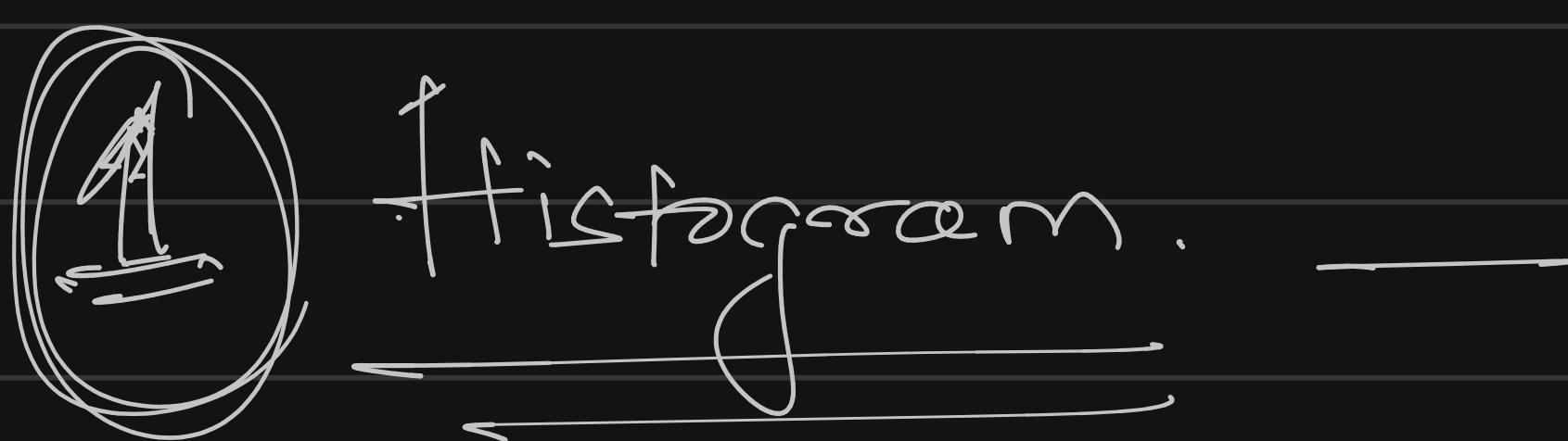


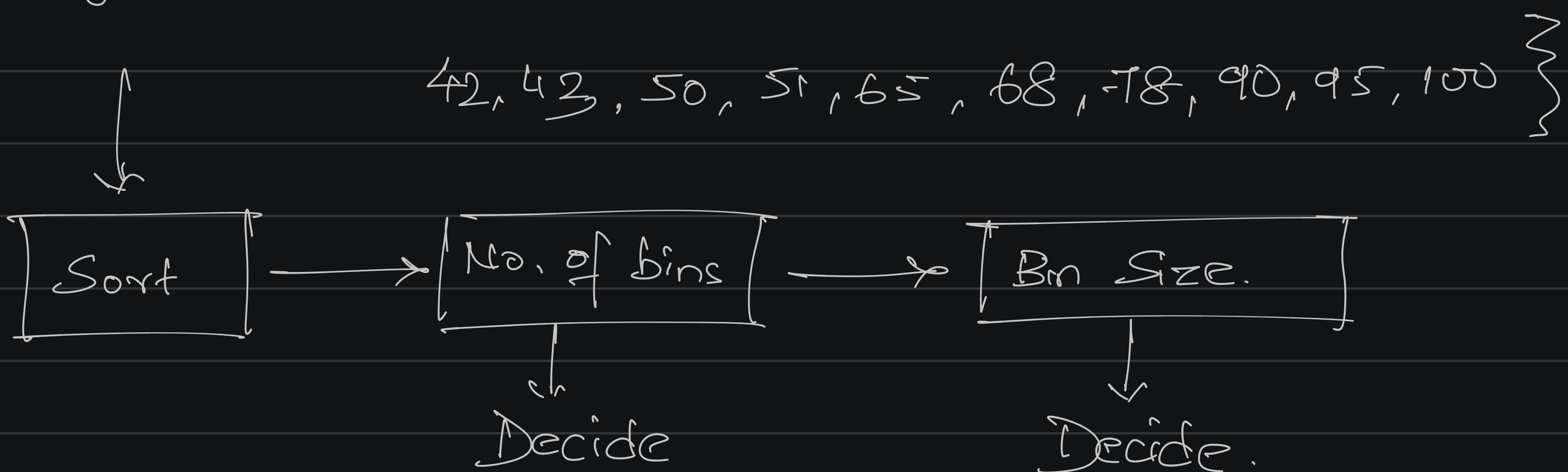
~~Day 2.~~

Agenda.

- (1) Histogram .
- (2.) Measure of central tendency .
- (3.) Measure of dispersion .
- (4.) Percentiles & Quartiles .
- (5.) 5 Number Summary (Box Plot) .



$$\text{Agl.} = \{ 10, 12, 14, 18, 24, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100 \}$$



