

Statistics

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

Data: "facts or pieces of information"

Eg: Height of students in a classroom
 $\rightarrow \{175\text{cm}, 150\text{cm}, 140\text{cm}, 130\text{cm}, 155\text{cm}\}$

Eg: Intelligence Quotient (IQ) of 5 randomly selected individuals ($109, 89, 129, 101, 105, 106$) \rightarrow Data.

Two Types

Statistics



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.

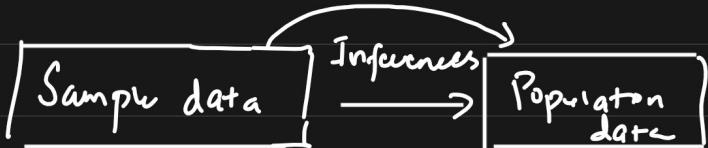
Ages $\{21, 20, 18, 34, 17, 22, 24, 25, 26, 23, 22\}$

$\min = \text{mode}$

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?

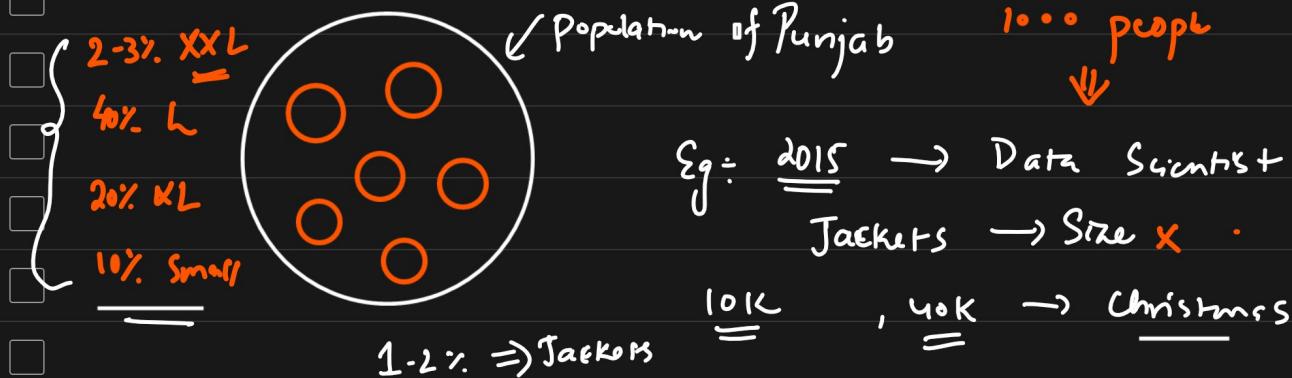


Population And Sample Data → Inferrential Statistics Results

Elections : Punjab

{ AAP, Congress }

Exit Polls ←

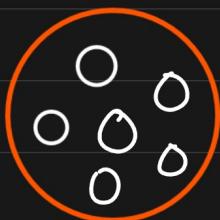


Population (N) ✓

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 → Layers
↓
Clusters
↑
Non overlapping groups

Gender → Male
Female

Blood Groups

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

Tax Slabs
Courses

Education Qualification

Thamno

Australias

③ Systematic Sampling

Snap

Customs

(N) → Select every n^{th} individual

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

\downarrow
Stratified

Eg: Survey → Mail (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 derive categorical variables }
Female]

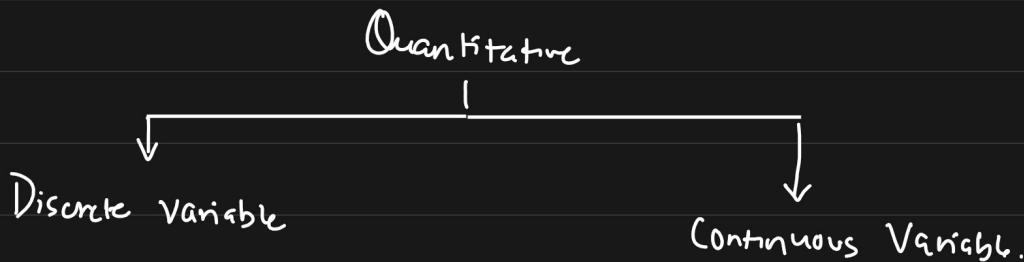
{ Quantitative → Qualitative Variable. }

Eg: IQ

0-10 10-50 50-100

↓ ↓ ↓

Low IQ Medium IQ Good IQ



Eg: whole number

Eg: No. of Bank accounts

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

Eg: Total No. of children in a family

Eg: Height = 172.5, 162.5 cm,

163.5 cm.

Rainfall: 1.35, 1.25, 1.75, 2.25 cm

Weight

Temp

Eg: 2, 3, 4, 5

Stock price.

25, 2.75

Eg: Total no. of Employees in a Company {e.g.: 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

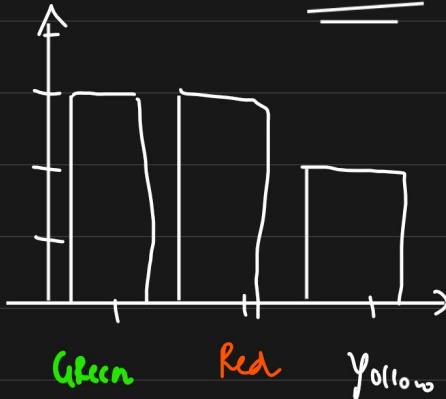
Sample Data : Green, Red, Yellow, Green, Red, Yellow, Green, Red

↓

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.

pHd
↓
{ NLP }
↓

Dance

phd

1

Sabu

✓

B.G

3

✓

Master

2

✓

BCA

4

✓

12

5

✓

Assignment

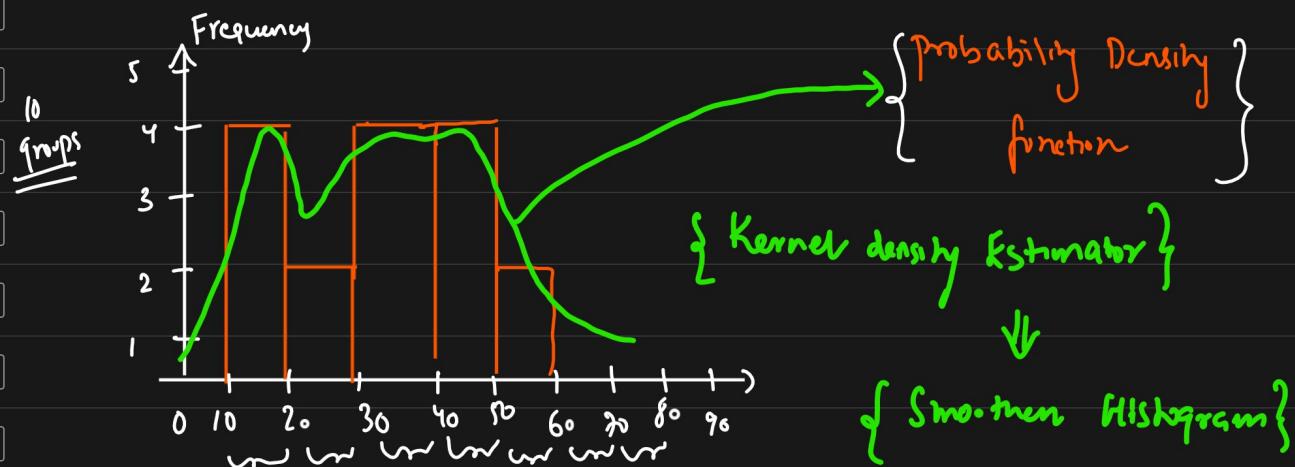
④ Ratio data ✓

③ Interval data ✓

⑤ Histograms ÷ Continuous

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

Histogram $\rightarrow \text{Bins} = 10$ \equiv Mean, Median, Mode.



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

Assignment

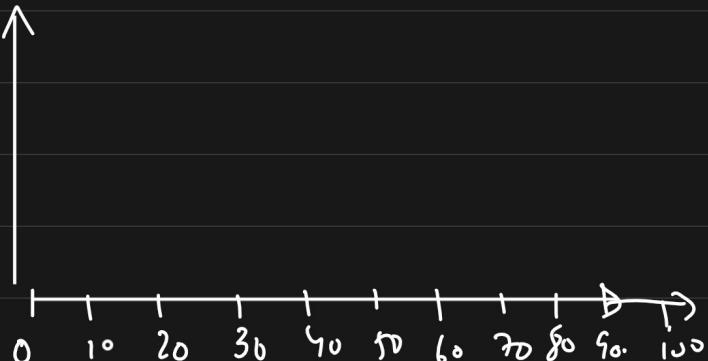
Eg: 10, 13, 18, 22, 27, 32, 38, 40, 45, 51, 56, 57, 88, 90, 92, 94, 99

bins \downarrow
10

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

40-50 50-60 60-70

70-80 80-90 90-100

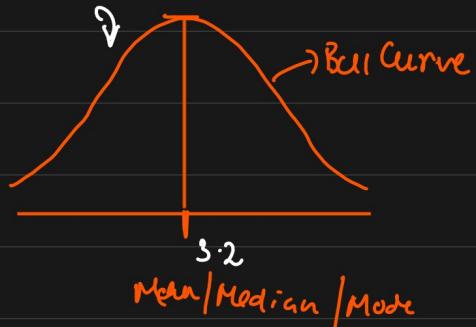


Intermediate Stats

- 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)

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

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

Population

mean

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

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

Sample (n)

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

Sample

Mean

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
↓ $\frac{2+3}{2} = 2.5$

2.5 $\approx 2.83 //$

Mode ÷ Highest frequency: Median

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

↓
3 ↓

EDA

1, 2, 2, 3, 3, 4, 4, 5, 5
↓
{mode}

[2, 3, 4]

Feature Engineering

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

⇒ Categorical Variable.

