

# Statistics

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

**Data** : "facts or pieces of information"

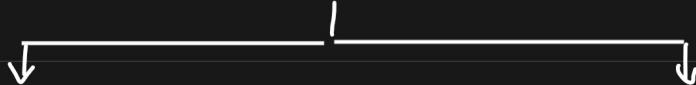
Eg: Height of students in a classroom

→ { 175cm, 150cm, 140cm, 130cm, 155cm }

Eg: Intelligence Quotient (IQ) of 5 randomly selected individuals (109, 89, 129, 101, 105) → Data.

Two Types

Statistics



① Descriptive Stats

② Inferential Stats

It consists of organizing and summarizing of data.

It consists of using that you've measured to form



Conclusions

Eg: Pdf, Histogram, Box plot, Bar chart, Pie charts

Eg: Hypothesis Testing, p value, Z test, t-test, Anova, Chi-square

Eg: 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 }

$\min = \text{mode}$

**Descriptive stats** : What is the average age of student in

your maths class?

Inferrential question : Are the ages of students in this maths classroom similar to what you would expect in a normal maths class at this university?

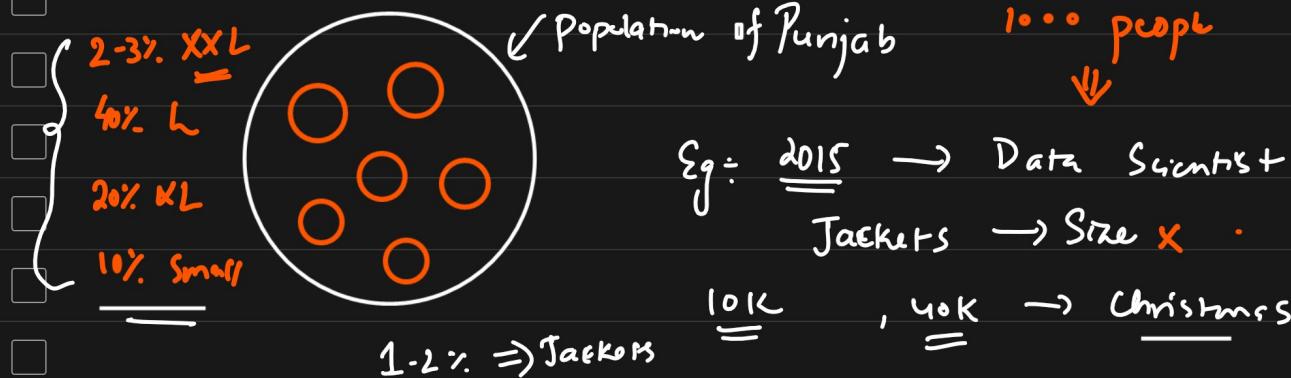


Population And Sample Data  $\rightarrow$  Inferrential Statistics Results

Elections : Punjab

{ AAP, Congress }

Exit Polls  $\leftarrow$

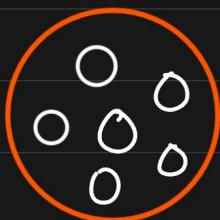


Population (N) ✓

Sample (n) ✓

Sampling Techniques

① Simple Random Sampling  $\downarrow$  : Every member of the population (N) has an equal chance of being selected for your Sample (n)



## ② Stratified Sampling

Strata → Layers  
 ↓  
 Clusters  
 ↑  
 Non overlapping groups

Gender → Male  
 Female

Age groups  
 0-18 }  
 18-35 }  
 35-60 }

Blood groups

Tax slabs  
 Courses

Education Qualification

Thanas

Illustration

③ Systematic Sampling

Snap

Customs

(N) → Select every  $n^{\text{th}}$  individual

$6^{\text{th}}$   
 =

↓  
Stratified

Eg: Survey → Mail (SBI credit card)

↓  
 (SBI credit card)

④ Convenience Sampling : Only those people who are interested will only be participating.

Healthcare Disease

Eg: Data Science → AI }  
 YouTube Survey → }

{ Blind people }

→ RBI → Household Survey → Female ←  $\frac{\downarrow \downarrow \downarrow \downarrow \downarrow}{\text{Economics}}$  → DATA Science }

Exit Poll : Stratified + Random Sampling

## Variable

A variable is a property that can take on any value

Eg: Height = 182  
150  
145  
160

{ 182, 170, 145, 160 }  
↓  
No

## Two kinds of Variable

① Quantitative Variable → Measured Numerically { Add, Subtract,  $\times$ ,  $\div$  }

② Qualitative Variable.

↳ Eg: Gender  
→ Male { Based on some { characteristics } we can  
→ Female derive categorical variables }

{ Quantitative → Qualitative Variable. }

Eg: IQ

0-10      10-50      50-100

↓      ↓      ↓

Low IQ      Medium IQ      Good IQ

Quantitative

↓  
Discrete Variable

↓  
Continuous Variable.

Eg: whole number

Eg: Height = 172.5, 162.5 cm,

163.5 cm.

Eg: No. of Bank accounts

{ 2, 3, 4, 5, 6, 7 }      2.5, x  
                                  2.75 x

Rainfall: 1.35, 1.25, 1.75, 2.25 cm

Weight

Eg: Total No. of children in a family

Temp

Eg: 2, 3, 4, 5

Stock price.

25, 2.75

Eg: Total no. of Employees in a Company Eg: 10k,

Ass:

- ① What kind of variable Marital Status is? Categorical
- ② What kind of variable Nile River length is? Continuous Quantitative
- ③ What kind of " Movie duration is? " "
- ④ What kind of variable IQ is? " "

Frequency Distribution

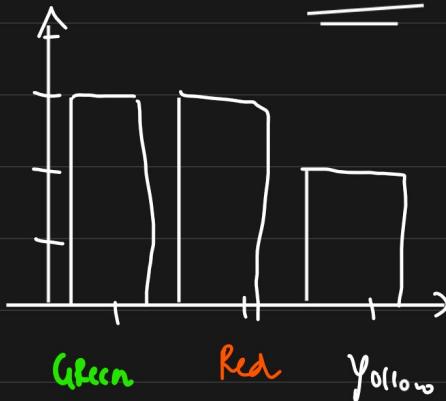
Sample Data: Green, Red, Yellow, Green, Red, Yellow, Green, Red

↓

Colors	Frequency
Green	3
Red	3
Yellow	2

① Bar Graph frequency

Bar Chart



## ① Variable Measurement Scales

4 types of Measured Variable.

① Nominal data {Categorical data}

Eg: Colors, Gender, Types of flowers

Ranking is not that important

② Ordinal data

Student (Marks)

→ 100

96

57

85

44

Rank

1

2

4

3

5

Percentiles

Ordinal Data.

$\frac{\text{Phd}}{\downarrow} \rightarrow \{ \text{NLP} \}$

Degree

Phd

B.G

Master

BCA

12

Salary

✓

✓

✓

✓

✓

Assignment

④ Ratio data ✓

③ Interval data ✓

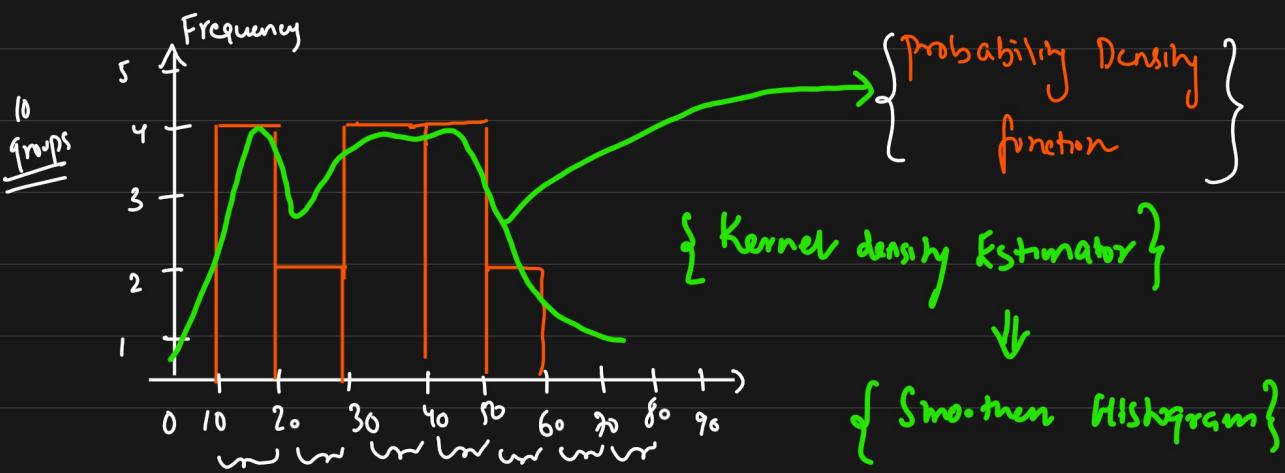
④ Histograms ÷ Continuous

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

Histogram

$\rightarrow \text{Bins} = 10$

Mean, Median, Mode



0 - 10  $\rightarrow$  0-5, 5-10, 10-15, 15-20, 20-25, 25-30, 30-35

Assignment

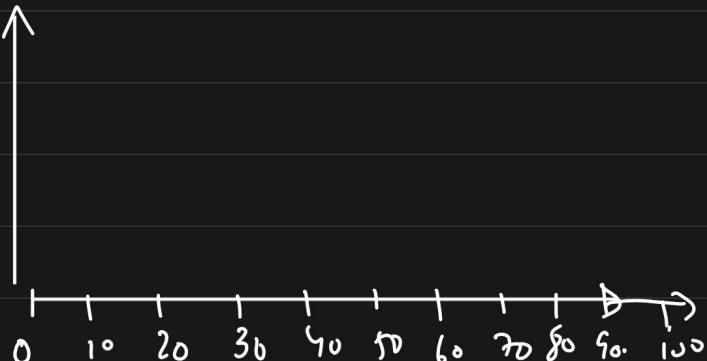
Ex: 10, 13, 18, 22, 27, 32, 38, 40, 45, 51, 56, 57, 88, 90, 92, 94, 99

bins = 10

0-10 10-20 20-30 30-40

40-50 50-60 60-70

70-80 80-90 90-100

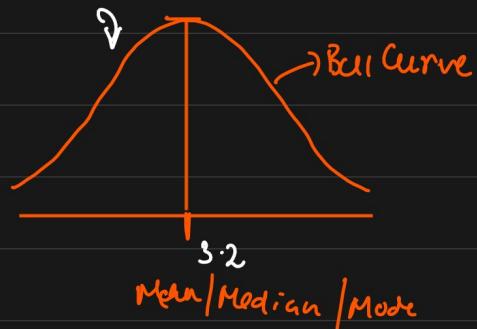


## Intermediate Stats

- 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 → Central position of the dataset
- 

