2, \alpha = 0.05, \text{ value} = 5.991$$

Step 4: Decision

If χ^2 is greater than value (5.991) then, Reject H_0 .

Otherwise accept H_0 .

Step 5:

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

$$\chi^2 = \frac{(121-100)^2}{100} + \frac{(288-150)^2}{150} + \frac{(91-250)^2}{250}$$

$$\chi^2 = \frac{141}{100} + \frac{19044}{150} + \frac{25281}{250}$$

where f_o = Observed value f_e = Expected value

$$\chi^2 = 1.41 + 126.96 + 101.12$$

$$\chi^2 = 232.494$$

Step 6:As $232.494 > 5.991$, Reject the Null Hypothesis.

- Q) A school principal would like to know which days of the week students are most likely to be absent. The principal expect the students will be absent equally during the 5-day school week. The principal selects a random sample of 100 teachers asking them which day of the week they had the highest number of student absences. The observed and expected results are shown in the table below. Based on the results, do the days for the highest number of absences occur with equal frequencies ($C.I = 95\%$)

	Monday	Tuesday	Wednesday	Thursday	Friday
Observed	23	16	14	19	28
Expected	20	20	20	20	20

 $n=100$ (sample)Step 1: H_0 : The data meets the expected distribution H_1 : The data does not meet the expected distributionStep 2:

$$C.I = 95\%$$

$$\alpha = 0.05$$

Step 3: n (Categories)

$$df = n-1 = 5-1 = 4$$

Chi-square-table $df=4, \alpha=0.05$, value = 9.488Step 4:If χ^2 greater than 9.488, Reject H_0 ; if χ^2 less than 9.488 accept H_0

M T W T F S S

Step 5:

$$\begin{aligned} \chi^2 &= \sum \frac{(f_0 - f_e)^2}{f_e} = \frac{(23-20)^2}{20} + \frac{(16-20)^2}{20} + \frac{(14-20)^2}{20} + \frac{(19-20)^2}{20} + \frac{(28-20)^2}{20} \\ &= \frac{9}{20} + \frac{16}{20} + \frac{36}{20} + \frac{1}{20} + \frac{64}{20} \\ &= \frac{126}{20} = 6.3 \end{aligned}$$

$$\therefore \chi^2 = 6.3$$

Step 6:

$\chi^2 < 9.488$, Accept the Null Hypothesis.

Try -

e) ANOVA (F-Test) [Empirical Statistics]

ANOVA basically means Analysis of Variance.

- ANOVA is a statistical method used to compare the means of 2 or more groups.

- Two important things in ANOVA

- a) Factor (also called variables)

- b) Levels (are subcategories)

Example: 1) Consider Medicine for reducing Anxiety with different dosages.

0 mg	50 mg	100 mg
9	7	1
8	6	3
7	6	2
8	7	3
6	8	6

Here,

Factor: Dosage

Levels: 0 mg, 50 mg, 100 mg

Example: 2) Factor: gender

Levels: Male, Female

--	--	--	--	--	--	--

- Types of ANOVA:

⇒ One way ANOVA:

One factor with at least 2 levels, levels are independent.
Example: 1) Medicine ...

2) Repeated Measures ANOVA:

One factor with at least 2 levels, but levels are dependent.

Example: Running & covering certain distance each day in km by people (consider 10 ppl). (⑥ Study hours of certain ¹⁰ ppl on each day.)

	Day 1	Day 2	Day 3	Day 4	
1	6	4	2	8	
2	5	7	2	4	
⋮	⋮	⋮	⋮	⋮	
10	7	5	3	2	

⑦ No. of hours spent each day in gym by 10 ppl

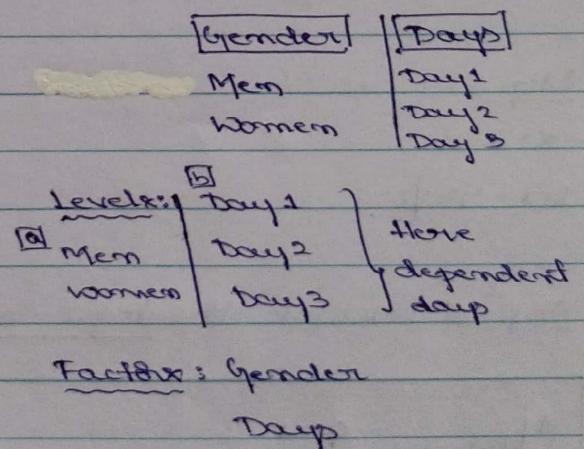
Here Day 1, Day 2, Day 3, Day 4 are dependent
(As a person will not have same energy every day)

3) Factorial ANOVA:

Two or more factors, with at least 2 levels with each factor.
levels can be either independent or dependent or both (mixed).