↓
Mode

Agnl

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

Measure of Dispersion → {Dispersion}

① Variance

② Standard deviation

} \Downarrow

Spread ⇒ How the data is spread



① Variance

Population Variance

{
Basis (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>X</u>	<u>\bar{x}</u>	<u>$x - \bar{x}$</u>	<u>$(x - \bar{x})^2$</u>
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

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

\uparrow

$n=6$

$n-1$

$$\mu = 2.83$$

{let consider

$$10.84$$

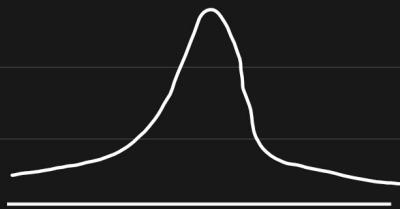
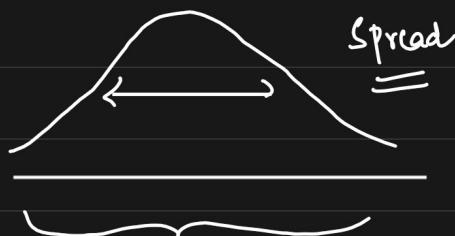
$$\frac{\sigma^2}{\text{Variance}} = \frac{6.42}{\text{as an example}}$$

Spread ↑↑

$$\frac{\sigma^2}{\text{Variance}} = 2.168$$

Variance ↑↑

Spread ↑↑



Standard deviation

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

$$= \sqrt{1.472}$$

Variance

\downarrow

Spreadness

1, 2, 2, 3, 4, 5

2.83

-1.472

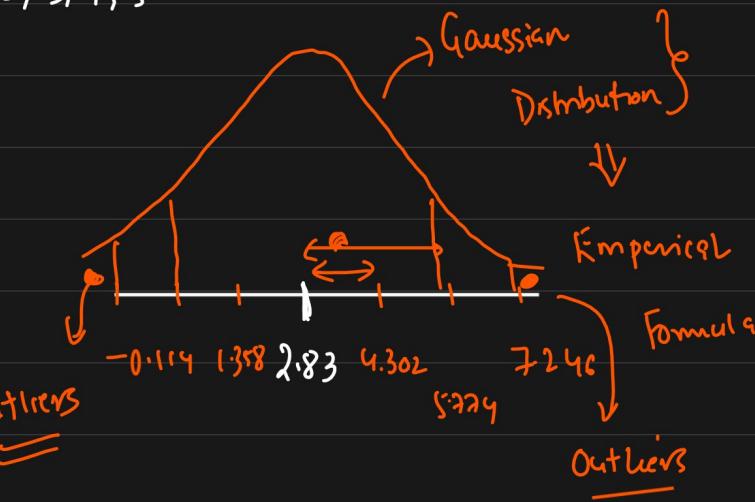
$\frac{1.358}{1.358}$

1.358

1.472

$\frac{1.358}{0.114}$

0.114



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

$$\begin{array}{r} 1.472 \\ 5.246 \\ \hline 7.246 \end{array}$$

(*) Percentiles And Quartiles



Percentages : 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? $n=20$

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

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

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

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]

$$\text{Lower Fence} = Q_1 - 1.5(\text{IQR}) \quad (25\%) \quad Q_1 = \frac{3+5}{2} \times 20^{\text{th}} \text{ index}$$

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

$$\text{IQR} = Q_3 - Q_1 = 7 - 3 = 4 \quad (75\%) \quad Q_3 = \frac{7+8}{2} \times 20^{\text{th}} = 15^{\text{th}} \text{ index}$$

$$Q_3 = 7$$

$$\text{Lower Fence} = 3 - 1.5(4) = 3 - 6 = -3$$

$$\text{Higher Fence} = 7 + 1.5(4) = 7 + 6 = 13$$

$$[-3 \longleftrightarrow 13] \quad -\text{ve} \quad \underline{\underline{3}}$$

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

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

5 Number Summary

Minimum = 1

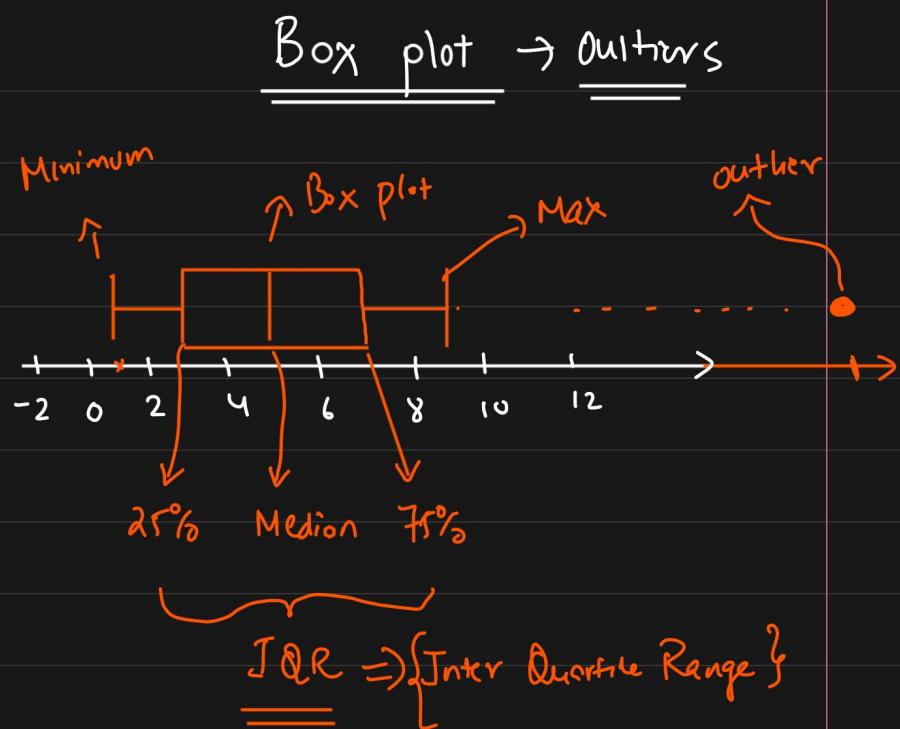
Q₁ = 3

Median = 5

Q₃ = 7

Max = 9

Use of Box plot



① Distributions

① Normal / Gaussian Distribution ✓

② Standard Normal Distribution ✓

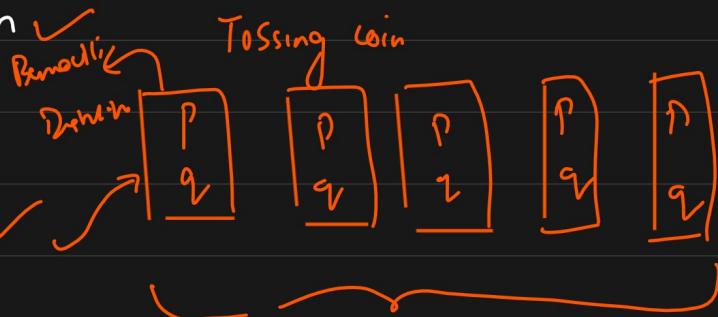
③ Z-Score ✓

④ Log Normal Distr ✓

⑤ Bernoulli's Distribution ✓

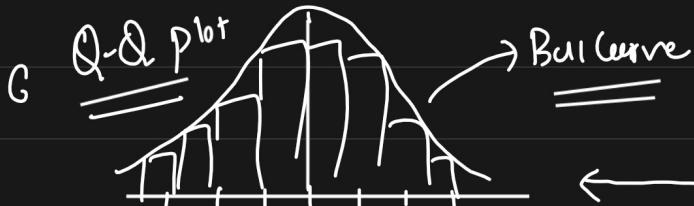
⑥ Binomial Distribution }

① Gaussian / Normal Distribution



Properties

{ Power Law }



① Empirical Rule of
Gaussian Distribution \Downarrow
80-20%

↳ DATASET → IRIS Dataset } → Petal, Sepal length
domain expansion }

② Weight of human brain

③ Height → Doctor

$$68.2, -95.4, -99.7$$

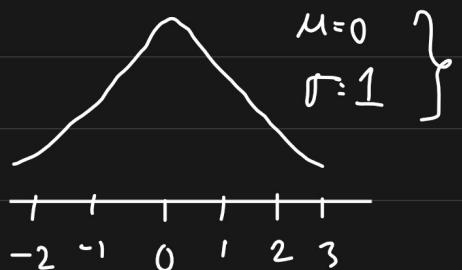
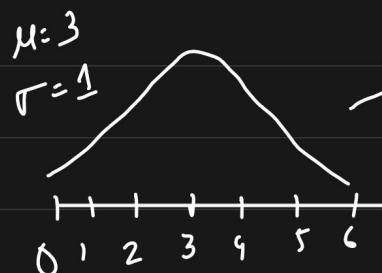
Outliers