- ① Mean ✓
- ② Median ✓ { EDA & Feature Eng. }
- ③ Mode ✓
- 



### Population (N)

### Sample (n)

$$X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\} \quad \overbrace{\quad}^{\text{Sample}} = \bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$

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

Sample  
Mean

Population  
mean

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

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

## Median

1, 2, 2, 3, 4, 5

1, 2, 2, 3, 4, 5, 100

$$\bar{x} = \frac{1+2+2+3+4+5}{6} = \frac{17}{6} = 2.83$$

$$\bar{x} = \frac{1+2+2+3+4+5+100}{7} = \frac{117}{7} = 16.71$$

Median ✓

1, 2, 2 3 4, 5, 100

$$\bar{x} = 16.71 //$$

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

1, 2, 2, 3, 4, 5 → odd or even  
↓  $\frac{2+3}{2} = 2.5$

$$\frac{2.5}{2} \approx 2.83 //$$

Mode ÷ Highest frequency: Median

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

↓  
3  
↓

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

[2, 3, 4]

EDA

Feature Engineering

↪ NAN values ⇒ Continuous Values + outlier  
= = Mean ↓  
= Median

⇒ Categorical Variable.

↓  
Mode

Agri

Lidley, Sunflower, Rock, - - - , Min, Max

Measure of Dispersion →  $\{ \text{Dispersion} \}$

① Variance

② Standard deviation

$\} \downarrow$

Spread ⇒ How the data is spread



≠



≠



① Variance

Population Variance

$\{$  Bessel's Correction

Degree of freedom  $\}$

Sample Variance

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

Population mean

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

Sample mean  
 $n-1$

Eg:

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

	<u><math>x</math></u>	<u><math>\bar{x}</math></u>	<u><math>x - \bar{x}</math></u>	<u><math>(x - \bar{x})^2</math></u>
1	2.83	2.83	-1.83	3.34
2	2.83	2.83	-0.83	0.6889
2	2.83	2.83	-0.83	0.6889
3	2.83	2.83	0.17	0.03
4	2.83	2.83	1.17	1.37
5	2.83	2.83	2.17	4.71

$$\left[ \frac{10.84}{5} \right] = 2.168$$

↑

$n=6$

$n-1$

$$\mu = 2.83$$

Let consider

$$10.84$$

$$\text{Variance} = \frac{\sigma^2}{6.42}$$

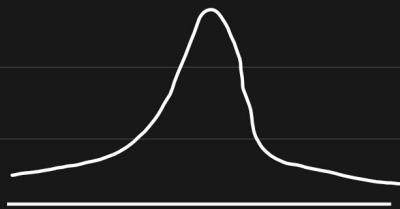
as an example

Spread ↑↑

$$\text{Variance} = 2.168$$

Variance ↑↑

Spread ↑↑



Standard deviation

$$\sigma = \sqrt{\text{Variance}} = \sqrt{2.168}$$

$$= \sqrt{1.472}$$

Variance

↓

Spreadness

1, 2, 2, 3, 4, 5



$$\begin{array}{r} 2.83 \\ 1.472 \\ \hline 4.302 \end{array}$$

$$\begin{array}{r} 1.472 \\ 5.774 \\ \hline 7.246 \end{array}$$

$$\begin{array}{r} 2.83 \\ -1.472 \\ \hline 1.358 \end{array}$$

$$\begin{array}{r} 1.358 \\ 1.472 \\ \hline 0.114 \end{array}$$

Outliers

Outliers

④ Percentiles And Quartiles



Percentiles : 1, 2, 3, 4, 5

% of the numbers that are odd?

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

Percentiles :  $\{CAT, GATE, SAT\} \Rightarrow \underline{\underline{99\%}}$

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

99 percentiles mean the person has got better marks than 99% of the students.

Data set : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

What is the percentile ranking of 10?  $\boxed{n=20}$

Percentile Rank of  $x = \frac{\text{# of values below } x}{n} \times 100$

$$= \frac{16}{20} \times \underline{\underline{0.8}} = 80 \text{ percentile}$$

$$= \frac{17}{20} = 85$$

② What value exists at percentile ranking of 25%?

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

$$= \frac{25}{100} \times (21) = \underline{\underline{5.25}} \rightarrow \text{Index}$$

Value = 5

Quartiles (25%)

Five Number Summary

- ① Minimum
- ② First Quartile (25%)  $Q_1$
- ③ Median
- ④ Third Quartile (75%)  $Q_3$
- ⑤ Maximum

Removing the Outliers

Inter Quartile Range: (75% - 25%)  
 $Q_3 - Q_1$

$\{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \cancel{10}\}\}$

Lower Fence  $\longleftrightarrow$  Higher Fence

Lower Fence =  $Q_1 - 1.5(IQR)$  (25%)  $Q_1 = \frac{3+5}{2} = 4$   $\times 1.5 = 6$   $\rightarrow 5^{\text{th}}$  index

Higher Fence =  $Q_3 + 1.5(IQR)$

$IQR = Q_3 - Q_1 = 7 - 3 = 4$  (75%)  $Q_3 = \frac{7+8}{2} = 7.5$   $\times 1.5 = 11.25$   $\rightarrow 15^{\text{th}}$  index

$Q_3 = 7$

Lower Fence =  $3 - 1.5(4) = 3 - 6 = -3$

Higher Fence =  $7 + 1.5(4) = 7 + 6 = 13$

$[-3 \longleftrightarrow 13]$

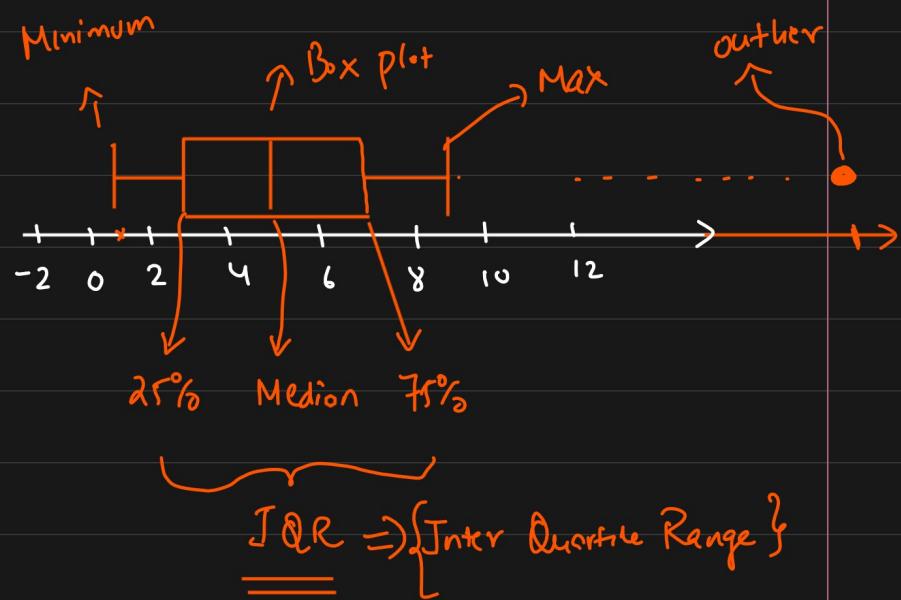
Remaining  $\frac{5+5}{2} = 5$

$1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \cancel{10}$

## 5 Number Summary

## Box plot $\rightarrow$ Outliers

- Minimum = 1
- $Q_1 = 3$
- Median = 5
- $Q_3 = 7$
- Max = 9



## Use of Box plot

## 1 Distributions

① Normal / Gaussian Distribution ✓

② Standard Normal Distribution ✓

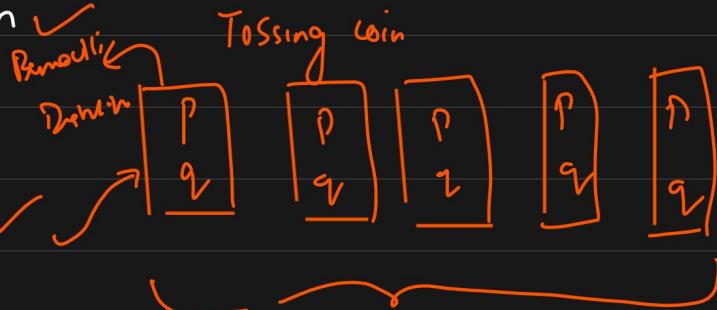
③ Z - Score ✓

④ Log Normal Distr ✓

⑤ Bernoulli's Distribution ✓

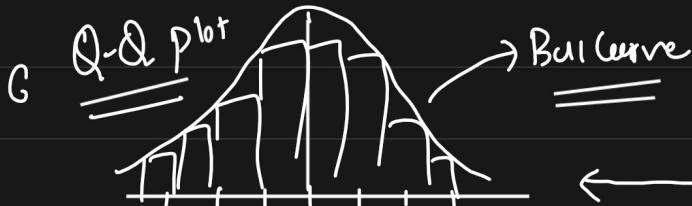
⑥ Binomial Distribution }

① Gaussian / Normal Distribution



Properties

{ Power Law }



① Empirical Rule of Gaussian Distribution  $\downarrow$   
80-20%

DATASET  $\rightarrow$  IRIS Dataset }  $\rightarrow$  Petal, Sepal length  $\downarrow$   
Domain (Exponne )

② Weight of human brain

③ Height  $\rightarrow$  Doctor

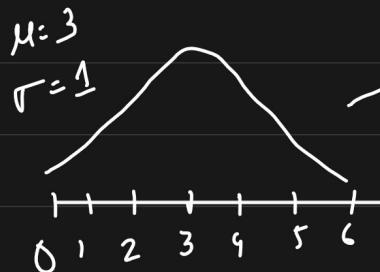
68.2, -95.4 - 99.7

Outliers

## Standard Normal Distribution

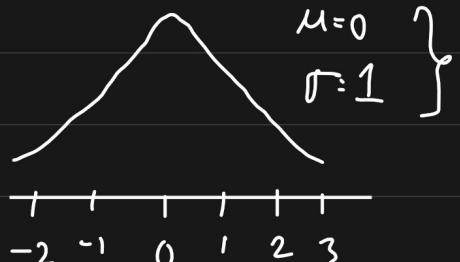
{1, 2, 3, 4, 5}

$\mu = 3$



$\sigma = 1.414 \approx 1$

$$\Rightarrow \begin{array}{c} \overline{\mu = 0} \\ \overline{\sigma = 1} \end{array}$$