$$\text{e.g. } \min = 10 .$$

$$\max = 40$$

$$\text{bins} = 10 .$$

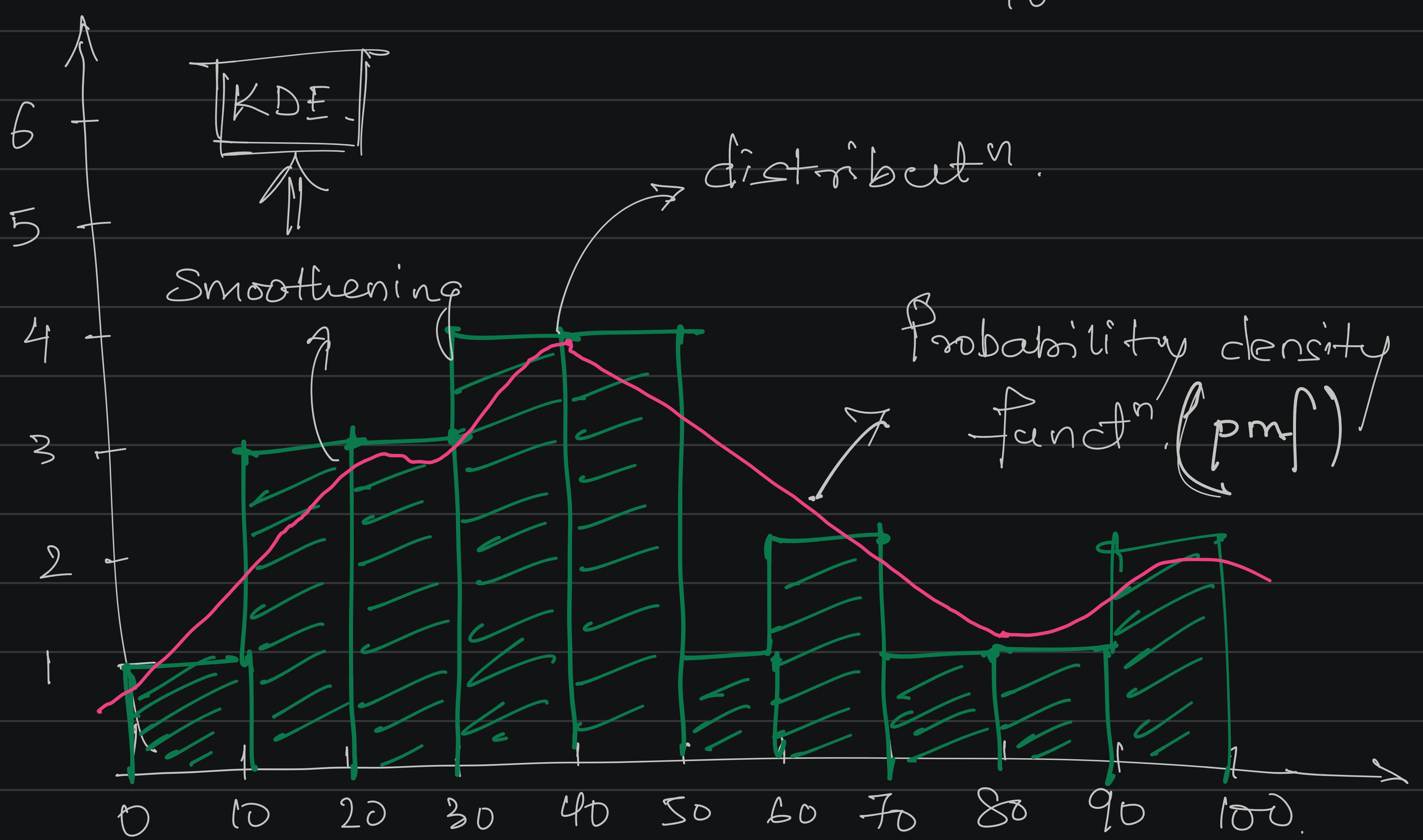
$$\text{binsize} = \frac{40}{10} .$$

$$\left[10, 20, 25, 30, 35, 40 \right]$$



For ages data, histogram :

$$\text{Bins} = 10, \quad \text{Bin size} = \frac{100}{10} = 10,$$



* Weight. = $\{ 30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 77, 80, 90, 95 \}$.

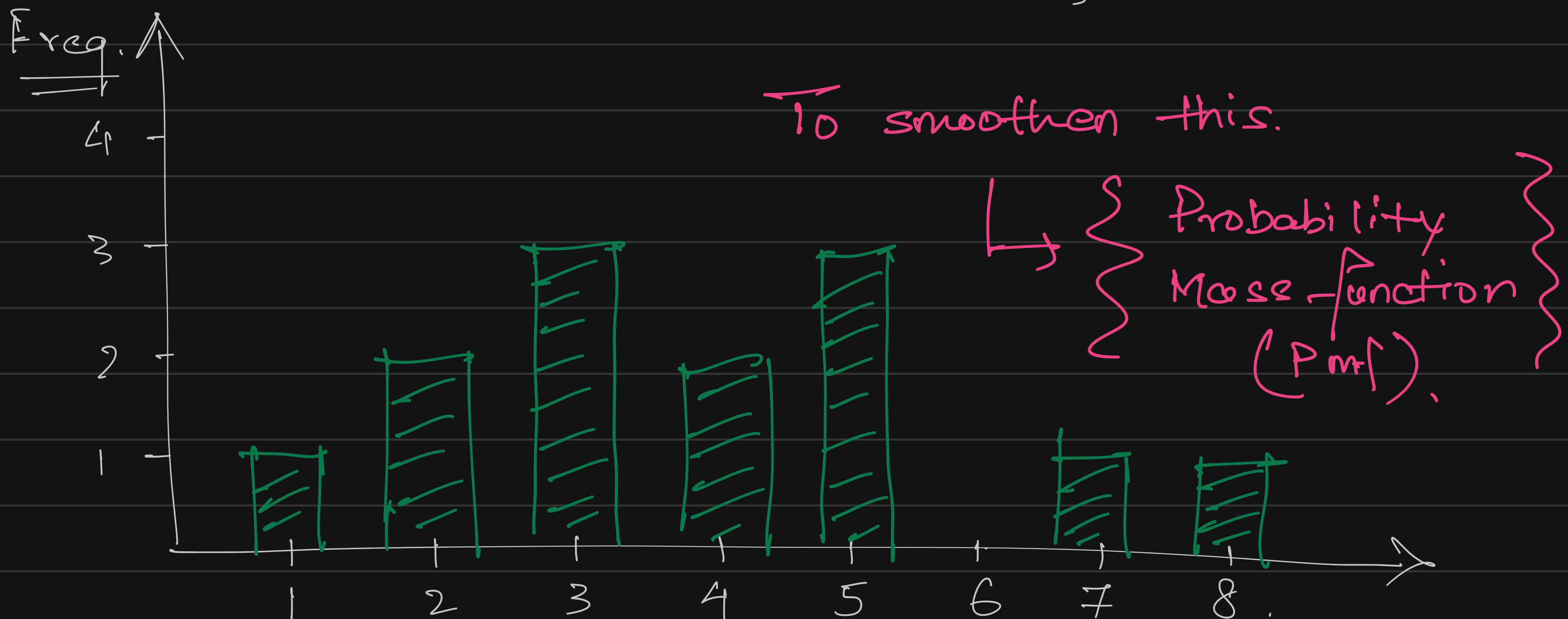
Bins = 10.

$$\text{Bin Size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5.$$

\Rightarrow Continuous variables.

* for discrete continuous variables.