Standard Normal Distribution

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

$$\mu = 3$$

$$\sigma = 1.414 \approx 1$$



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

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

$$\begin{aligned} &= \frac{3-3}{1} = 0 \\ &= \frac{2-3}{1} = -1 \\ &= \frac{1-3}{1} = -2 \end{aligned}$$

$$\text{Why } 22 \quad \boxed{\begin{array}{l} \mu = 0 \\ \sigma = 1 \end{array}}$$

✓

Standardization vs Normalization

Years
Age ↑
Different unit
Weight ↗ kg

$$\begin{array}{ll} \text{Age} & \text{Weight} \\ 25 & 75 \end{array}$$

$$26 \quad 80$$

$$28 \quad 85$$

$$30 \quad 60$$

$$32 \quad 70$$

$$\text{un} \quad \text{un}$$

INR
Salary
25K
30K

$$40K$$

$$80K$$

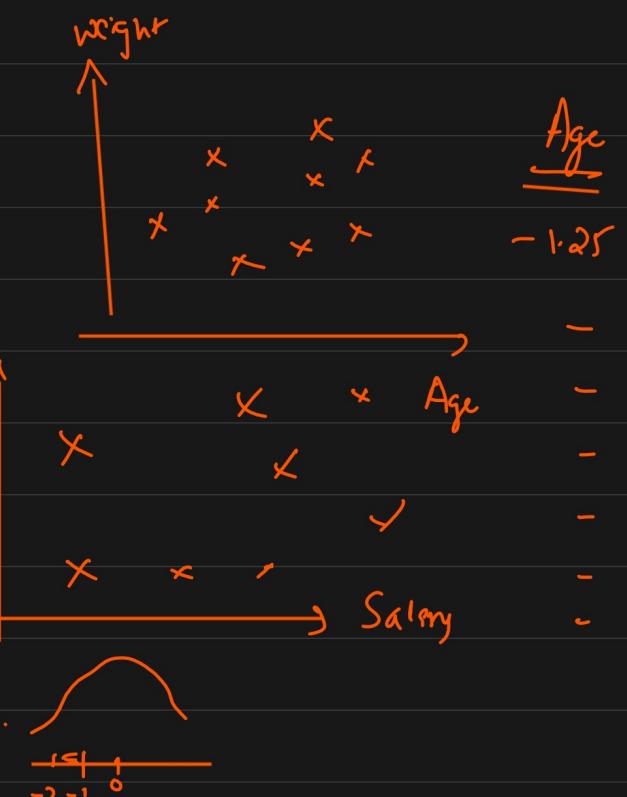
$$\text{un}$$

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

Same unit scale ??

Maths → Scale

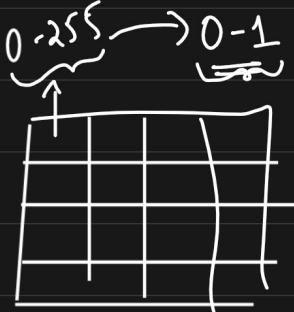
Standardization



Normalization

[Min Max Scalar]

$$\begin{matrix} \downarrow \\ 0 \text{ to } 4 \end{matrix}$$
$$\left. \begin{matrix} \downarrow \\ 0 \text{ to } 1 \end{matrix} \right\}$$



Convolutional

Neural Netw.

ML Disease
✓
Standardization

g

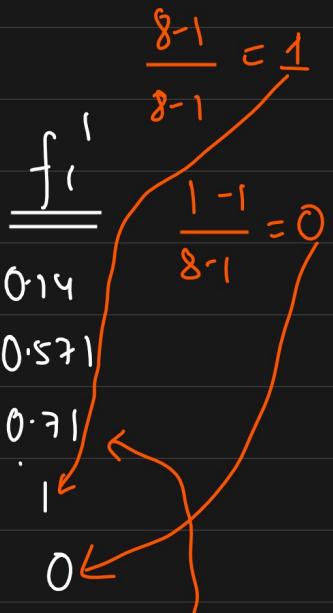
Normalization

$\leftarrow \boxed{\text{CNN}}$

f1

Normalization

(0 - 1)



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

1

Min Max Scalar

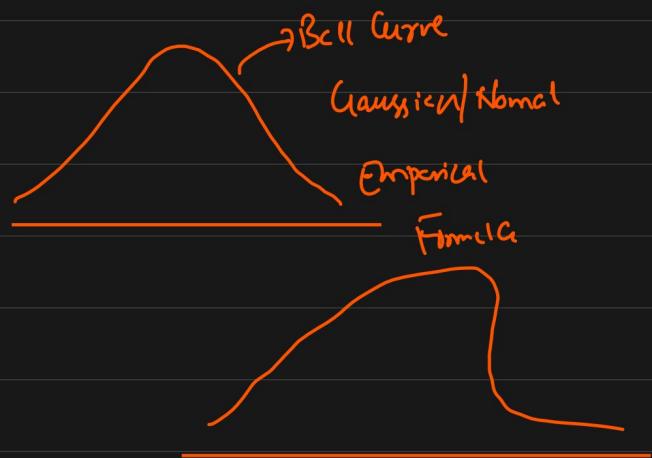
0 to 1

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

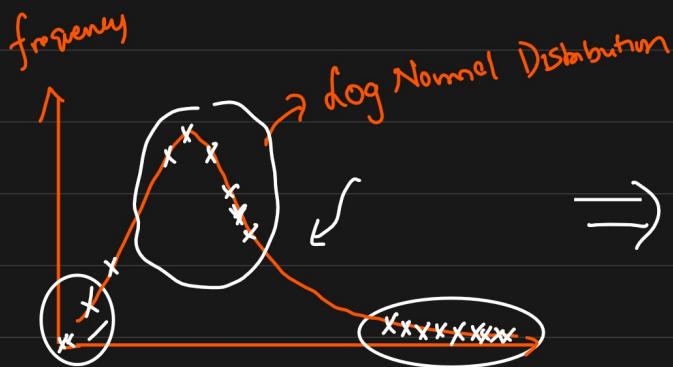
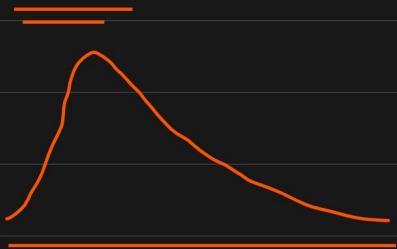
$$\frac{6 - 1}{8 - 1} = \frac{5}{7}$$

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

Log Normal Distribution

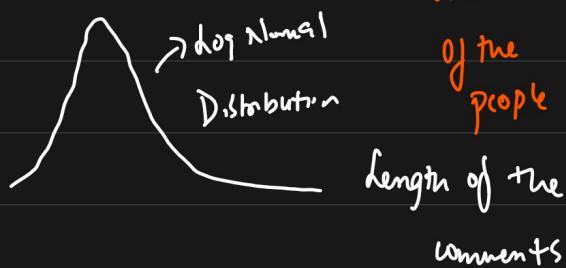
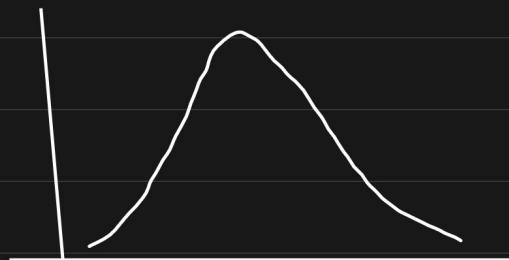


Skewed Curve



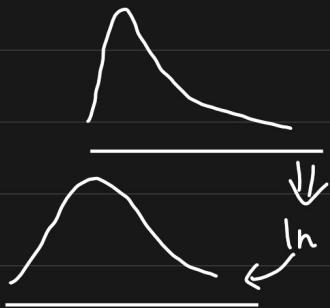
Gaussian Distribution

Normal Distn.



$X = \text{log Normal Distributed}$

$$\left\{ Y = \ln(X) \right\} \quad \begin{array}{l} \text{Gaussian} \\ \text{Distribution} \end{array}$$



$$\left\{ X = \exp(Y) \right\} \rightarrow c^y$$

X

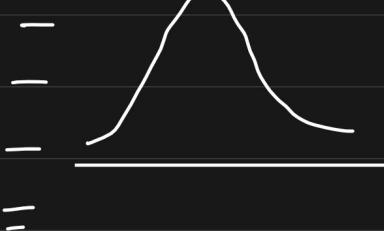
$\mathcal{Y} = \ln(x)$

25

30

40

45

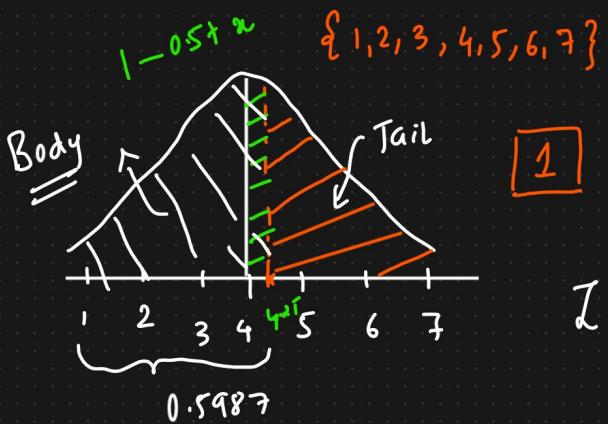


① Bernoulli's Distribution

Day 2 - Stats

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

Stats Interview Question



How many standard deviation

4.25 fall from the mean??