{1, 2, 3, 4, 5}

$$\left\{ Z\text{-Score} = \frac{x - \mu}{\sigma} \right\}$$

=

$$\begin{array}{l} \frac{3-3}{1} = 0 \\ \frac{2-3}{1} = -1 \\ \frac{1-3}{1} = -2 \end{array}$$

Why 22

$$\begin{array}{c} \boxed{\mu = 0} \\ \boxed{\sigma = 1} \end{array}$$

✓

## Standardization vs Normalization

Years ↑  
Age ↑  
25  
26  
28  
30  
32  
un

different unit

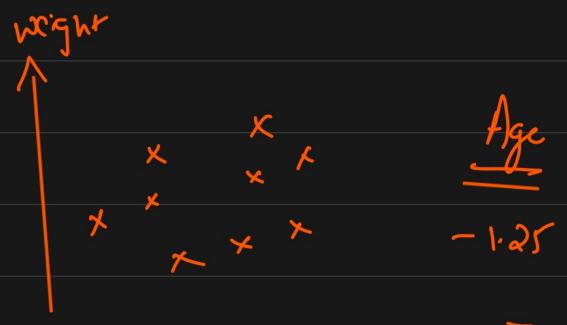
kg  
75  
80  
85  
60  
70  
un

INR

Salary  
25K  
30K  
40K  
50K  
un

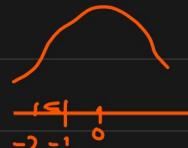
$$\frac{25 - 28.2}{25.6}$$

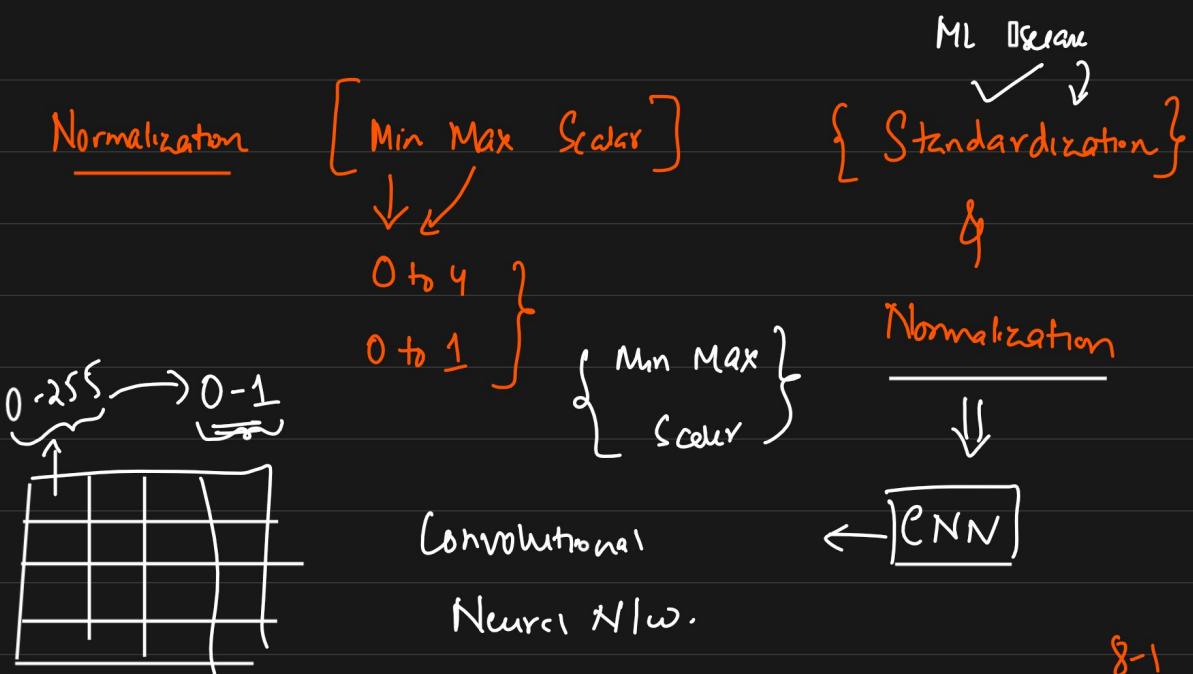
Same unit scale ??



Maths → Scale

Standardization





f1

2

5

6

8 ✓

1 ✓

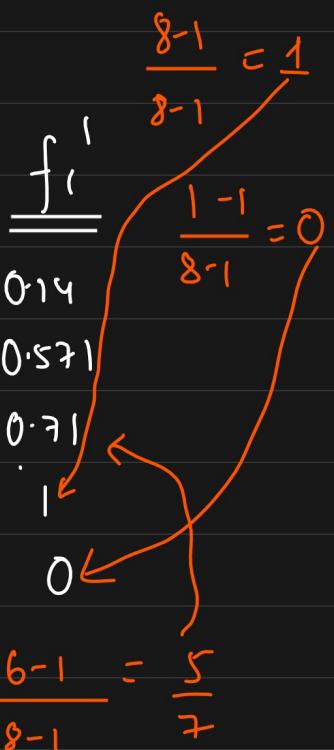
Normalization

$$\left\{ \begin{array}{l} x_{\text{Norm}} = \frac{x_i - x_{\min}}{x_{\max} - x_{\min}} \\ \downarrow \end{array} \right.$$

$$\begin{array}{c} \text{Min Max Scalar} \\ \hline \end{array} = \begin{array}{c} 0 \text{ to } 1 \\ \hline \end{array}$$

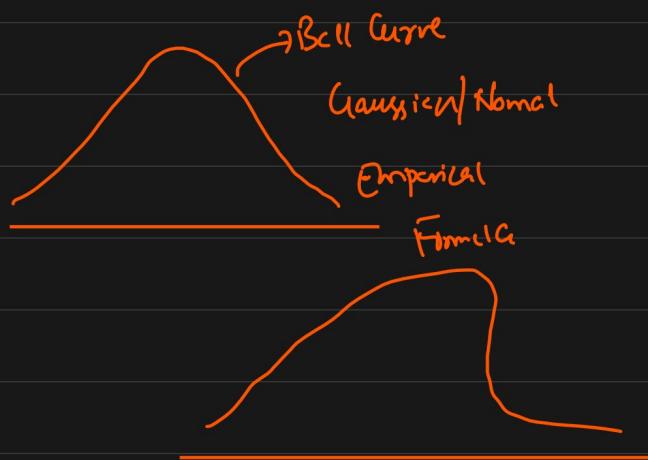
$$= \frac{2 - 1}{8 - 1} = \frac{1}{7} = 0.142$$

(0 - 1)

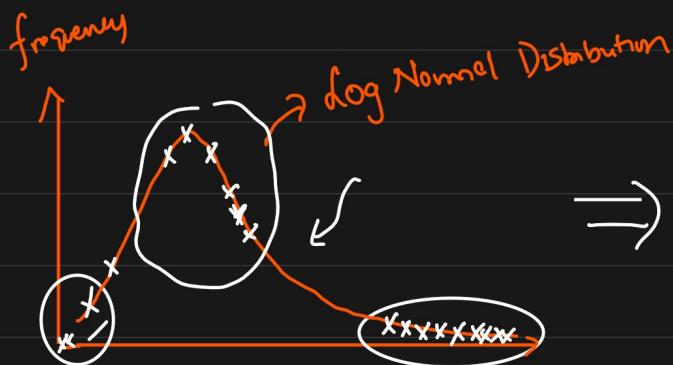
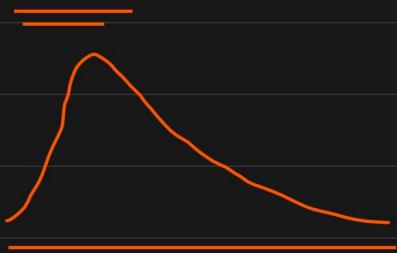


$$\frac{5-1}{8-1} = \frac{4}{7} = 0.571$$

# Log Normal Distribution

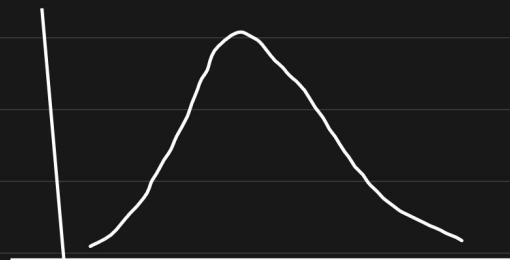


## Skewed Curve



## Gaussian Distribution

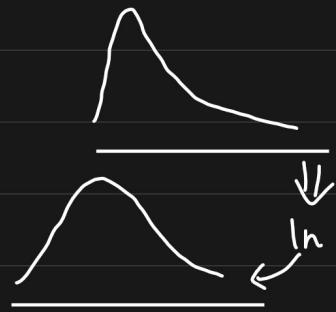
### Normal Distribution



$$X = \text{log Normal Distributed}$$

$$\{ Y = \ln(X) \}$$

Gaussian Distribution



$$\{ X = \exp(Y) \}$$

$X$

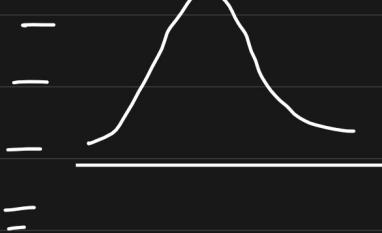
$Y = \ln(x)$

25

30

40

45



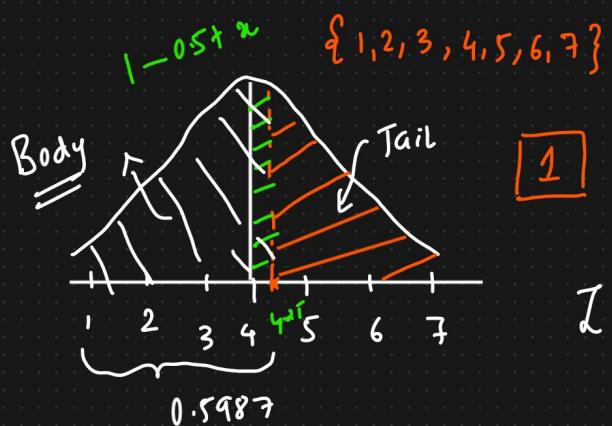
① Bernoulli's Distribution

## Day 2 - Stats

$$\textcircled{1} \quad Z\text{-Score} = \frac{x_i - \mu}{\sigma}$$

Stats Interview Question

How many standard deviation



$$\mu = 4$$

$$\sigma = 1$$

4.25 fall from the mean??