No. of bank. accounts. = $\{ 2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5 \}$



2. Measure of central tendency

i) Mean.

ii) Median.

iii) Mode

A measure of central tendency is a single value that attempts to describe a set of data identifying the central position.

i) Mean. $\Rightarrow \bar{x} = \{1, 2, 3, 4, 5\}$

$$\text{Mean / Avg} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

Population (N)

Sample (n)

Mean.

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

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

Population age = $\{24, 23, 2, 1, 28, 27\}$

$$N = 6, \mu = \frac{24+23+2+1+28+27}{6}$$

$$\mu = 17.5$$

Sample = $\{24, 2, 1, 27\}$.

Sample mean (\bar{x}) = $\frac{24 + 2 + 1 + 27}{4}$.

$$\bar{x} = 13.5$$

Here, $N > n$.

For means,

$$\mu \geq \bar{x}$$

or $\bar{x} \geq \mu$.

Practical Application

Age

Salary

Family Size.

—
—
—
—
NAN

—
—
—
—
NAN

—
—
—
—
NAN

→ loss of info.
↓
Replace NAN with mean

NAN



Median.

Age

Salary.

24

45

Mean Age

$$= \underline{\underline{29.6}}$$



28

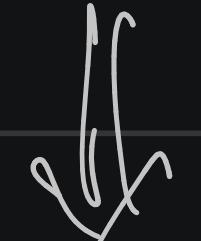
50

Mean salary

29

NAN.

$$= \underline{\underline{62}}$$



Can replace NAN.

31

60

Can replace.
NAN.

36

80

NAN.

NAN.

Additional =
data.

$\left[\begin{array}{c} 80 \\ 200 \end{array} \right]$

outliers.

New mean,

$$= 38.$$

New mean

$$= 85.$$

Causing deviatⁿ.

This is why

Median is used.

Example.

Outlier.

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

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

$$\bar{x} = 3.$$

$$\bar{x} = 19.16.$$

deviation due to outlier.

Steps to find out median

1. Sort the numbers

2. Find the central number.

(i) No. of elements are even:

We find avg. of central elements.

(ii) No. of elements are odd:

We find the central element.

Solved.

=====

Central elements.

$$\{1, 2, 3, 4, \boxed{5, 6}, 7, 8, 9, 10\} \rightarrow \text{Even. no.}$$

$$\text{Median} = \frac{5+6}{2} = 5.5.$$

of
elements.

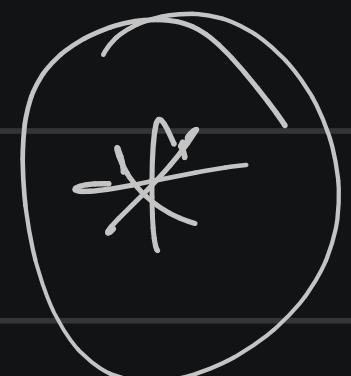
Sorted .

$\{ 0, 1, 2, 3, 4, \boxed{5}, 6, 7, 8, 100, 120 \}$

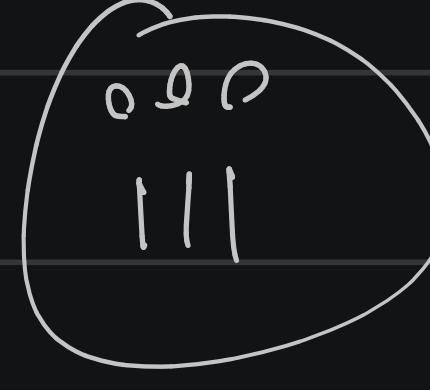
↓
Central element
↓
Median .

Odd no.
of elements

$$\text{Mean} = \underline{\underline{25.6}}.$$



When there are outliers \rightarrow Use median.

( Mode \rightarrow Most frequent occurring element .

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

Mode = 3 .

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

Mode = 2, 3

Practical Importance

e.g. Dataset -

Types of flower.

Lily

Sunflower

Rose

[NAN]

Sunflower

Rose.

[NAN.]

Rose.

Mode = Rose.

Replace.

→ Mostly used
for categorical
variables.

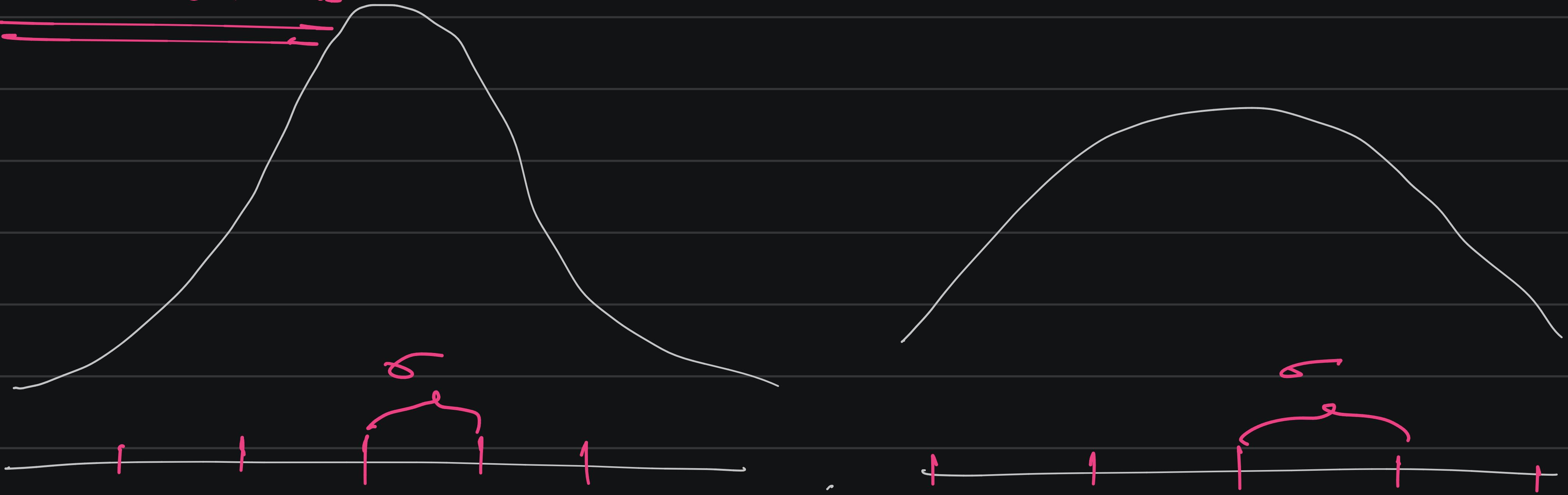
3. Measure of dispersion

(i) Variance (σ^2)

(ii) Standard deviation (σ)

(i) Variance \rightarrow Shows the spread of data.