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

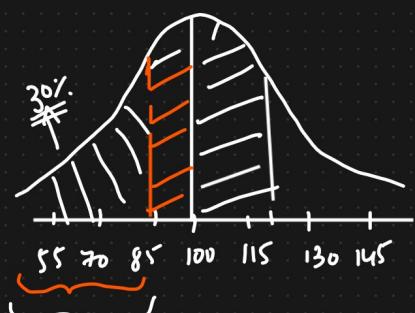
Question : What percentage of scores fall above 4.25?

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

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?

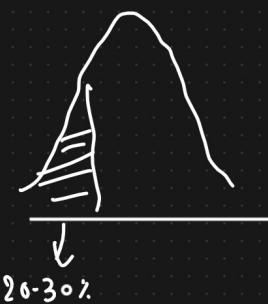
Ans)



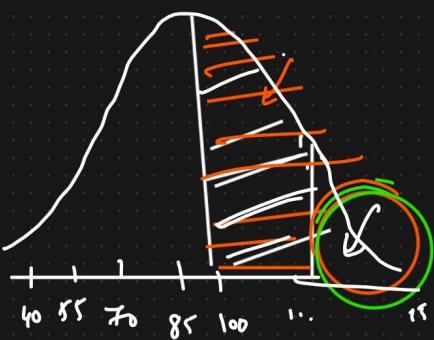
$$Z\text{-Score} = \frac{85 - 100}{15} = \frac{-15}{15} = \boxed{-1}$$

① Area under this curve

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



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

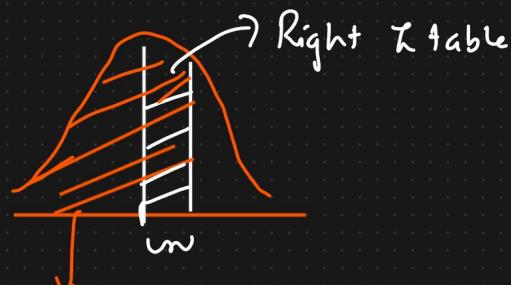


$$Z\text{score} = \frac{125 - 100}{15} = \frac{25}{15} = 1.667$$

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

1.667

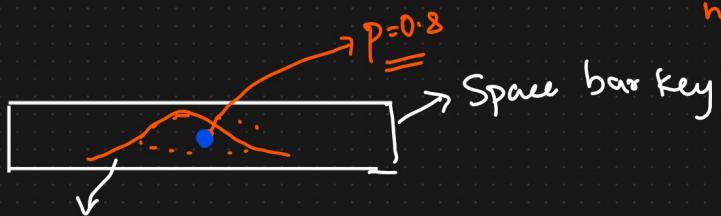
$$\underline{0.5 - 0.4515 = 0.0485} \Rightarrow 4.8\%$$



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 times 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

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

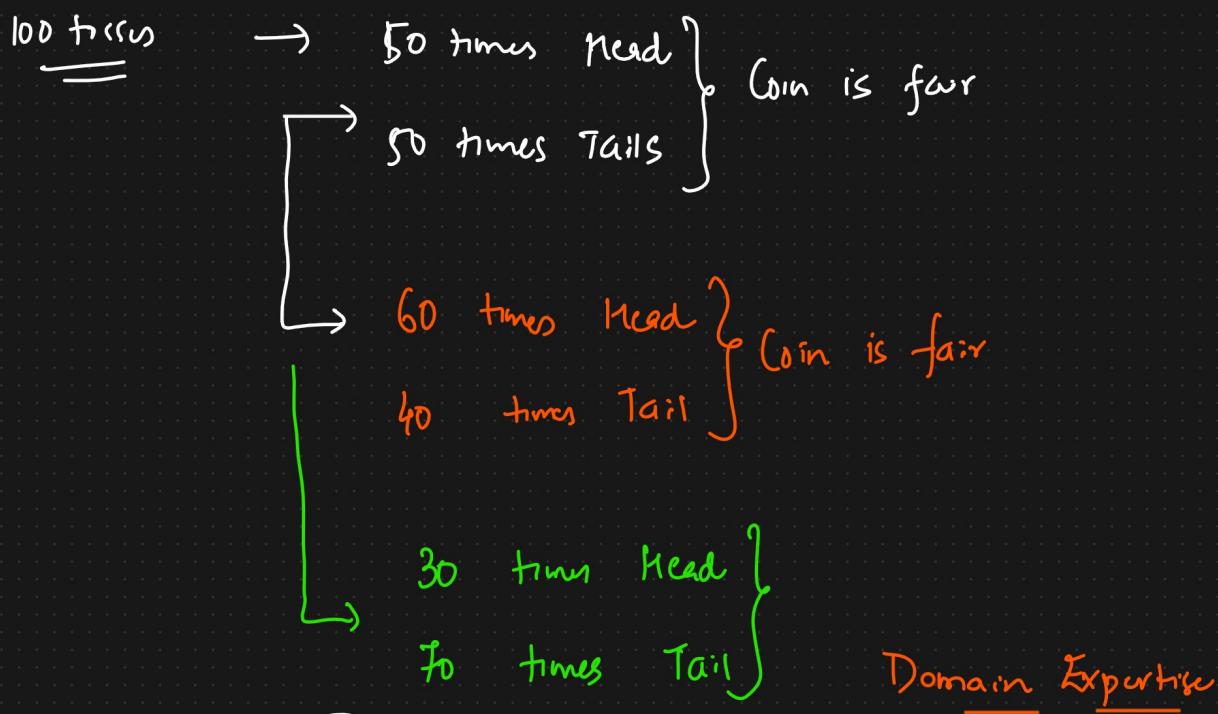
Hypothesis Testing

Criminal is → Court

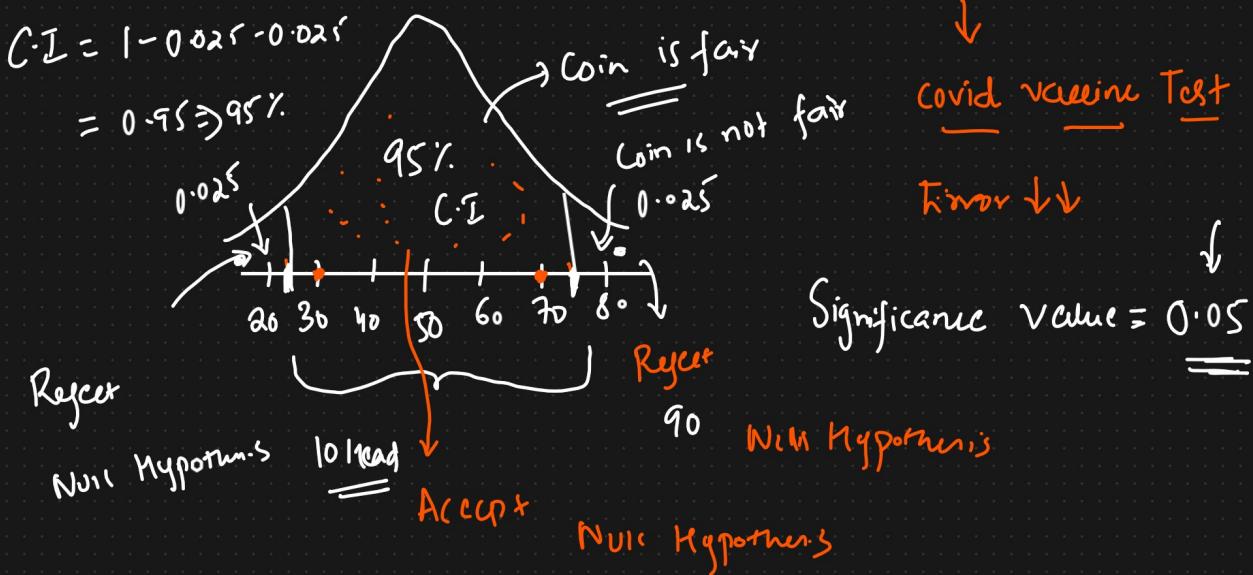
S110LAY

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

- ① Null Hypothesis — Coin is fair $\rightarrow (H_0)$
 - ② Alternate Hypothesis — Coin is not fair $\rightarrow (H_1)$
 - ③ Experiments
 - ④ Reject or Accept the Null Hypothesis