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$$

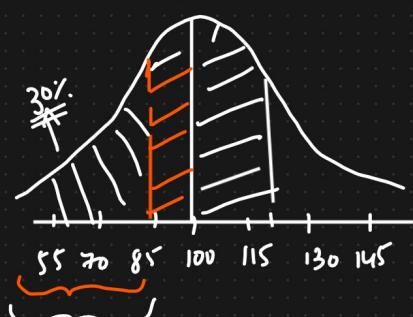
Question : What percentage of scores fall above 4.25?

$$1 - 0.59871 = 0.4013 \Rightarrow 40.13\%$$

2 In India the average IQ is 100, with a standard deviation of 15.

What is the percentage of the population would you expect to have an IQ lower than 85?

Ans)



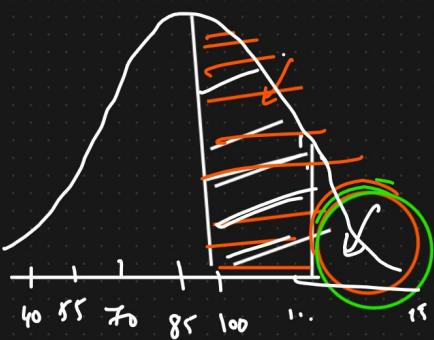
$$Z\text{-Score} = \frac{85 - 100}{15} = \frac{-15}{15} = \boxed{-1}$$

① Area under this curve

$$0.5 - 0.15866 = 0.34143 \Rightarrow \boxed{34.14\%}$$



$$\{ \text{Growth} = 100 \text{ less than } 125 \}$$

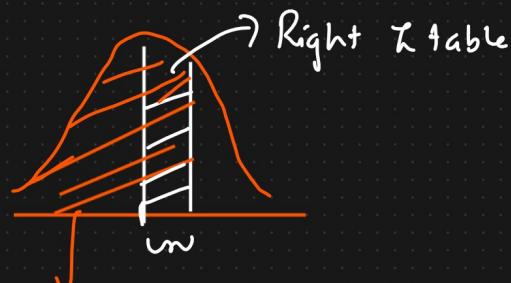


$$Z\text{score} = \frac{125 - 100}{15} = \frac{25}{15} = \frac{5}{3} = 1.667$$

$$\text{Ans} = 0.4515 \Rightarrow 45.15\%$$

1.667

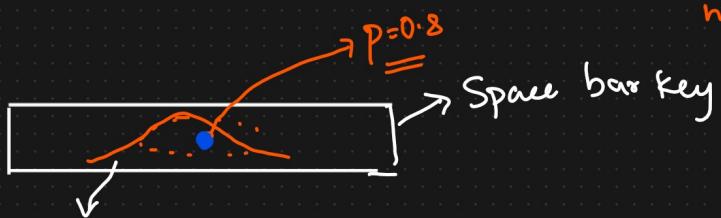
$$\underline{0.5 - 0.4515 = 0.0485} \Rightarrow 4.8\%$$



Left Z table

P value, Hypothesis Testing, Confidence Interval

Out of all 100 touches, the no. of touches is 80



$$P = 0.4$$

Out of all 100 touches, the no. of touches 40 times.

Hypothesis Testing, C.I., Significance value Together Fair Coin

Coin  $\rightarrow$  Test whether the coin is a fair coin or not by performing 100 tosses

$$\begin{array}{c} \text{P(H)} = 0.5 \\ = \end{array} \quad \begin{array}{c} \text{P(T)} = 0.5 \\ = \end{array}$$

## Hypothesis Testing

Criminal is  $\rightarrow$  Court

SMOLAY

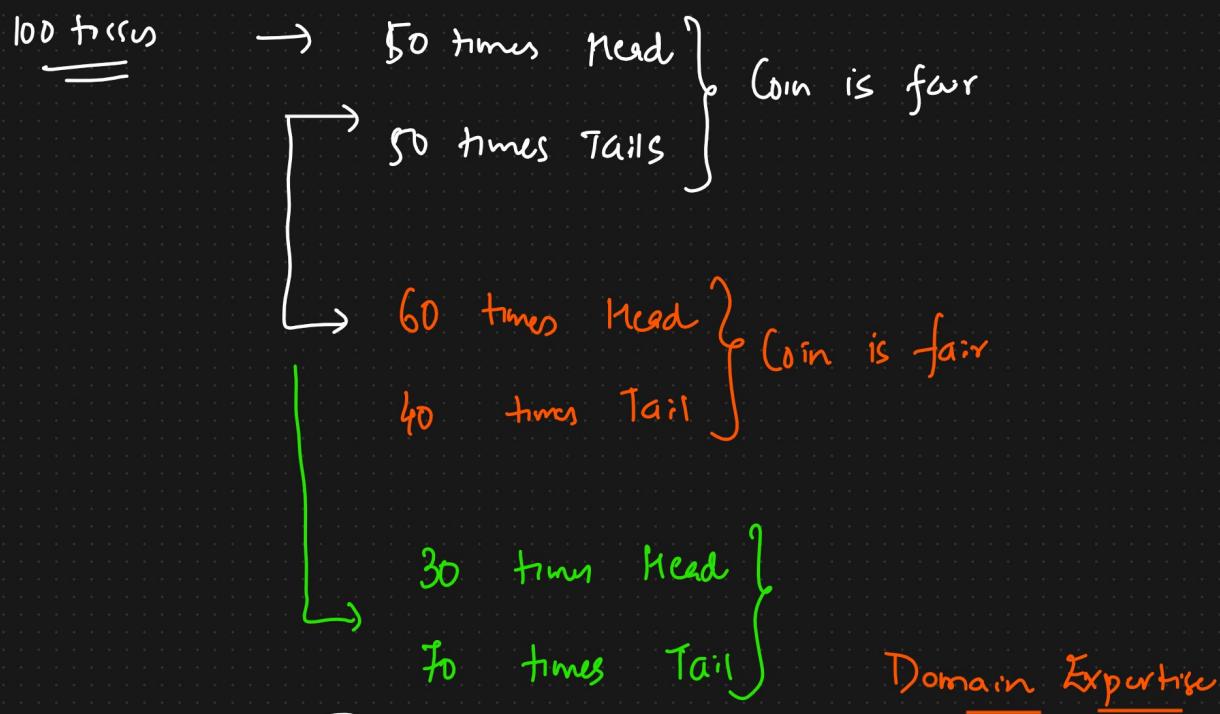
$P(H) = 100\%$   $P(T) = 0\%$

① Null Hypothesis — Coin is fair  $\rightarrow (H_0)$

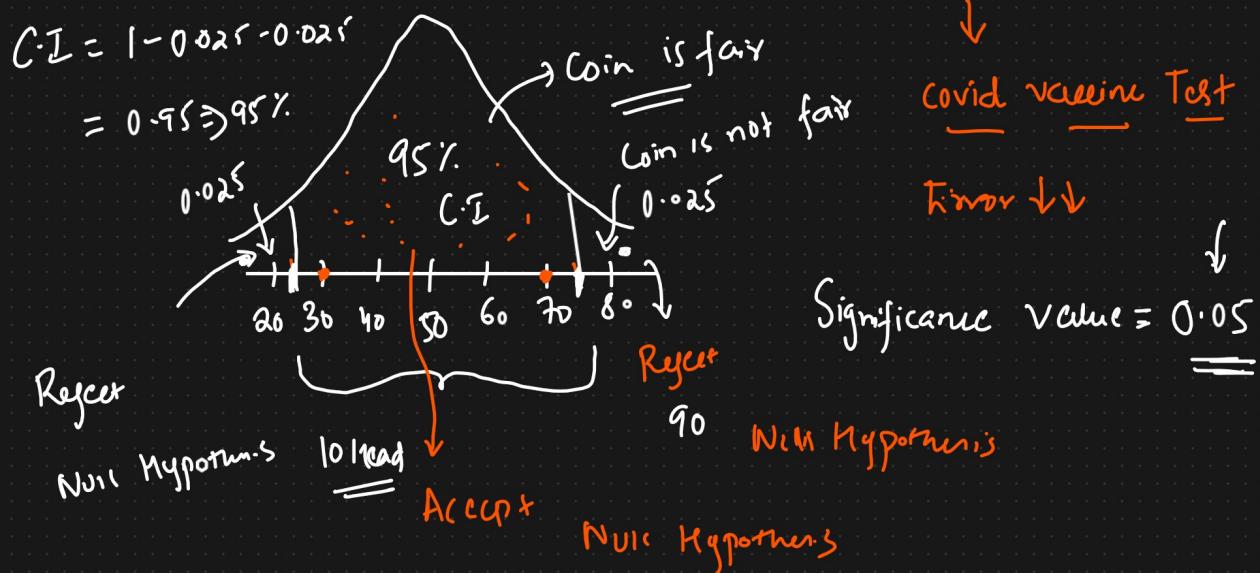
② Alternative Hypothesis — Coin is not fair  $\rightarrow (H_1)$