Distributions



Variance

< Variance.

SD

< SD.

4

Percentiles & Quartiles.

Percentage = { 1, 2, 3, 4, 5, 6, 7, 8 }

Percentage (Even nos.) = No. of even nos.

Total no. of nos.

$$= \frac{4}{8} . = 0.5 = 50\%$$

Percentile → Percentile is a value below
 which a certain percentage
 of observations lie.

$$Q_1 = \frac{25}{100} * (n+1) = \frac{25}{100} * 21$$

≈ 5.25 . \rightarrow av. of 5th & 6th index.

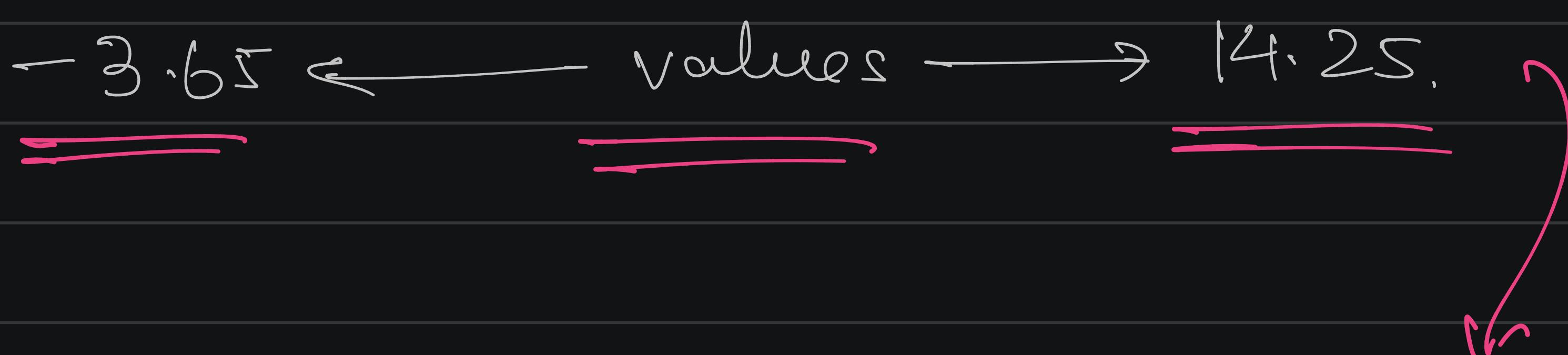
$$Q_1 = 3.$$

$$Q_3 = \frac{75}{100} * 21 = 15.75. \rightarrow \text{av. of } 15^{\text{th}} \& 16^{\text{th}} \text{ index.}$$

$$Q_3 = 15.5.$$

$$\text{Lower Fence} = 3 - 1.5(4.5) = -3.65.$$

$$\text{Higher Fence} = 14.25.$$



So, 27 is the outlier.

5 number Summary

1 Minimum = 1

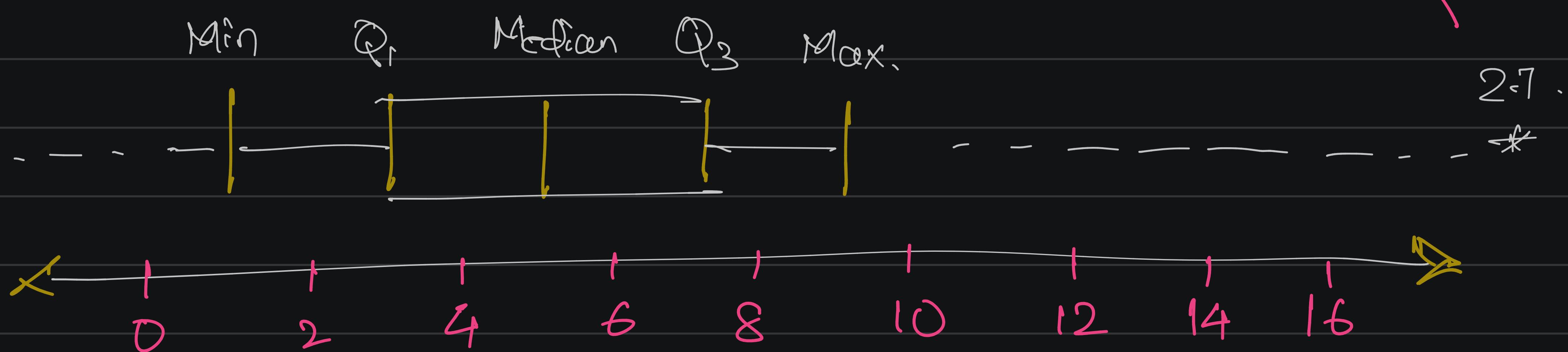
2 $Q_1 = 3.$

3 Median = 5.

4 Maximum = 7.5.

5 Maximum = 9.

Outlier.