Confidence Interval, Significance Values

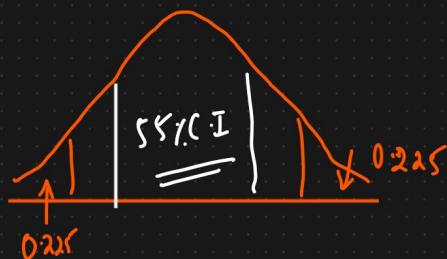


$$\lambda = 0.45$$

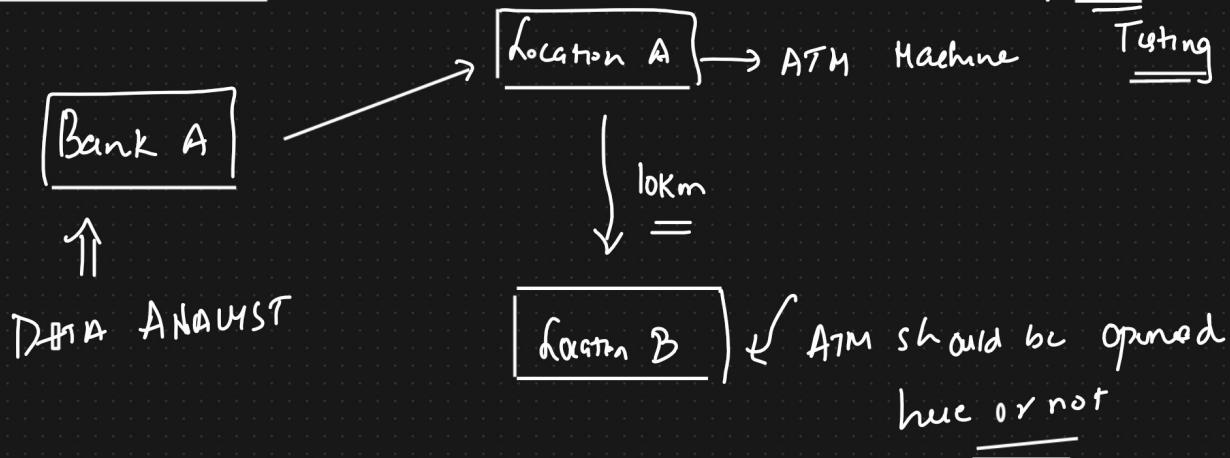
Medical

$f \uparrow \uparrow$

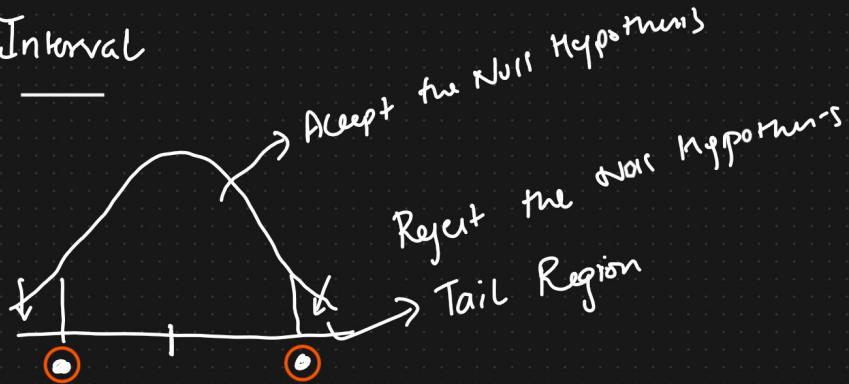
$$\frac{0.45}{2} = 0.225$$



Real World Project

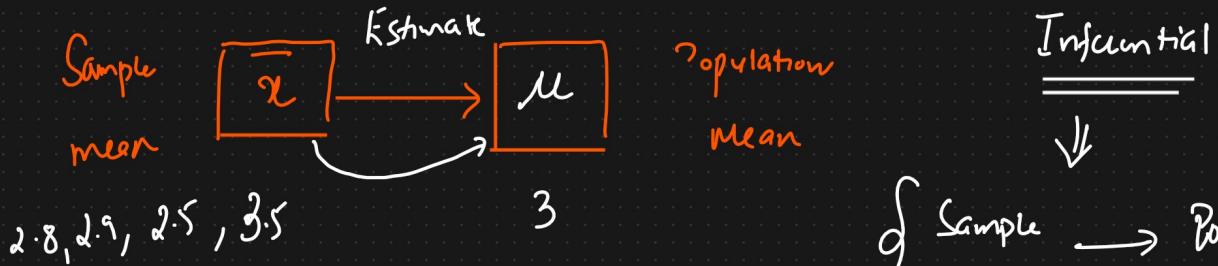


① Confidence Interval



Point Estimate

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

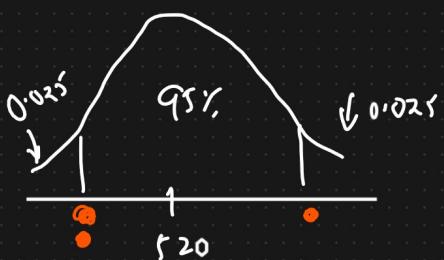


Confidence Interval

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?

$$\text{Ans) } \sigma = 100 \quad n = 25 \quad \bar{x} = 520 \quad (\cdot I = 95\%) \quad \alpha = 0.05$$



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

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

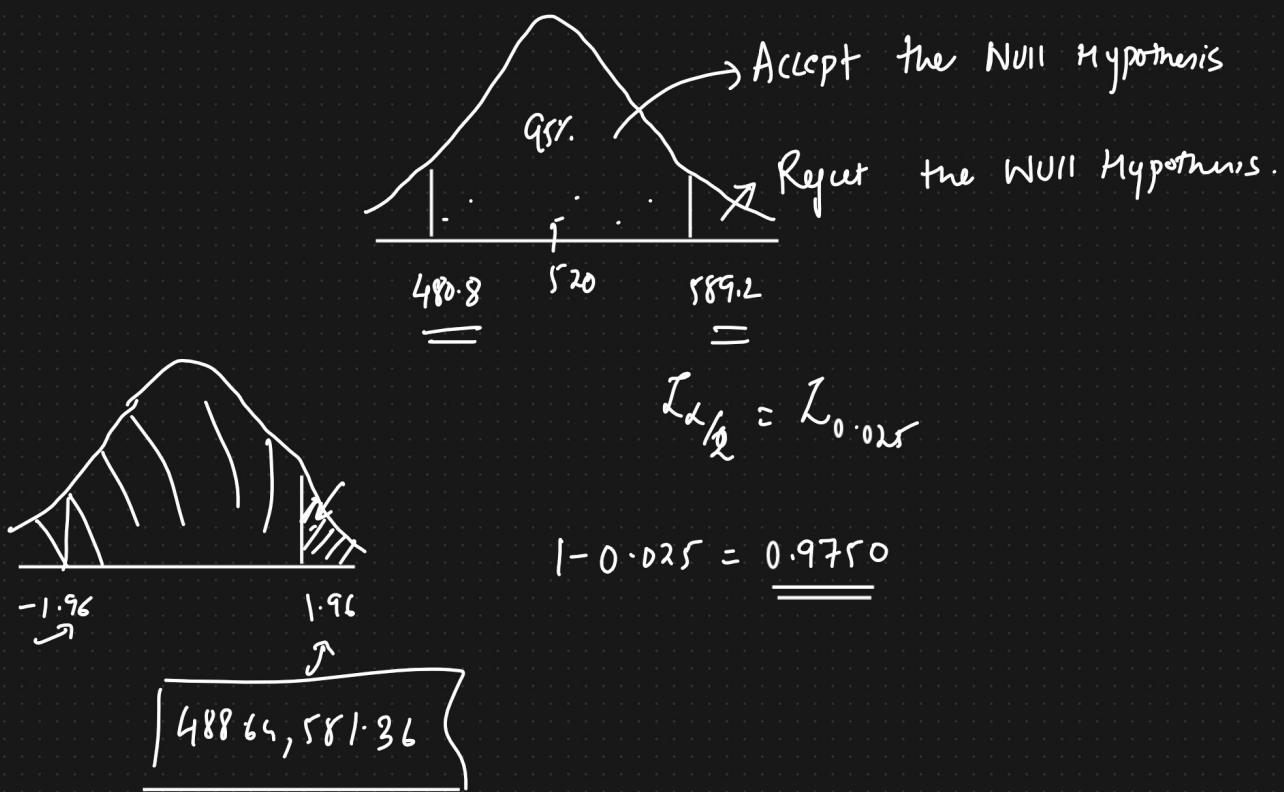
$$\bar{x} \pm Z_{\alpha/2} \left[\frac{\sigma}{\sqrt{n}} \right] \rightarrow \text{Standard Error}$$

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

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

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

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



- 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

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

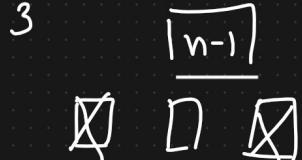
t -test $\Rightarrow t$ - table { Because population S_d is not given }

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

$$t_{0.025}$$

t -test

$$\textcircled{1} \text{ Degree of freedom} = n-1 = 25-1 = 24 \quad \underline{\underline{=}}$$



3 people

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

- (f) Type 1 and Type 2 Error.
- (g) 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 X

Outcome 3:

We accept the Null Hypothesis, \Rightarrow Type 2 Error X

When in reality it is false

Confusion Matrix

Outcome 4: We accept the Null Hypothesis

when in reality it is True



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

{ \rightarrow Stock market is going to crash }

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?

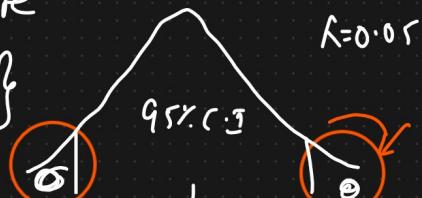
$$\alpha = 0.05 \Rightarrow 95\% \text{ C.I} \rightarrow 85\%$$

of placement rate

less than 85% }



1 tail



{ Placement rate greater than 85% }

2 tail Test

1 tail Test

Saturday

10 min probability

Sunday

① Z test Hypothesis Testing

EDA \rightarrow 3-4 projects

② J Test Hypothesis Testing

FE \rightarrow _____

③ Significance value of P value.

Machine Learning

④ ANOVA TEST

⑤ CHI SQUARE TEST

⑥ Practical

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.