③ Experiments

④ Reject or Accept the Null Hypothesis



Confidence Interval, Significance Values

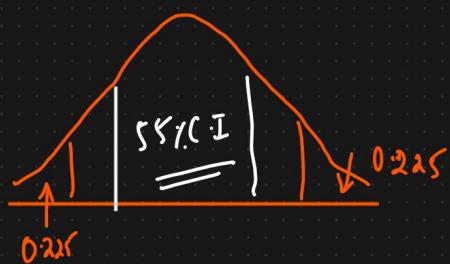


$$\lambda = 0.45$$

Medical

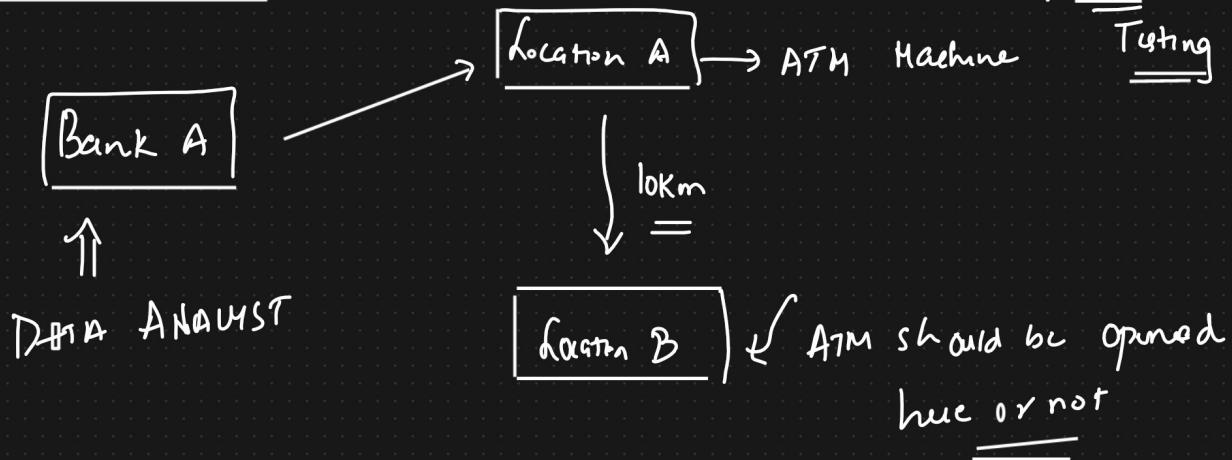
$f \uparrow \uparrow$

$$\frac{0.45}{2} = 0.225$$

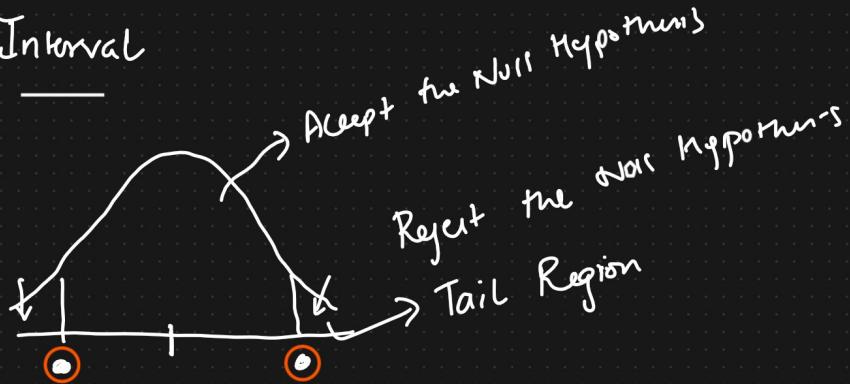


Real World Project

Hypothesis  
Testing

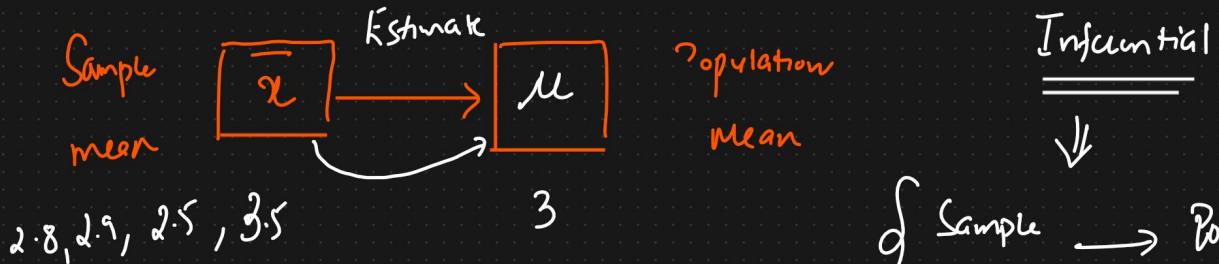


## ① Confidence Interval



## Point Estimate

{ The value of any statistic that estimates the value of a parameter is called Point Estimate.



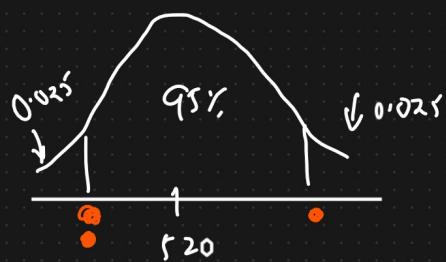
Inferring  
↓  
{ Sample Mean → Population }

## Confidence Interval

t test Point Estimate  $\pm$  Margin of Error  $\Rightarrow$  Population.

Q) On the quant test of CAT Exam, the standard deviation is known to be 100. A sample of 25 test takers has a mean of 520. Construct 95% CI about the mean?

Ans)  $\sigma = 100$   $n = 25$   $\bar{x} = 520$  ( $\alpha = 0.05$ )  $\alpha/2 = 0.025$



① Population std is given  $\{Z \text{ score}\} \rightarrow Z \text{ table}$

Point Estimate  $\pm$  Margin of Error  $\Rightarrow C.I.$

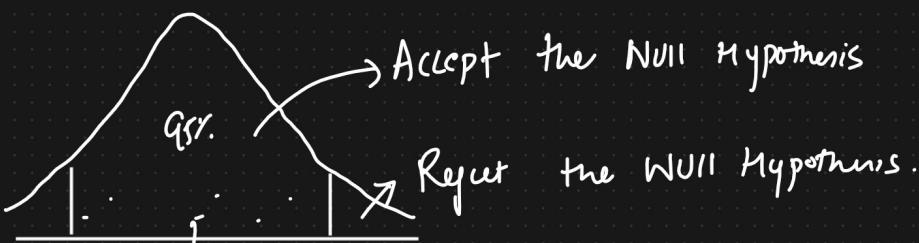
$\bar{x} \pm Z_{\alpha/2} \left[ \frac{\sigma}{\sqrt{n}} \right]$   $\rightarrow$  Standard Error

Lower fence  $C.I. = \bar{x} - Z_{\alpha/2} \left[ \frac{\sigma}{\sqrt{n}} \right]$   $\Rightarrow Z_{0.05} = 1.96$

Higher fence  $C.I. = \bar{x} + Z_{\alpha/2} \left[ \frac{\sigma}{\sqrt{n}} \right]$

Lower fence  $= 520 - (1.96) \times \frac{100}{\sqrt{25}} = 520 - (1.96) \times 20 = 480.8$

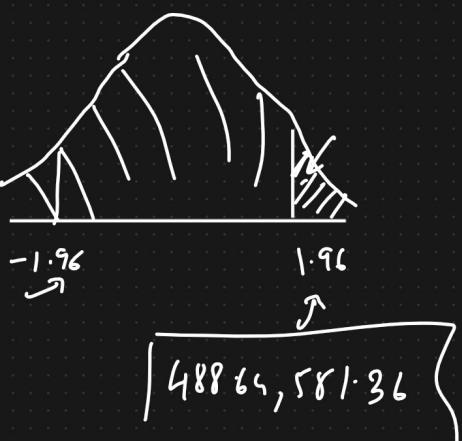
Higher fence  $= 520 + (1.96) \times 20 = 559.2$



$$480.8 \quad 520 \quad 589.2 \\ = \quad =$$

$$I_{\alpha/2} = Z_{0.025}$$

$$1 - 0.025 = 0.9750 \\ =$$



$$\boxed{488.64, 581.36}$$

- ④ On the quant test of CAT exam, a sample of 25 test-takers has a mean of 520 with a sample standard deviation of 80. Construct 95% C.I. about the mean? 2

Ans)  $\bar{x} = 520 \quad S = 80 \quad \alpha = 0.05 \quad n = 25$

$t$ -test  $\Rightarrow t$  - table { Because population   
  $sd$  is not given }

$$\bar{x} \pm t_{\alpha/2} \left( \frac{S}{\sqrt{n}} \right) \rightarrow \text{Standard Error}$$

$$t_{0.025}$$

$t$ -test

① Degree of freedom =  $n-1 = 25-1 = 24$   $\equiv$

$$3 \quad \frac{1}{n-1} \quad \boxed{3} \quad \boxed{1} \quad \boxed{3}$$

3 people

$$\bar{x} \pm 2.064 \left( \frac{80}{5} \right) \Rightarrow 486.976 \leftrightarrow 553.024$$

- ⑦ Type 1 and Type 2 Error.
- ⑧ One Tailed vs 2 Tailed Test

## Type 1 and Type 2 Error

### Reality Check

$H_0 \Rightarrow$  Coin is Fair

① Null Hypothesis is True or Null

$H_1 \Rightarrow$  Coin is not Fair

Hypothesis is False

#### Outcome 1:

#### Decision of Experiments?

We reject the Null Hypothesis Null Hypothesis is True or False.

in reality if it is false  $\rightarrow$  Yes



Null Hypothesis



$H_0 \rightarrow$  The Criminal is not guilty

$H_1 \rightarrow$  " " is guilty

#### Outcome 2:

We reject the Null Hypothesis

when in reality it is true  $\Rightarrow$  No  $\Rightarrow$  Type 1 Error  $\times$

#### Outcome 3:

We accept the Null Hypothesis,  $\Rightarrow$  Type 2 Error  $\times$

When in reality it is false

#### Outcome 4: We accept the Null Hypothesis

when in reality it is True

#### Confusion Matrix

$\begin{bmatrix} \downarrow & \\ \text{Cancer} & \rightarrow \\ \text{True} & \end{bmatrix} \rightarrow \begin{bmatrix} \text{Not Cancer} \\ \hline \end{bmatrix}$

$\left\{ \begin{array}{l} \rightarrow \text{Stock market is going to crash} \\ \end{array} \right\}$

## ② 1 Tail and 2 Tail Test

Eg: College is 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 has a different placement rate?

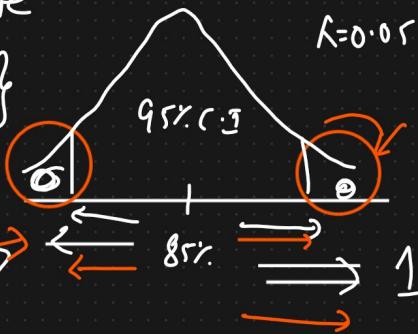
$$\alpha = 0.05 = 95\% \text{ C.I} \rightarrow 85\%$$

of placement rate

less than 85% }

↓

1 tail



of placement rate greater than 85% }

2 Tail Test

1 tail Test

Saturday

10 min probability

Sunday

① Z test Hypothesis Testing

② J Test Hypothesis Testing

③ Significance value of P value.