Box Plot

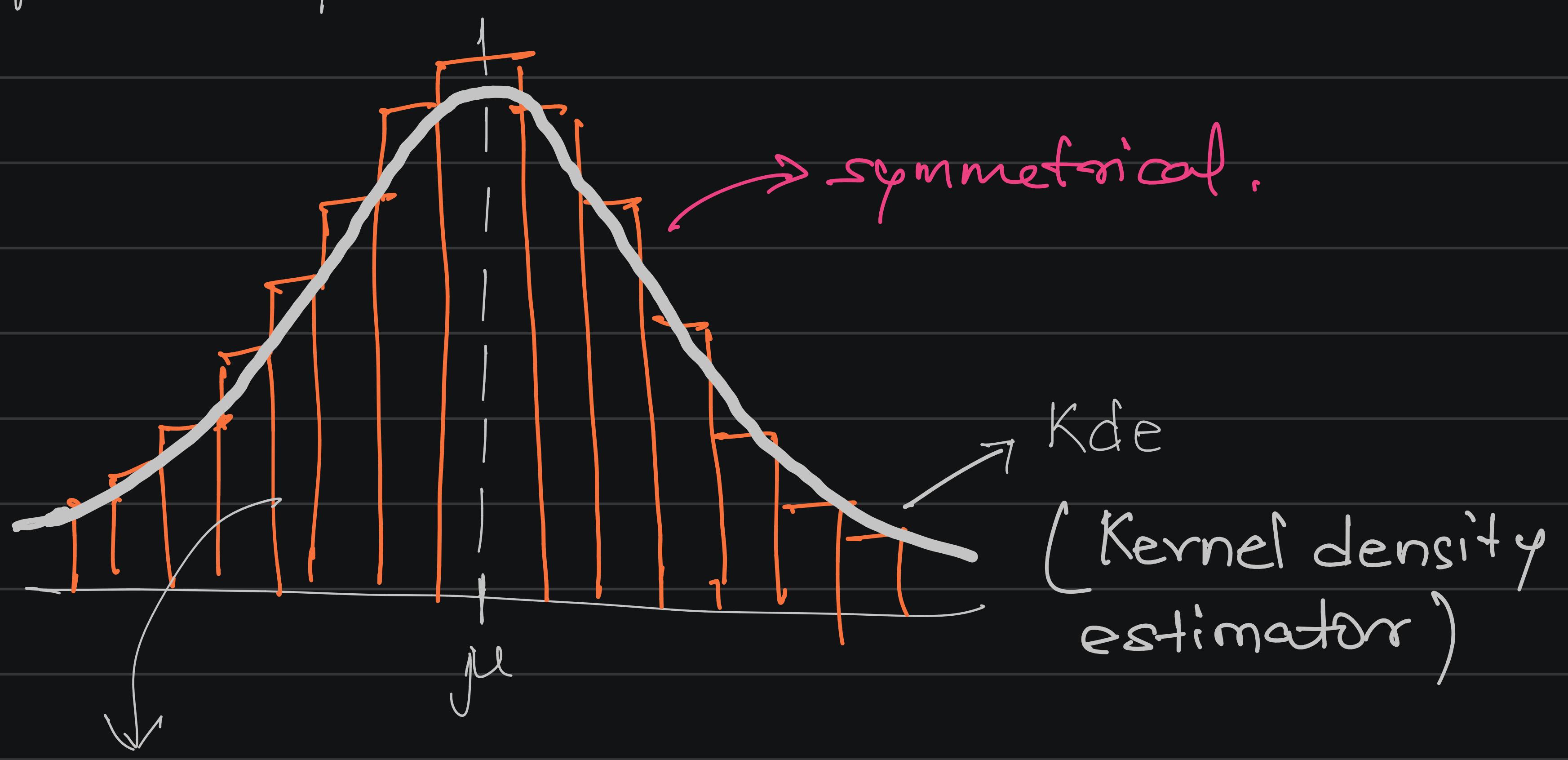
Day 5

Agenda.

- (1) Normal distribution .
- (2) Standard Normal distribution -
- (3) Z - Score
- (4.) Standardization & Normalization



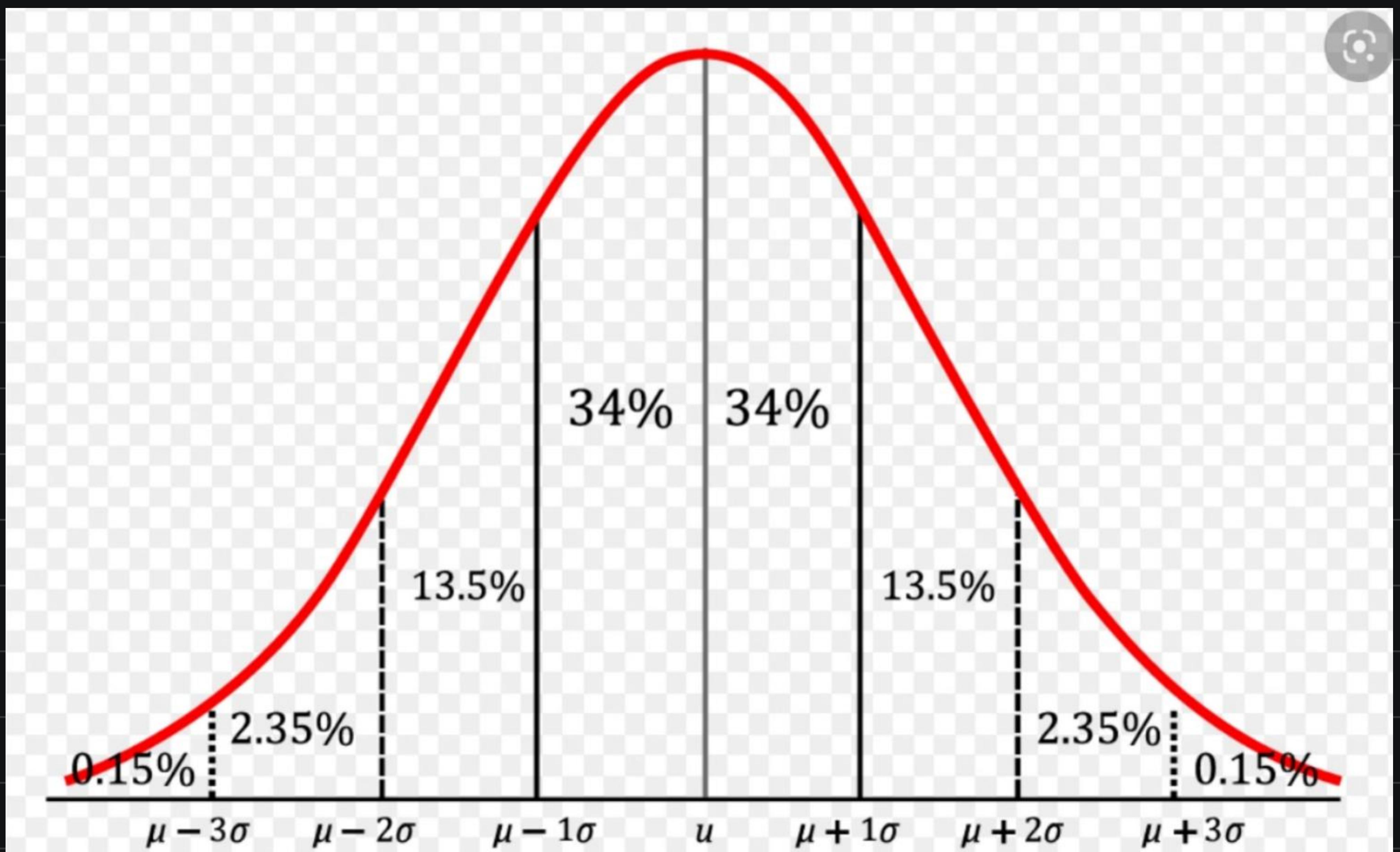
Gaussian / Normal distribution .



$$\text{Area under curve} = 1 (100\%)$$

e.g. Age, Weight, Height . \rightarrow follow
Normal distribution .