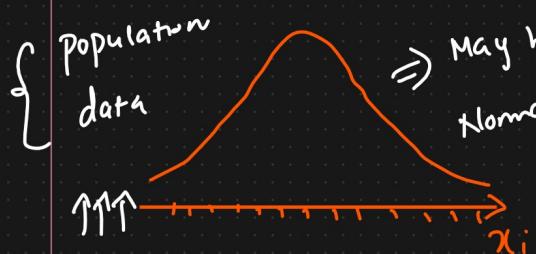
d) Chi Square (Categorical Test)

e) ANNOVA (F Test)

Influential



① Central Limit Theorem



\Rightarrow May be Gaussian
Normal Distr → $[x_1, x_2, x_3, x_4, \dots, x_{30}] \rightarrow \bar{x}_1$

$n > 30$

Sample mean distribution

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

$\rightarrow \bar{x}_3$

$\rightarrow \bar{x}_4$

\vdots

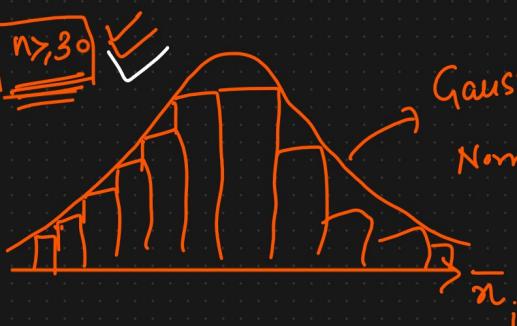
$\rightarrow \bar{x}_m$

\Rightarrow It may not

Sample m =

Sample mean

distribution

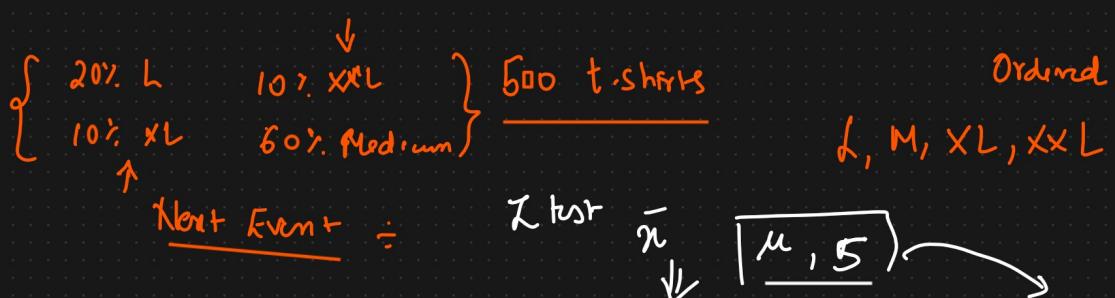


Gaussian Distribution
Normal Distribution

(2) Influential Statistics {Data Analyst, Data Scientist}

① 100K \Rightarrow T-shirt \Rightarrow No \Rightarrow Sample data \Rightarrow X_L, L, Small

② iNeuron \rightarrow Meetup \rightarrow Hitesh \Rightarrow 300-400 people \rightarrow T-shirts



③ ATM ④ Measure the size of entire sharks CI []

⑤ Amazon delivery {Percentile, Quantiles} \Rightarrow

(*) Hypothesis Testing

① A factory has a machine that fills 80ml of baby medicine in a batch. 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

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

$H_0: \mu = 80$ {Null Hypothesis}

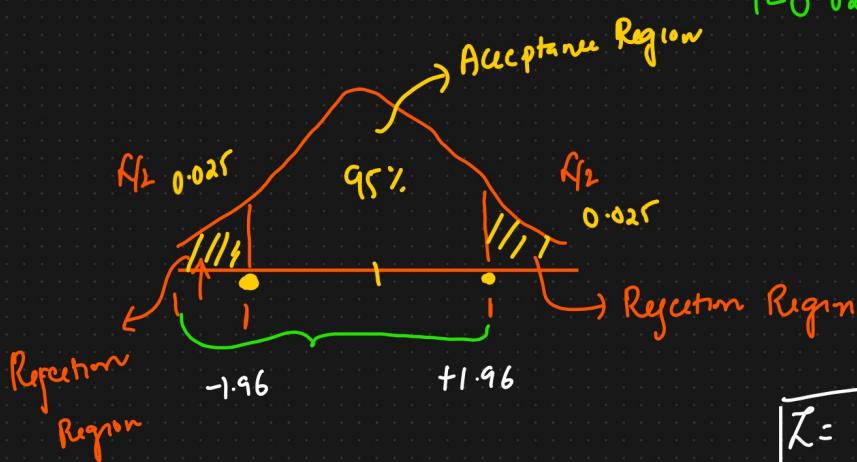
$H_1: \mu \neq 80$ {Alternate Hypothesis} Why Z test?

Step 2:

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

$n > 30$ $n \leq 30$
 (1) (2)
 population std or sample
 std

Step 3: Decision Boundary



$$1 - 0.025 = 0.9750$$

Why $+ z_{1-\alpha}$
 (1) Sample std
 (2) $n < 30$

$n=1$

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

$$Z = \frac{\bar{x} - \mu}{$$

$$\text{Sample Standard Deviation} = \frac{S / \sqrt{n}}{\sigma / \sqrt{n}} \Rightarrow \text{Standard Error}$$

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

(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.

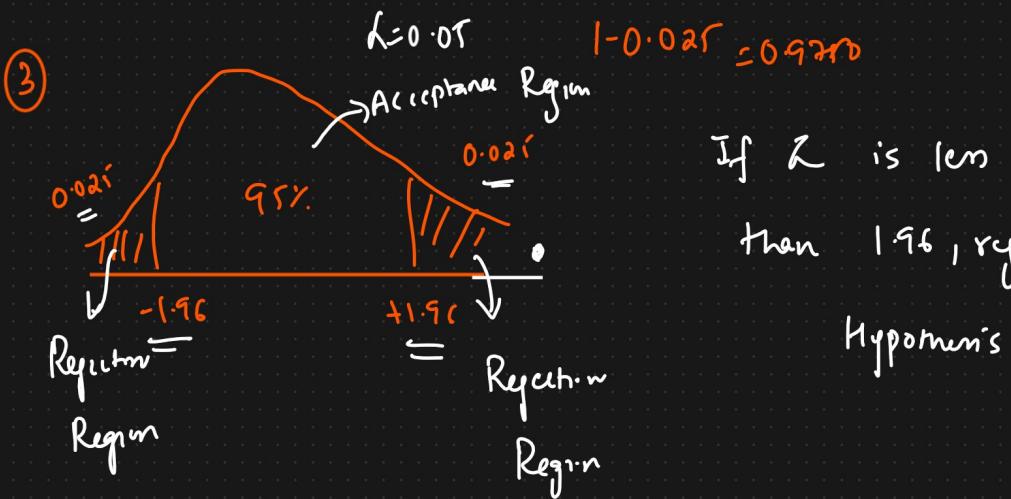
A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect Intelligence? $\left\{ \begin{array}{c} 95\% \\ \hline \downarrow \\ C.I. \end{array} \right\}$

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

① $H_0 : \mu = 100$

$H_1 : \mu \neq 100$

② $\alpha = 0.05 \quad C.I = 95\%$



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

$14.60 > 1.96$ Reject the Null Hypothesis

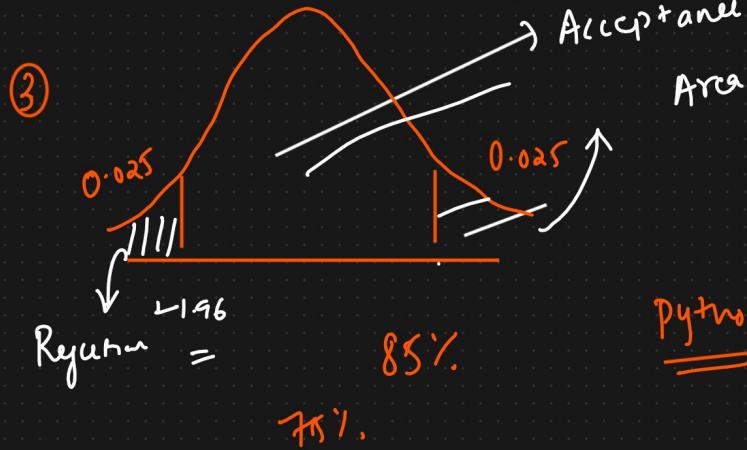
(*) A Complain was registered , the boys in the Municipal Primary School are underfed. Average weight of boys of age 10 is 32kgs with $S.D = 9\text{ kgs}$. A sample of 25 boys was selected from the municipal school and the average weight was found to be 29.5 kgs ? With $C.I = 95\%$ Check whether it is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