Example:

	Day 1	Day 2	Day 3
Men	9	7	1
	8	6	3
Women	7	5	3
	8	7	3
Women	8	9	4
	9	9	3



M T W T F S S

One Way Anova (F-test)

Q) Researchers want to test a new anxiety medication. They split participants into 3 conditions (0 mg, 50 mg, 100 mg), then ask them to rate their anxiety level on the scale of 1-10. Are there any differences b/w the 3 conditions using $\alpha = 0.05$?

0 mg	50 mg	100 mg
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

Step 1:

$$H_0: \mu_{0 \text{ mg}} = \mu_{50 \text{ mg}} = \mu_{100 \text{ mg}}$$

H_1 : Not all means are equal

Step 2:

$$\alpha = 0.05$$

$$1 - \alpha = 95\%$$

Step 3: Calculate the degree of freedom.

$$df_{\text{Between}} = a - 1$$

where,

$$df_{\text{Within}} = N - a$$

- N = Population

dataset

Combining all the datapoints in a

$$df_{\text{Total}} = df_{\text{Between}} + df_{\text{Within}}$$

- n = Sample

No. of datapoints in sample
in each category/level

- a = no. of levels



Here, $N = 7 \times 3 = 21$; $n = 7$; $a = 3$

$$df_{\text{Between}} = a - 1 = 3 - 1 = 2$$

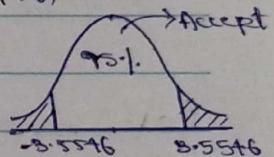
$$df_{\text{Within}} = N - a = 21 - 3 = 18$$

$$df_{\text{Total}} = N - 1 = 21 - 1 = 20 \quad (\text{this is nothing but } 2 + 18)$$

Step 4: State Decision Rule

$$\begin{aligned} df_{\text{Between}} &= df_1 \\ df_{\text{Within}} &= df_2 \\ \alpha \text{ value} & \end{aligned} \left. \begin{array}{l} \text{Required} \\ \text{for F-table} \end{array} \right\}$$

$$\begin{aligned} df_1 &= 2 \\ df_2 &= 18 \\ \alpha &= 0.05 \end{aligned} \quad \begin{array}{l} \text{Decision Rule range} \\ (2, 18) \end{array}$$



$$F\text{-table value} = 3.5546$$

If F-test is greater than F-table ~~test~~ value, Reject Null Hypothesis,
if lesser than F-table value, Accept Null Hypothesis.

Step 5: Calculate F-Test Statistic.

	SS	df	MS	F Test
Between	98.67	2	49.34	86.56
Within	10.29	18	0.57	
Total	108.96	20		

[a] Sum of squares: SS

$$\text{SS}_{\text{Between}} = \frac{\sum (\sum a_{ij})^2}{n} - T^2 \quad \text{where, } T = \left[\sum (\sum a_{ij}) \right]^2$$

$$\text{SS}_{\text{Within}} = \sum y^2 - \frac{\sum (\sum a_{ij})^2}{n} \quad \text{where, } y^2 = \text{Each } y_i \text{ every value squared}$$

[B] Degree of Freedom (df)

[C] Mean square (MS)

$$\boxed{MS = \frac{SS}{df}}$$

[d] F Test

$$\boxed{F = \frac{MS_{\text{Between}}}{MS_{\text{Within}}}}$$

Calculations:

[a] SS

$$\begin{aligned} SS_{\text{Between}} &= \frac{\sum (\sum a_{ij})^2 - \bar{T}^2}{n} \\ &= \frac{57^2 + 47^2 + 21^2}{7} - \frac{(57 + 47 + 21)^2}{21} \end{aligned}$$

$$\begin{aligned} \sum (\sum a_{ij})^2 &= (9+8+7+8+8+9+8)^2 \\ &\quad + (7+6+6+7+8+7+6)^2 \\ &\quad + (4+3+2+3+4+3+2)^2 \\ \sum (\sum a_{ij})^2 &= 57^2 + 47^2 + 21^2 \end{aligned}$$

$$SS_{\text{Between}} = 98.67$$

$$\begin{aligned} SS_{\text{Within}} &= \sum y^2 - \frac{\sum (\sum a_{ij})^2}{n} \\ &= 853 - 842.7 \end{aligned}$$

$$\begin{aligned} \sum y^2 &= 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + 9^2 + 8^2 + \\ &\quad 7^2 + 6^2 + 6^2 + 7^2 + 8^2 + 7^2 + 6^2 \\ &\quad 4^2 + 3^2 + 2^2 + 3^2 + 4^2 + 3^2 + 2^2 \end{aligned}$$

$$SS_{\text{Within}} = 10.29$$

$$\sum y^2 = 853$$

Step 6:

As $86.56 > 3.5516$, we Reject the Null Hypothesis.