④ ANOVA TEST

⑤ CHI SQUARE TEST

⑥ Practical

EDA → 3-4 projects

FE → \_\_\_\_\_

Machine Learning

## ① Central Limit Theorem

## ② Influential Statistics

a) Z test {Z table} [5-6 problems]

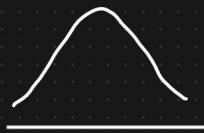
b) t test {t table}

c) Z test proportion population.

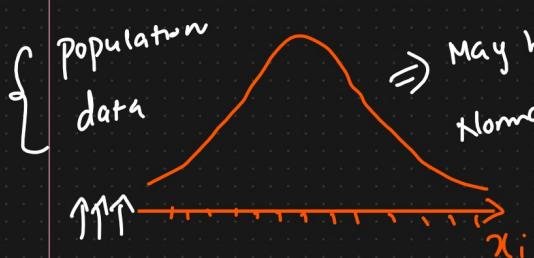
d) Chi Square (Categorical Test)

e) ANNOVA (F Test)

Influential



## ① Central Limit Theorem



May be Gaussian

$n > 30$

Normal Dist

Sample mean distribution

Sample 1  $[x_1, x_2, x_3, x_4, x_5, \dots, x_{30}] \rightarrow \bar{x}_1$

Sample 2  $[x_1, x_2, x_3, x_4, x_5, \dots, x_{30}] \rightarrow \bar{x}_2$

$\rightarrow \bar{x}_3$

$\rightarrow \bar{x}_4$

$\vdots$

$\rightarrow \bar{x}_m$

It may not

Sample m



10,

$n > 30$

Sample mean

distribution



Gaussian Distribution  
Normal Distribution

## ② Influential Statistics {Data Analyst, Data Scientist}

$$\text{Next Event} \div \text{Last } \bar{x} \downarrow \boxed{\underline{\mu, 5}} \curvearrowright$$

③ ATM    ④ Measure the size of entire sharks (I [ ]

⑤ Amazon delivery { Percentile, Quartiles } =

## \* Hypothesis Testing

① 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 80ml. Using 40 samples, he measures the average amount dispersed 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 machine is not working properly.

Ans) Step 1

$$n=40 \quad \bar{x}=78 \quad s=2.5$$

Step 1  $\Downarrow$   $H_0: \mu = 80$  {null hypothesis}  $\Downarrow$

2

$H_1: \mu \neq 80$  {Alternate Hypothesis} Why  $Z$  test?

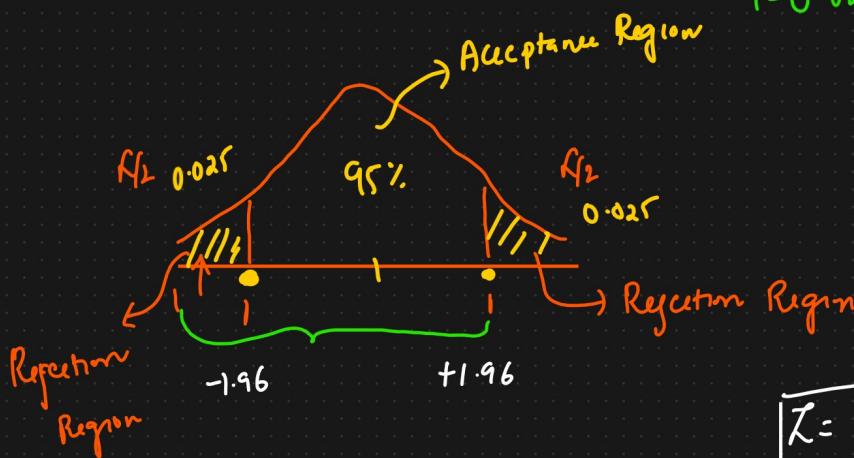
Step 2:

$$\alpha = 0.05$$

$$C.I = 95\%$$

$n > 30$   $\frac{n \leq 30}{}$   
 ①  $n > 30$   $\frac{n \leq 30}{}$   
 ② population std or sample std  
 $\frac{n > 30}{}$  std

Step 3: Decision Boundary



$$1 - 0.025 = 0.9750$$

Why  $Z$  test?  
 ① Sample std  
 ②  $n < 30$

$$n=1$$

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

④ Calculate Test Statistics

$$Z = \frac{\bar{x} - \mu}{\frac{S}{\sqrt{n}}} \Rightarrow \text{Standard Error}$$

Sample Standard deviation

$$\text{deviation} = \frac{78 - 80}{2.5 / \sqrt{40}} = \frac{-2 \times \sqrt{40}}{2.5} = \frac{-2}{2.5} \times 6.32 = \underline{\underline{-5.05}}$$

⑤ State the Results

Decision Rule: If  $Z = -5.05$  is less than  $-1.96$  or greater than  $1.96$ , then reject the null hypothesis with  $95\% C.I$

Reject  $H_0$  Null hypothesis  $\{$  There is some fault in the machine  $\}$

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.

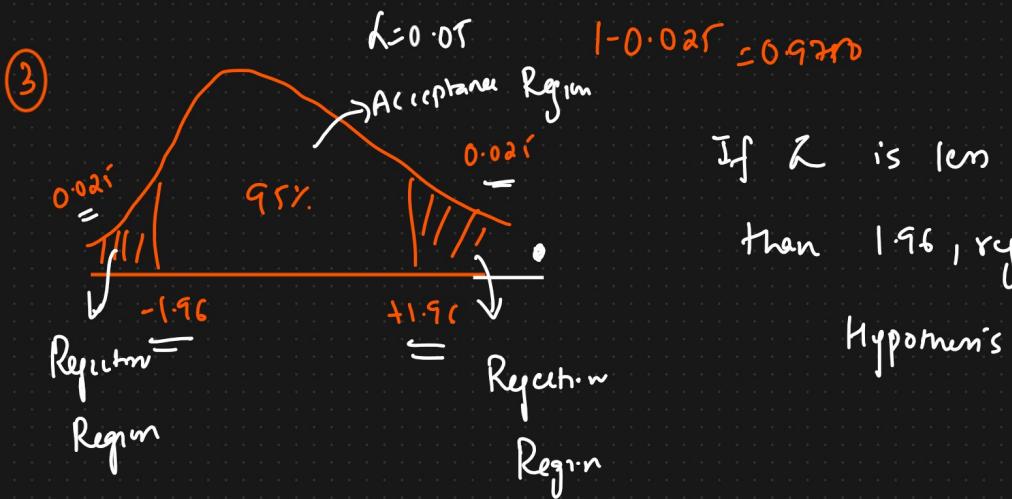
A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect Intelligence?  $\left. \begin{array}{c} 95\% \\ \hline \downarrow \\ C.I. \end{array} \right\}$

Ans)  $\sigma = 15 \quad n = 30 \quad \bar{x} = 140$

①  $H_0: \mu = 100$

$H_1: \mu \neq 100$

②  $\alpha = 0.05 \quad C.I = 95\%$



If  $Z$  is less than -1.96 or greater than 1.96, reject the null

④  $Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{140 - 100}{15 / \sqrt{30}} = 14.60 \quad \text{Reject Null Hypothesis}$

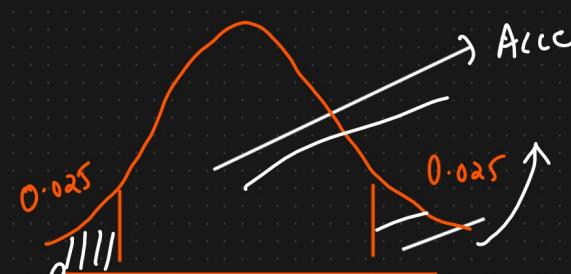
$14.60 > 1.96 \quad \text{Reject the Null Hypothesis}$

(\*) A Complain was registered, the boys in the Municipal Primary School are underfed. Average weight of boys of age 10 is 32 kgs with  $S.D = 9 \text{ kgs}$ . A sample of 25 boys was selected from the municipal School and the average weight was found to be 29.5 kgs? with  $C.I = 95\%$ . Check whether it's True or False?

Ans)  $\mu = 32 \text{ kgs}$   $\sigma = 9 \text{ kg}$   $n = 25$   $\bar{x} = 29.5$   $\alpha = 0.05$

1)  $H_0: \mu = 32$  }      2)  $\alpha = 0.05$        $1 - 0.95 = 0.05$

$$H_1: \mu < 32$$

3) 

Rejection =  $7.75\%$

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

$$= \frac{29.5 - 32}{9 / \sqrt{25}} = -1.39$$

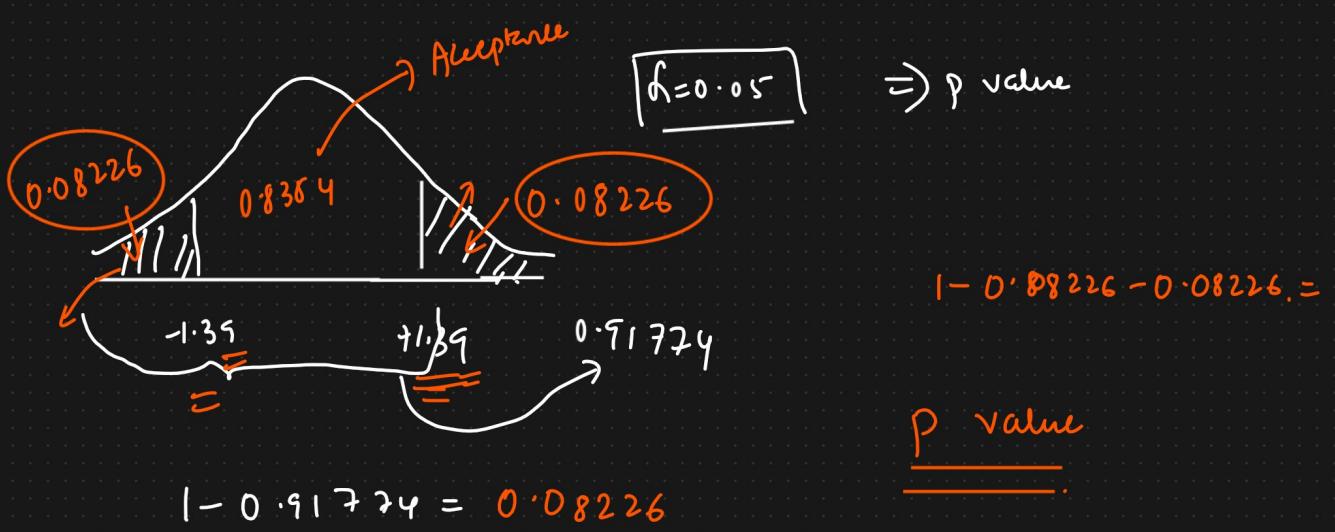
python  $\{Z \text{ stat, p-value}\}$

Conclusion :  $-1.39 > -1.96$  therefore we accept the Null Hypothesis

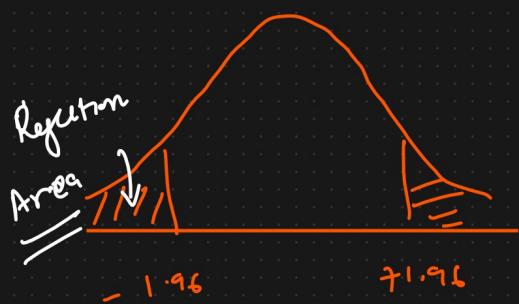
So, the boys are not underfed.