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

=

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

$$H_1 = \mu < 32$$



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

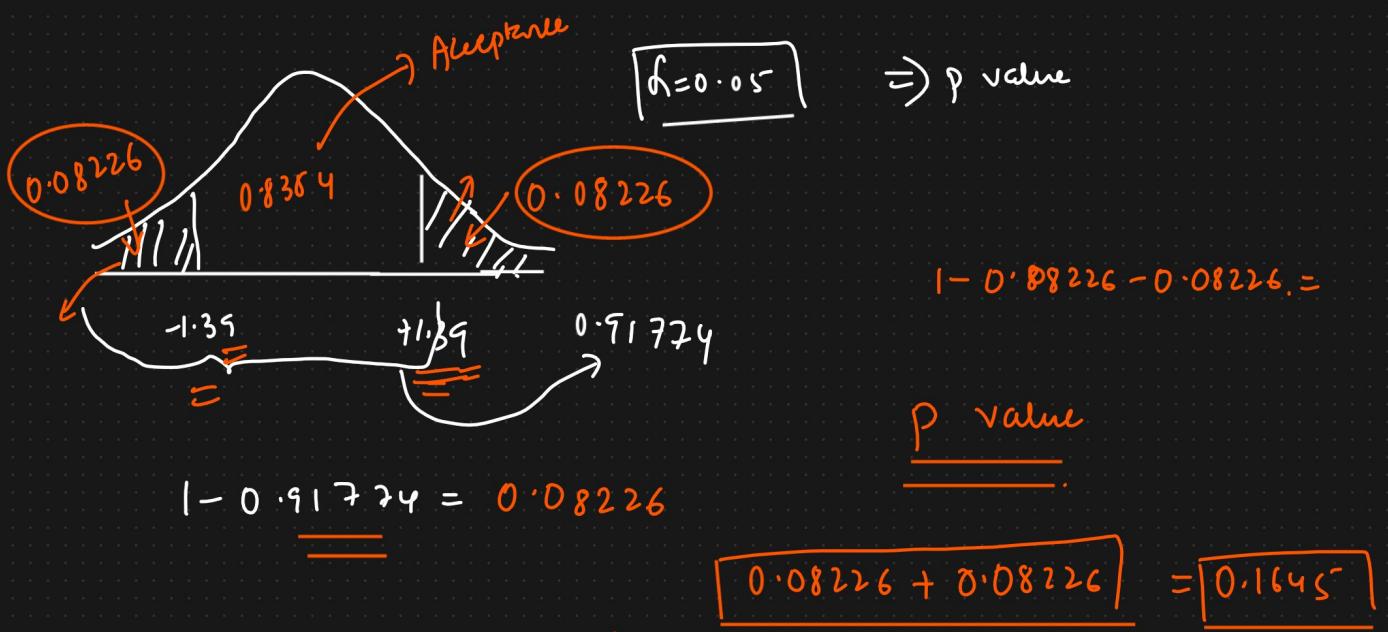
$$= \frac{29.5 - 32}{9 / \sqrt{25}} = -1.39.$$

Z-test, p-value}

Conclusion : $-1.39 > -1.92$ therefore we accept the Null Hypothesis

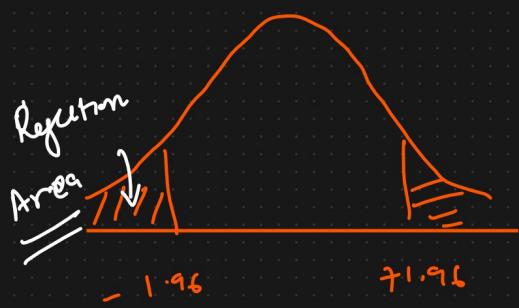
So, the boys are not underfed.

So, the boys are not underfed.



Significance value
 $0.1645 > 0.05$

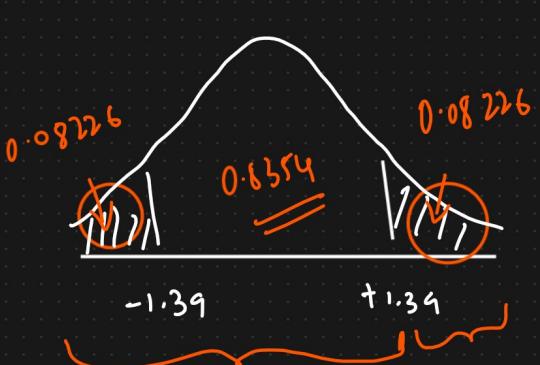
\downarrow
 $p \geq 0.05 \Rightarrow \text{Accept the Null Hypothesis}$



$Z \text{ test}$

\downarrow

$p \text{ value}$



$\Rightarrow p \text{ value} = 0.08226 + 0.08226$
 $= 0.16$

$1 - 0.08226 - 0.08226$

Domain

\downarrow

$0.1645 > \text{Significance}$

\downarrow
 value

$1 - 0.91774 = 0.08226$

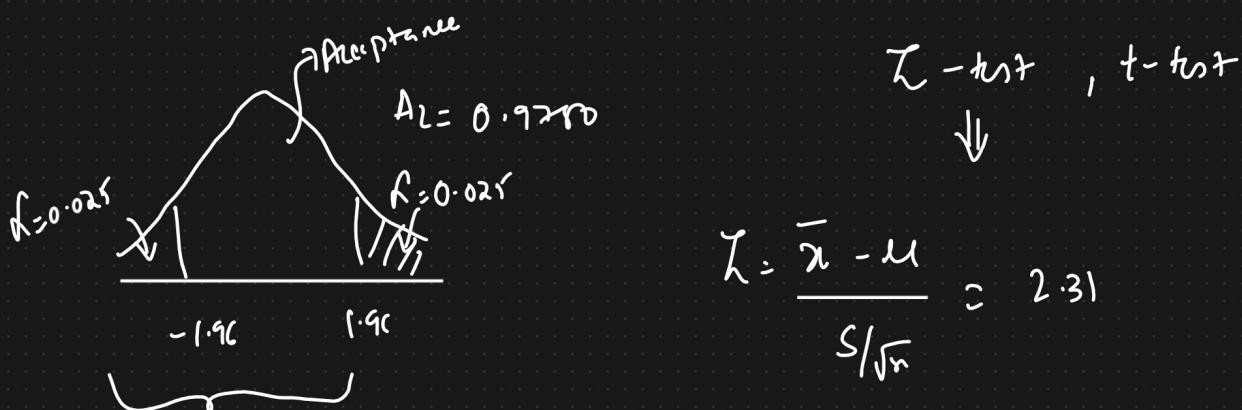
$\Rightarrow \text{Accept the Null Hypothesis.}$

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??

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

$$H_1 : \mu \neq 168 \quad \underline{\quad} \quad c = 0.95 \quad \alpha = 1 - c \cdot I = 0.05$$



$2.31 > 1.96$ Reject the Null Hypothesis

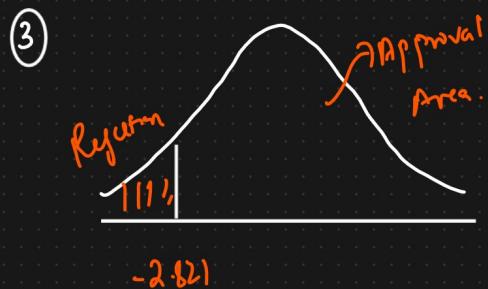
⑤ A company manufactures bike batteries with an average life span of 2 or more years. An engineer believes this value to be low. Using 10 samples, he measures the average life span to be 1.8 years with a standard deviation of 0.15.

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
 $H_1 : \mu < 2$ ≤ 3.0 $t-tst??$ Std is
 $\{$ given }

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



Degrees of freedom: $n-1$

$= 9$

④ Calculate t-test Statistic:

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

⑤ Conclusion

$-4.216 < -2.821$ Reject the Null Hypothesis. }
 \Downarrow

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) $H_0: p_0 = 0.70.$

$H_1: p_0 \neq 0.70$

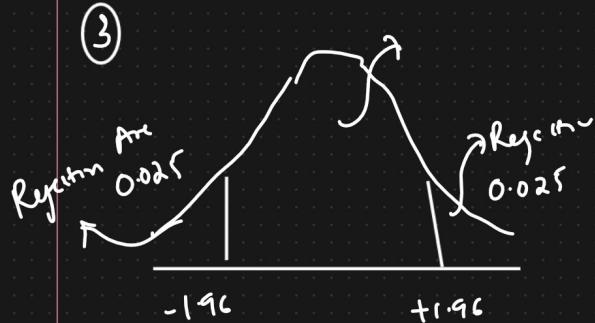
$$n = 200 \quad X = 130 \\ \hat{P} = \frac{X}{n} = \frac{130}{200} = \frac{13}{20} = 0.65$$

$$q_0 = 1 - p_0$$

② $\alpha = 0.05 \quad C.I = 95\%$

$$Z_{test} = \frac{\hat{P} - P_0}{\sqrt{\frac{p_0 q_0}{n}}}$$

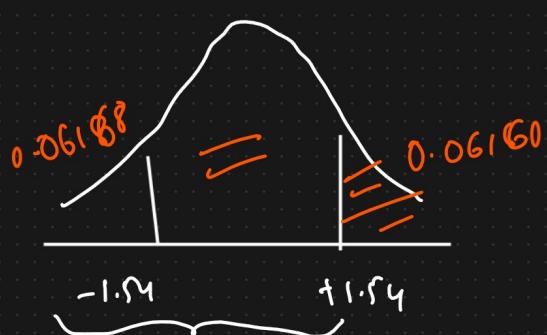
$$= \frac{0.65 - 0.70}{\sqrt{\frac{0.7 \times 0.3}{200}}} \approx -1.54$$



At 95% C.I there is

$-1.54 > -1.96$, So we accept

the Null Hypothesis



$$1 - 0.93822 = 0.06168$$

p-value
 \downarrow
 $2 \times 0.06168 > 0.05$

Accept Null Hypothesis

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?

$$p\text{ value} = .014$$

Statistics

{ 11:30 - 12pm }

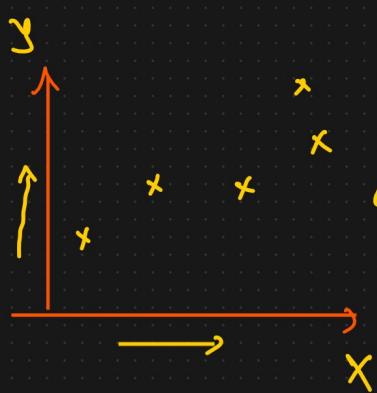
- ① Covariance
- ② Pearson Correlation Coefficient
- ③ Spearman Rank Correlation Coefficient
- ④ CHI SQUARE TEST
- ⑤ ANNOVA (F-Test)