(*) Empirical Rule of Normal distribution.



→ Assumptions

(i) 68% of data falls within 1 SD (Standard deviation)

(ii) 95% of data falls within 2 SD.

(iii) 99.7% of data falls within 3 SD

68 - 95 - 99.7% rule

→ Empirical Rule.

~~Q-Q~~

Q-Q plot → Can show if a distribution
is Gaussian or not.

(2.)

Standard Normal Distribution.

if X belongs to Gaussian Distribution (μ, σ)

↓ can be converted into. using Z Score

Y \sim Standard Normal distribution
 $(\mu = 0, \sigma = 1)$

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

$$\mu = 3$$

$$\sigma = 1.41$$

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma / \sqrt{n}}$$

$n = 1$, applied on every element

⇒ Standard Error.

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

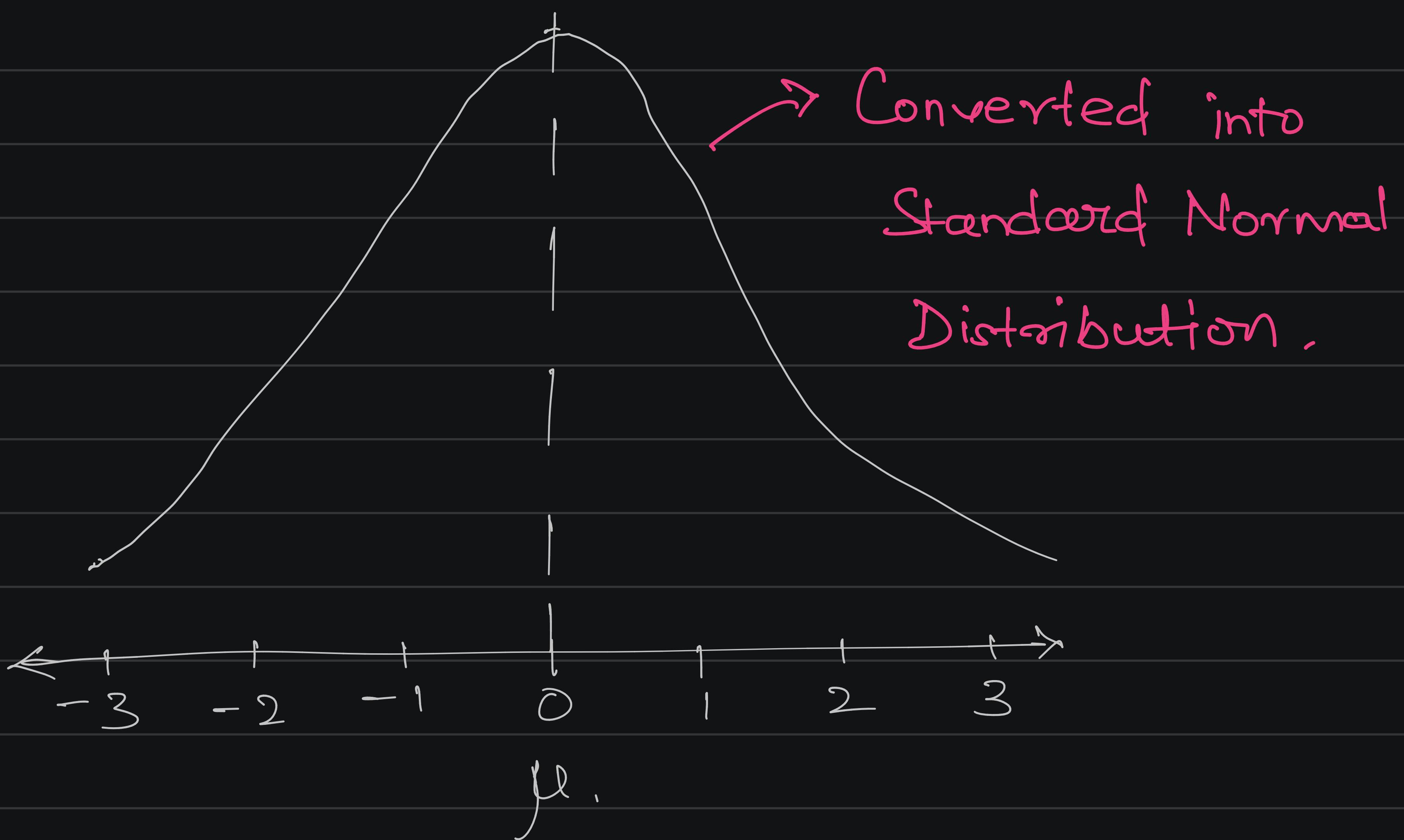
$\textcircled{*}$ $y = \{-1.414, -0.707, 0, 0.707, 1.414\}$

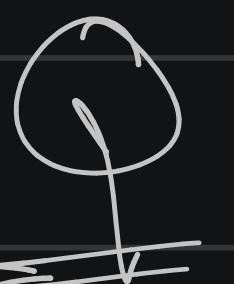
for 1st element $Z\text{-score} = \frac{1-3}{1.414} = -1.414$

for 2nd element $Z\text{-score} = \frac{2-3}{1.414} = -0.707$

for 3rd element $Z\text{-score} = \frac{3-3}{1.414} = 0$

Distribution of y .



 Why do we need to convert ?

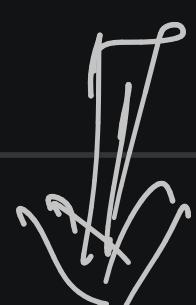
3 features →	Age (year)	height (kg)	height (cm)
	24	72	150
	26	78	160
	32	84	165
	33	92	170
	34	87	150
	28	83	180

29.