Significance value  $\Rightarrow$

P value



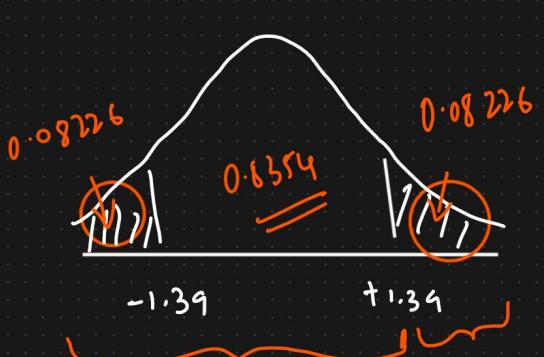
Significance value  
↑  
 $0.1645 > 0.05$



Z tut

14

P value



$$\text{New } \hat{P} \text{ value} = 0.08226 + 0.08226$$

|-0.08226 -0.08226

## Domain

2/2

0.1645 > Significance

value

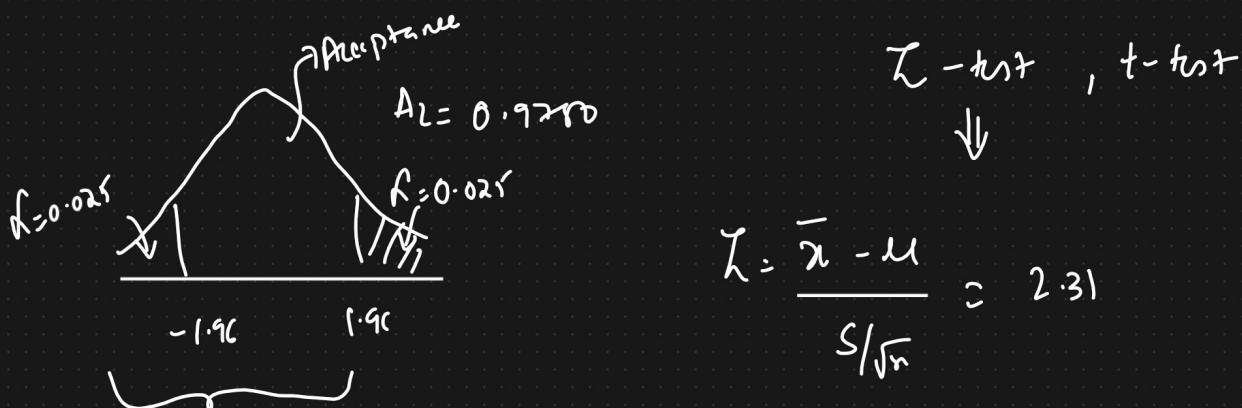
Accept the Null Hypothesis.

4) 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% CI is there enough evidence to discard the Null Hypothesis??

Ans)  $H_0 : \mu = 168$        $n = 36$        $\bar{x} = 169.5$        $s = 3.9$

$H_1 : \mu \neq 168$        $\bar{x} =$        $\alpha = 0.95$        $\beta = 1 - C.I = 0.05$



$$Z = \frac{\bar{x} - \mu}{s/\sqrt{n}} = 2.31$$

$2.31 > 1.96$  Reject the Null Hypothesis

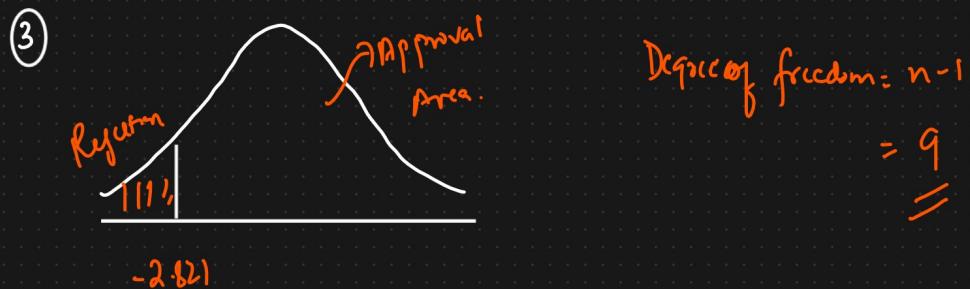
5) A company manufactures bike batteries with an average life span of 2 or more years. An engineer believes this value to be less. Using 10 samples, he measures the average life span to be 1.8 years with a standard deviation of 0.15.

a) State the Null and Alternative Hypothesis

b) At a 99% CI, is there enough evidence to discard the  $H_0$ ?

Ans)  $H_0 : \mu \geq 2$        $n=10$        $\bar{x}=1.8$        $S=0.15$        $\{$  of sample  
 $H_1 : \mu < 2$        $\leq 30$        $t-tst??$       Std is  
 $\{$  given }

②  $\alpha = 0.01$        $\alpha = 1 - C.I = 1 - 0.99 = 0.01$   $\equiv$



④ Calculate  $t$ -test Statistic:

$$t = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{1.8 - 2}{0.15/\sqrt{10}} = \frac{-0.2}{0.15/\sqrt{31.622}} = -4.216 \equiv$$

⑤ Conclusion

$-4.216 < -2.82$       Reject the Null Hypothesis.  $\{$   
 $\Downarrow$   $\equiv$

Z test with proportions

⑥ 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

(a) State the Null and Alternative Hypothesis?

(b) At a 95% C.I, is there enough evidence to reject the Null Hypothesis?

Ans)  $H_0: p_0 = 0.70$ .

$H_1: p_0 \neq 0.70$

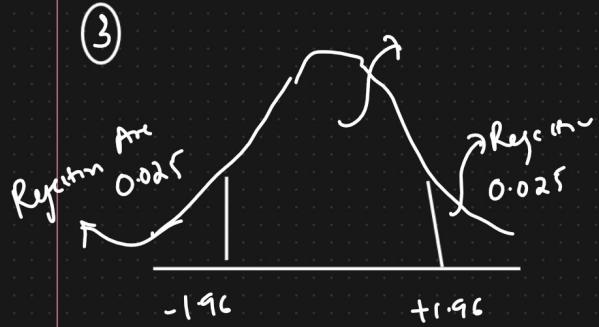
$$n = 200 \quad \chi = 130$$
$$\hat{p} = \frac{\chi}{n} = \frac{130}{200} = \frac{13}{20} = 0.65$$

$$q_0 = 1 - p_0$$

②  $\alpha = 0.05 \quad C.I = 95\%$

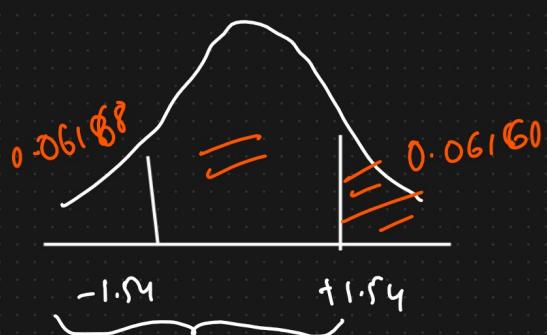
$$Z_{\text{test}} = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0 q_0}{n}}}$$

$$= \frac{0.65 - 0.70}{\sqrt{\frac{0.7 \times 0.3}{200}}} \approx -1.54$$



At 95% C.I there is

$-1.54 > -1.96$ , So we accept  
the Null Hypothesis



pvalue  
 $2 \times 0.06168 > 0.05$

Accept Null Hypothesis

$$1 - 0.93822 = 0.06168$$

④ A car company believes that the percentage of residents in City ABC that owns a vehicle is 60% or less. A sales manager disagrees with this. He conducts a hypothesis testing surveying 250 residents and found that 170 responded yes to owning a vehicle.