✓ practicals
✓

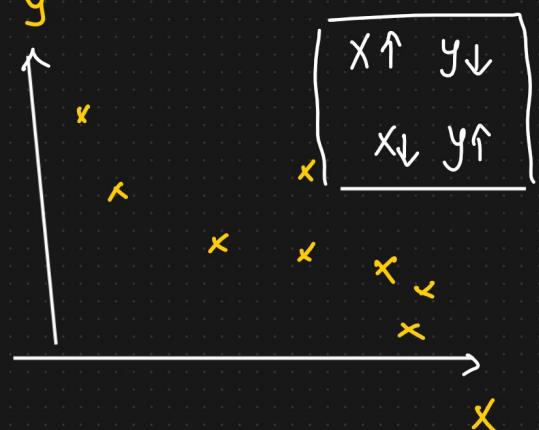
Covariance

$X \uparrow$	$y \uparrow$	-	-
$X \uparrow$	$y \downarrow$	-	-
$X \downarrow$	$y \uparrow$	-	-
$X \downarrow$	$y \downarrow$	-	-

↓ ↓ { quantity the relationship
X = Y = between X & Y }



$\left\{ \begin{array}{l} X \uparrow \quad Y \uparrow \\ X \downarrow \quad Y \downarrow \end{array} \right.$



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

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

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

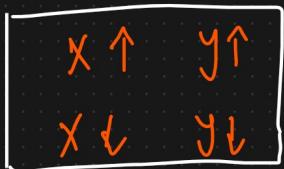
$$= \frac{\sum (x_i - \bar{x}) \times (x_i - \bar{x})}{N-1}$$

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

↓

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

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



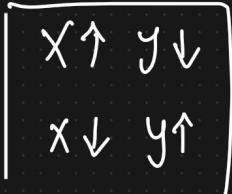
$$\text{Cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\left\{ \begin{array}{l} = (2-4)(3-5) + (4-4)(5-5) \\ \quad + (6-4)(7-5) \end{array} \right.$$

x	y
2	3
4	5
6	7

2

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



\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

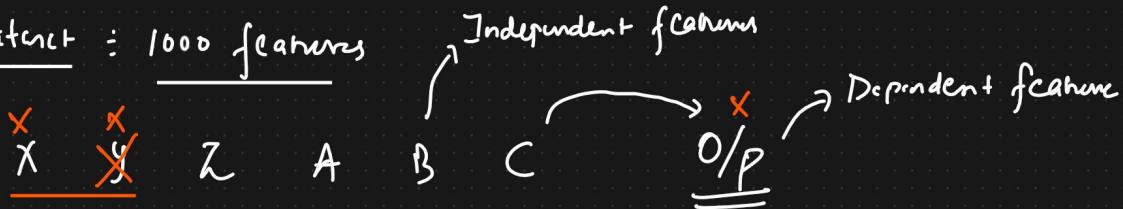
OP

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

$$x, y \Rightarrow 99\% \quad \underline{=}$$

$$\underline{90\%} \quad \underline{0.9}$$

-ve Correlation

↓
Keep it

③ Spearman Rank Correlation

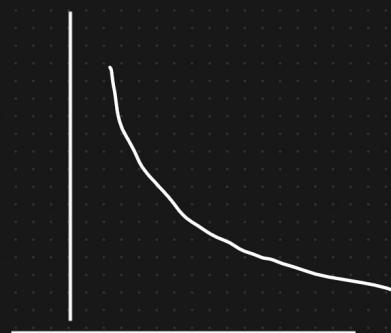
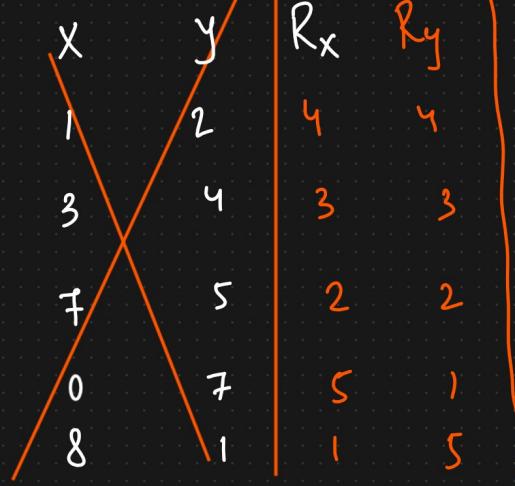
$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sqrt{R(x)} \sqrt{R(y)}}$$

Marks



Spearman Rank

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



$$\underline{-1}$$

(f) 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.

- f) 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)

Expected	<18	18-35	>35
20%	30%	50%	$95\% \text{ C.I}$

$n=500$

Observed : 121 288 91

Expected 100 150 250

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

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 χ^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.

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

Practicals + EDA + Feature Engineering }

Statistics

- ① ANOVA (F-Test) → 1 hour }
 ② FDD → { 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)
Medicine { Dosage } Anxiety reducing { Gender }

0mg	50mg	100mg
\bar{x}	\bar{y}	\bar{z}

factor : Dosage

9	6	3
8	6	2
7	7	3
8	8	3
8	8	3

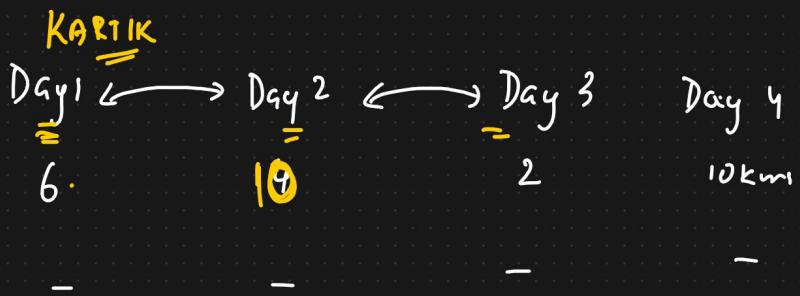


Types of ANOVA

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 levels	<u>Running Kms</u>
1	6.
2	-

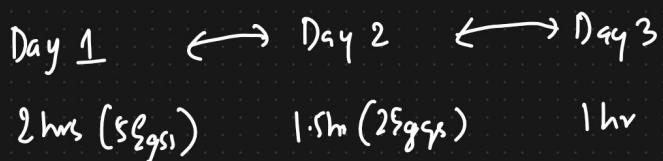


Ques. Study hours of KARTIK

④ Factorial ANOVA



Gym



⑥ Factorial ANOVA : 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
Men	9	7	4	
	8	6	3	
Women	7	5	2	
	8	7	3	
	8	8	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 $H_0 = \mu_{0\text{mg}} = \mu_{50\text{mg}} = \mu_{100\text{mg}}$ }
 $H_1 = \text{not all } \mu's \text{ are equal}$ }

② State α and C.I

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

③ Calculate the Degree of freedom

Statistics

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

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

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

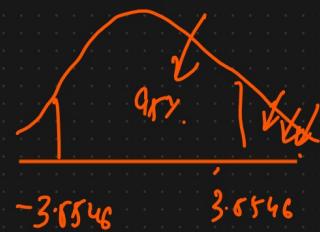
$$\rightarrow df_{\text{Within}} = N - a = 21 - 3 = 18$$

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

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		

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

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

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

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

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

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

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

Final Conclusion

Accept

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



$$\begin{aligned} H_0 : \mu &= \text{Some value} \\ H_1 : \mu &\neq \text{Some value} \end{aligned} \rightarrow \underline{95\% \text{ CI}}$$

Virginia

=

=

Pctz1 width

-

-

-

-

-

-

-

-

-

-

-

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

$H_1 = \cdot \neq \text{p-value.} \neq \text{reject the Null Hypothesis}$

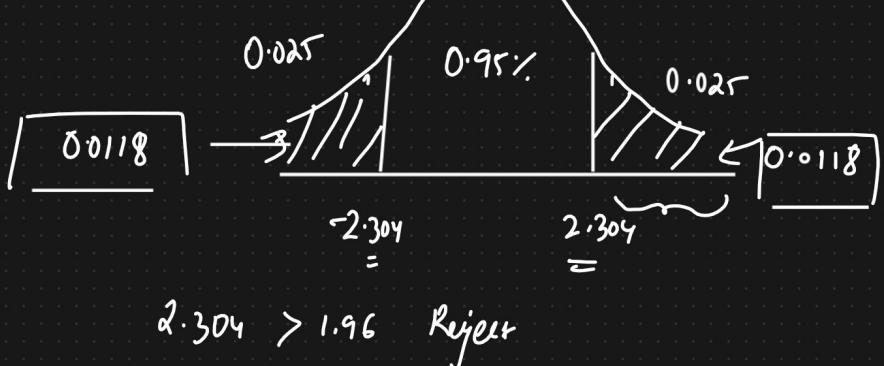
$$0.0118 \quad 0.0228 < 0.05 \quad 1 - 0.025 = 0.975$$

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

Z_{test}

$d =$

$Z_{\text{test statistic}}$



$$Z = 2.304$$

$2.304 > 1.96 \text{ Reject}$