80.

175.

Units are different



Value ranges differ a lot.



Not favorable of ML models.



* Apply Z-score on all Features to bring them all to the same scale (-3 to 3)

$$\mathcal{X} \mu = 0 \quad \mathcal{X} \sigma = 1$$

Feature Scaling

① Standardization :- { z-score applied }

$$\mu = 0, \sigma = 1$$

$$\text{data} = \{-3 \text{ to } +3\} \quad 99.7\% \text{ data}$$

2 In Normalization → We give the range
within which we need to
convert the values.

e.g. (1) Min-max scaler → transform values between
0 to 1.

Formula,

$$x_{\text{scaled}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

x

1

2

3

4

5

y.

0

0.25

0.5

0.75

1

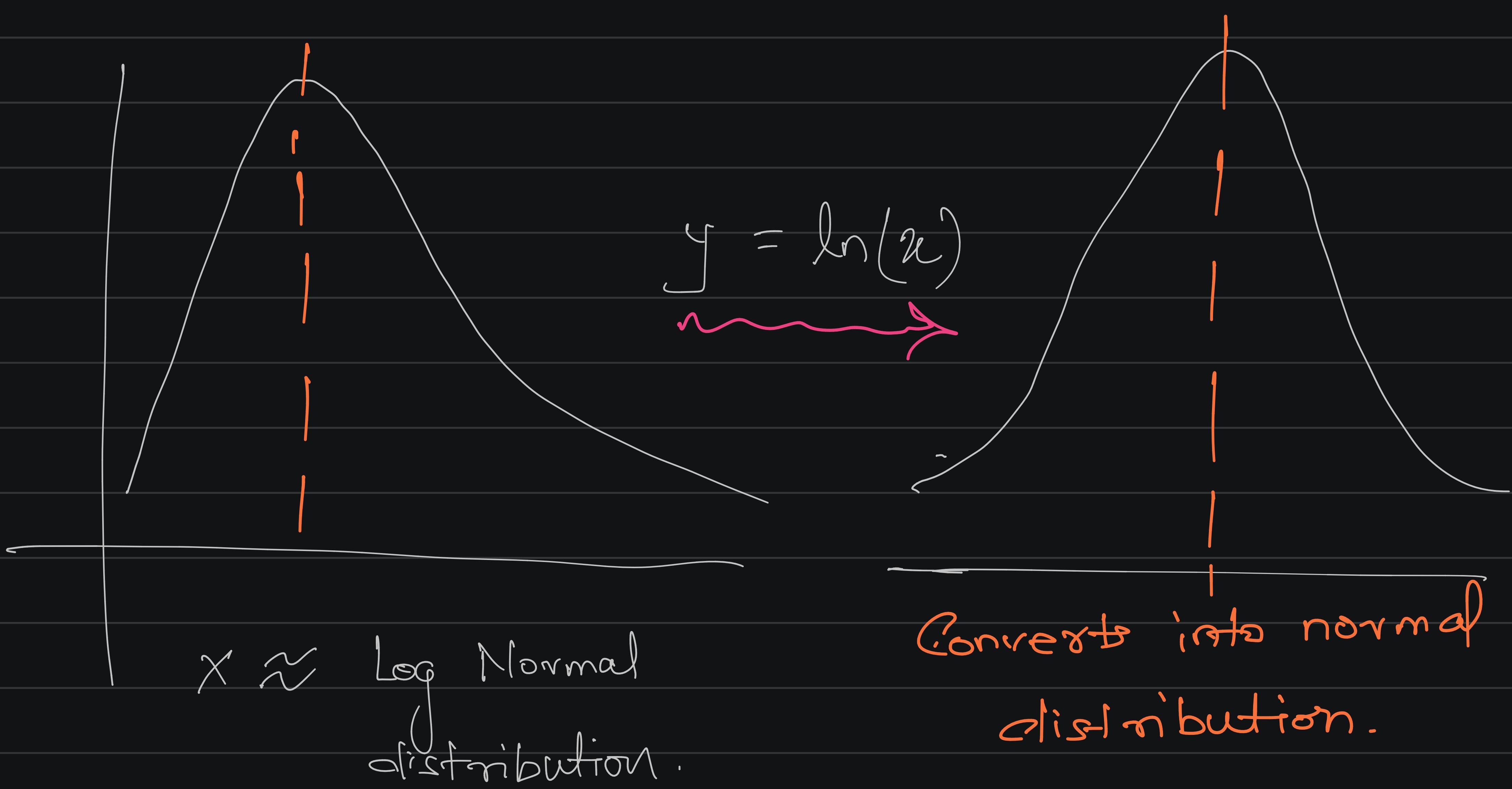
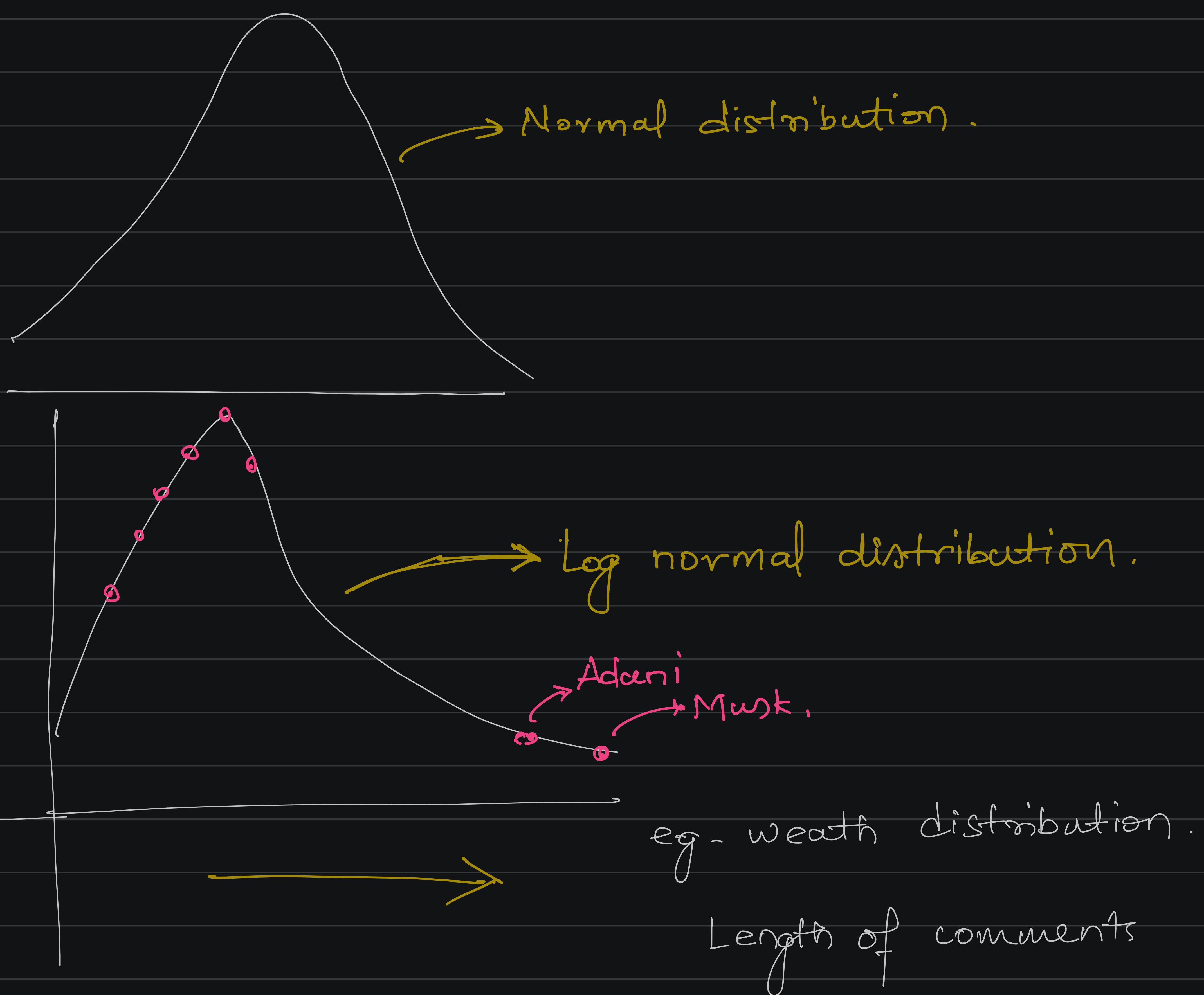
where,
 $x_{\max} = 1$