- (a) State the Null & Alternate Hypothesis
- (b) At 10% significance level, is there enough evidence to support the idea that vehicle ownership in City ABC is 60% or less?

$$p\text{ value} = 0.014$$

---

---

# Statistics

{ 11:30 - 12pm }

- ① Covariance
- ② Pearson Correlation Coefficient
- ③ Spearman Rank Correlation Coefficient
- ④ CHI SQUARE TEST
- ⑤ ANNOVA (F-Test)

✓  
✓ practicals  
✓

## Covariance

$x \uparrow \quad y \uparrow$

$x =$

$y =$

$x \uparrow \quad y \downarrow$

-

-

$x \downarrow \quad y \uparrow$

-

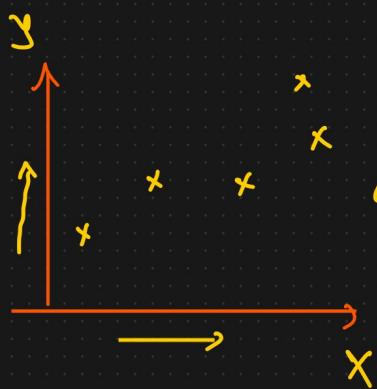
-

$x \downarrow \quad y \downarrow$

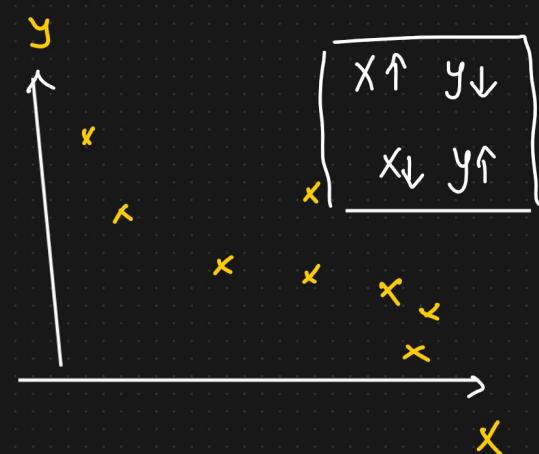
-

-

{ quantity the relationship  
between  $x$  &  $y$  }



$\left\{ \begin{array}{l} x \uparrow \quad y \uparrow \\ x \downarrow \quad y \downarrow \end{array} \right.$



$$\text{Cov}_{x,y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N-1} \Leftrightarrow \text{Var}_x(x) = \frac{\sum (x_i - \bar{x})^2}{N-1}$$

$\text{Cov}(x,y)$

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

$$= \frac{\sum (x_i - \bar{x}) \times (x_i - \bar{x})}{N-1}$$

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

↓

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

+ve  $\Rightarrow \Rightarrow \Rightarrow \Rightarrow$   
 $\Rightarrow$  Positively Correlation

$x \uparrow y \uparrow$   
 $x \downarrow y \downarrow$

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

$$\left. \begin{aligned} &= (2-4)(3-5) + (4-4)(5-5) \\ &\quad + (6-4)(7-5) \end{aligned} \right. \overline{x} = 4 \quad \overline{y} = 5$$

2

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

$x \uparrow y \downarrow$   
 $x \downarrow y \uparrow$

$\Rightarrow$  -ve Correlation  $\Rightarrow$  -ve value.

Disadvantage Covariance

$\text{Cov}(x, y) \Rightarrow$  +ve value  
 or -ve value

↓

Relationship  $[-1 \rightarrow 1]$

$$\text{Cov}(x, y) = 500$$

$$\text{Cov}(y, z) = 600$$

Limit  $-400$   
 $+500$   $-300$   
 $-400$   $+1000$

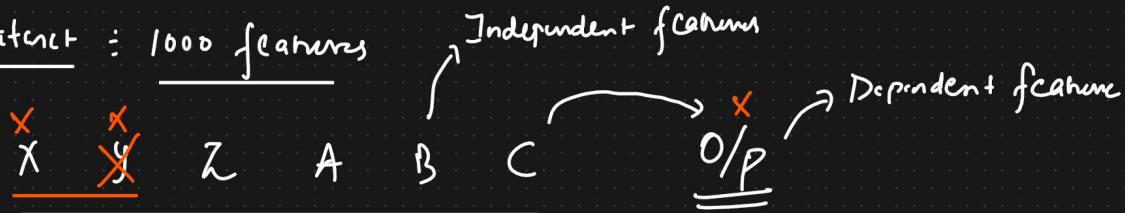
$\boxed{\infty}$

## ② Pearson Correlation Coefficient

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

The more the value towards 1 more the it is correlated

Dataset : 1000 features



+ve Correlation

$$x, y \Rightarrow 99\% \quad =$$

$$9 \text{ or } 0.9 \quad =$$

-ve Correlation

↓  
Keep it

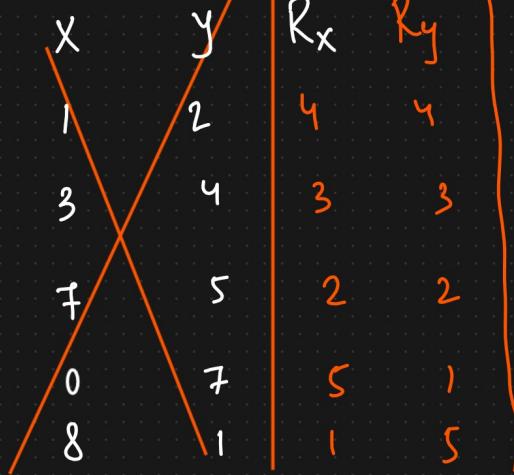
## ③ Spearman Rank Correlation

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma_{R(x)} \sigma_{R(y)}}$$

Marks

Spearman Rank

$$\text{Corr} = 1 \quad =$$



$$-1 \quad =$$

## Chi Square

The Chi Square Test claims about population proportions.

It is a non parametric test that is performed on categorical (nominal 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 the population distribution of ages has changed in the last 10 years?

Ans)	<18	18-35	>35
Expected	20%	30%	50%

$n=500$

Observed : 121 288 91

95% C.I  
=

Expected 100 150 250

- ①  $H_0$  = the data meets the expected distribution  
 $H_1$  = the data does not meet the expected distribution

- ② State Alpha  $\therefore \alpha = 0.05$

- ③ Calculate the degree of freedom

$$df = n - 1 = 3 - 1 = 2 \Rightarrow 3 \text{ categories.}$$

- ④ Decision Chi Square Table.

If  $\chi^2$  is greater  $\underline{\underline{5.99}}$  than, Reject  $H_0$

- ⑤ Calculate Chi square Test

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e} = \frac{(121 - 100)^2}{100} + \frac{(288 - 150)^2}{150} + \frac{(91 - 250)^2}{250} \\ \chi^2 = 232.494$$

$232.494 > 5.99$  Reject the null hypothesis.

- ⑥ 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 these results, do the days for the highest number of absences occur with equal frequencies (use 95% C.I.)

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

$$\text{Ans} = \frac{6.3}{\text{---}} \quad \left\{ \text{Accept the Null Hypothesis} \right\}$$

Practicals + EDA + Feature Engineering }

## Statistics

- ① ANOVA (F-Test)  $\rightarrow$  1 hour }  
② FDD  $\rightarrow$  { Solve Some Examples } 

ANOVA : { Analysis of Variance }

ANOVA IS a statistical method used to compare the means of 2 or more group

ANOVA :

① Factors      ② Levels  
(variables)      { Dosage }      Anxiety reducing      Gender  
Medicine           0mg      50mg      100mg      M      F  
                     $\bar{x}$        $\bar{x}$        $\bar{x}$       =      =

factor : Dosage

levels : 0mg, 50mg,  
100mg

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

Types of ANOVA

One Way ANOVA : One factor with at least 2 levels, levels are independent.

② Repeated Measures ANOVA - One factor with at least 2 levels, but levels are dependent

Factor	Running Kms
Levles	1
	2

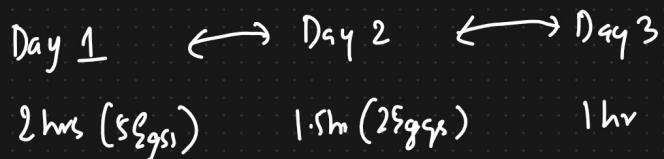


Ques. Study hours of KARTIK

#### ④ Factorial ANOVA



Gym



⑤ Factorial ANOVA  $\div$  Two or more factor (each of which with at least 2 levels), levels can be either independent, dependent or both (mixed)

factor

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

One Way ANOVA ( $F$ -test)  $\Rightarrow$  Inferential stats



Comparing means of 2 or more groups

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

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

①  $H_0 = \mu_{0\text{mg}} = \mu_{50\text{mg}} = \mu_{100\text{mg}}$  }  
 $H_1 = \text{not all } \mu\text{'s are equal}$  }

② State  $\alpha$  and C.I

$$\alpha = 0.05 \quad C.I = 95\%$$

③ Calculate the Degree of freedom

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

$$a = 3 \rightarrow \{ \text{No. of levels} \}$$

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

$$\rightarrow df_{\text{Total}} = N-1 = 21-1 = 20$$

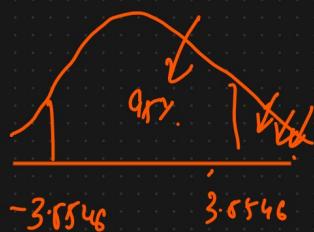
$$N = 21 \quad n = 7$$

Statistics

#### ④ State Decision Rule

$$df_{\text{Between}} = a-1 = 3-1 = 2 \quad \{(2, 18)\}$$

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



If F test is greater than 3.8846, Reject the Null Hypothesis

If F test is less than -3.8846 " " " "

#### ⑤ Calculate F Test Statistics

$$F_{\text{test}} = \frac{MS_{\text{between}}}{MS_{\text{within}}} = \frac{49.34}{0.57} =$$

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

$$SS_{\text{between}} = \frac{\sum (\sum a_i)^2}{n} \quad \frac{T^2}{N} \quad N=21 \quad n=7 \text{ //} \quad T^2 = [57 + 47 + 21]^2 \\ = (125)^2$$

$$\sum (\sum a_i)^2 = (9+8+7+8+8+9+8)^2 + (7+6+6+7+8+7+6)^2 + (4+3+2+3+4+3+2)^2 \\ = 57^2 + 47^2 + 21^2$$

$$SS_{\text{Between}} = \frac{57^2 + 47^2 + 21^2}{7} - \frac{125^2}{21} = 98.67 =$$

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

$$\left. \begin{array}{l} P = 0.48 \\ d.f. = 0.05 \end{array} \right\} = \sum y^2 - \frac{\sum \downarrow}{853} \left[ \frac{57^2 + 47^2 + 21^2}{7} \right] = 10.29$$

$$\sum y^2 = 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + 9^2 + \dots + 2^2 = 853$$

$$\frac{0.75 > 0.05}{\Downarrow}$$

Final Conclusion

Accept

~~86.56 > 35846~~ So we reject the Null hypothesis?

$$\left. \begin{array}{l} H_0: \mu = \text{Some value} \\ H_1: \mu \neq \text{Some value} \end{array} \right\} \rightarrow 95\% \text{ CI}$$

Virginia

—

—

Petal width

—

—

—

—

—

—

—

—

—

—

—

—

—

—

$$\rightarrow H_0 = \mu_{\text{virgin}} = \mu_{\text{swiss}} = \mu_{\text{...}}.$$

$H_1 = \cdot \neq \text{pvalue.} \neq \cdot \rightarrow \text{Reject the Null Hypothesis}$

$$0.0118 \quad 0.0228 < 0.05 \quad 1-0.025 = 0.975$$

$$0.0118 \quad \underline{\underline{0.0228}}$$

$Z_{\text{test}}$

$d =$

$Z_{\text{test}} \neq \text{standard}$