$x_{\min} = 0$

All values
converted
to 0 to 1.

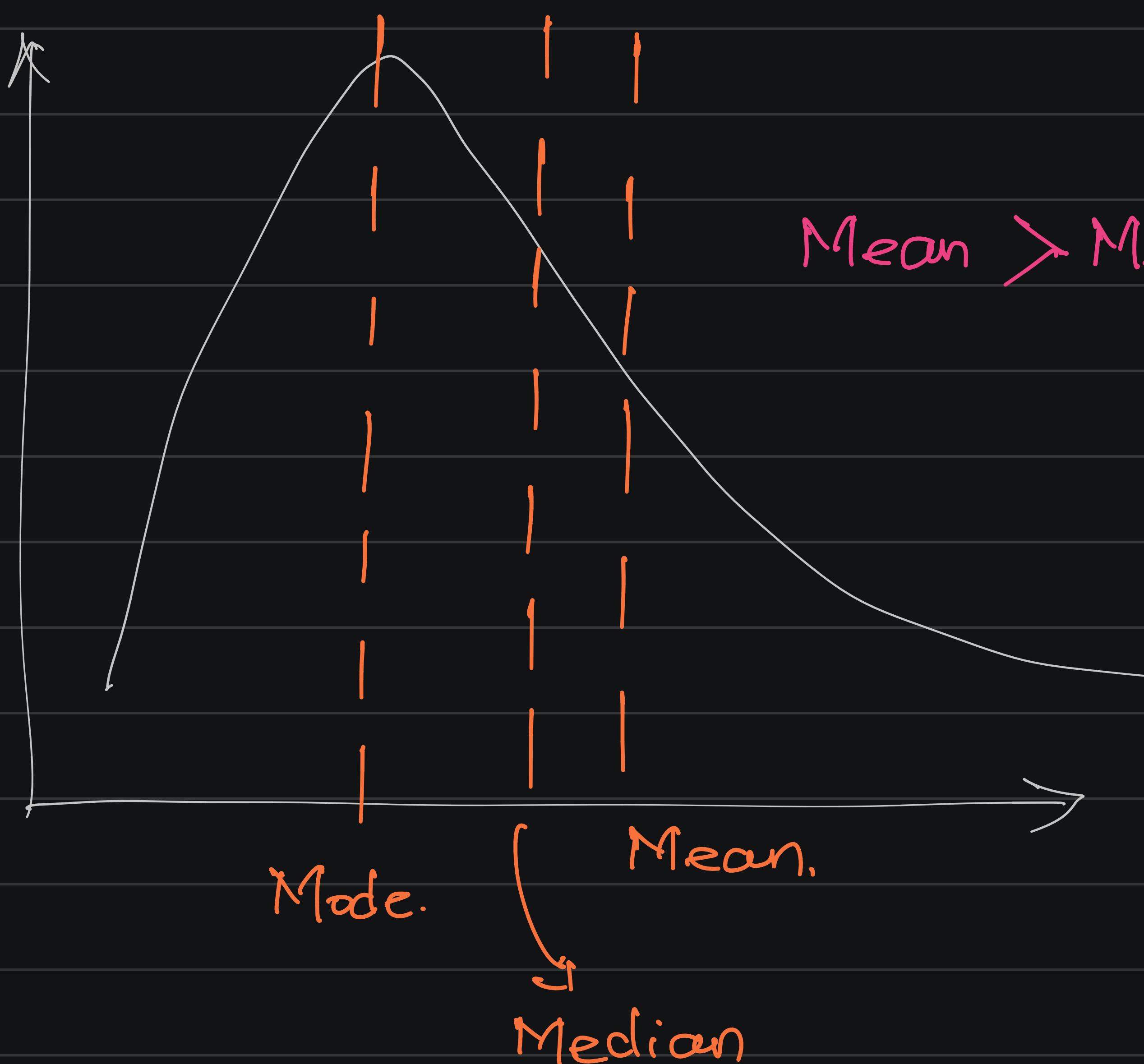
Applied in Deep Learning

Log Normal Distribution.



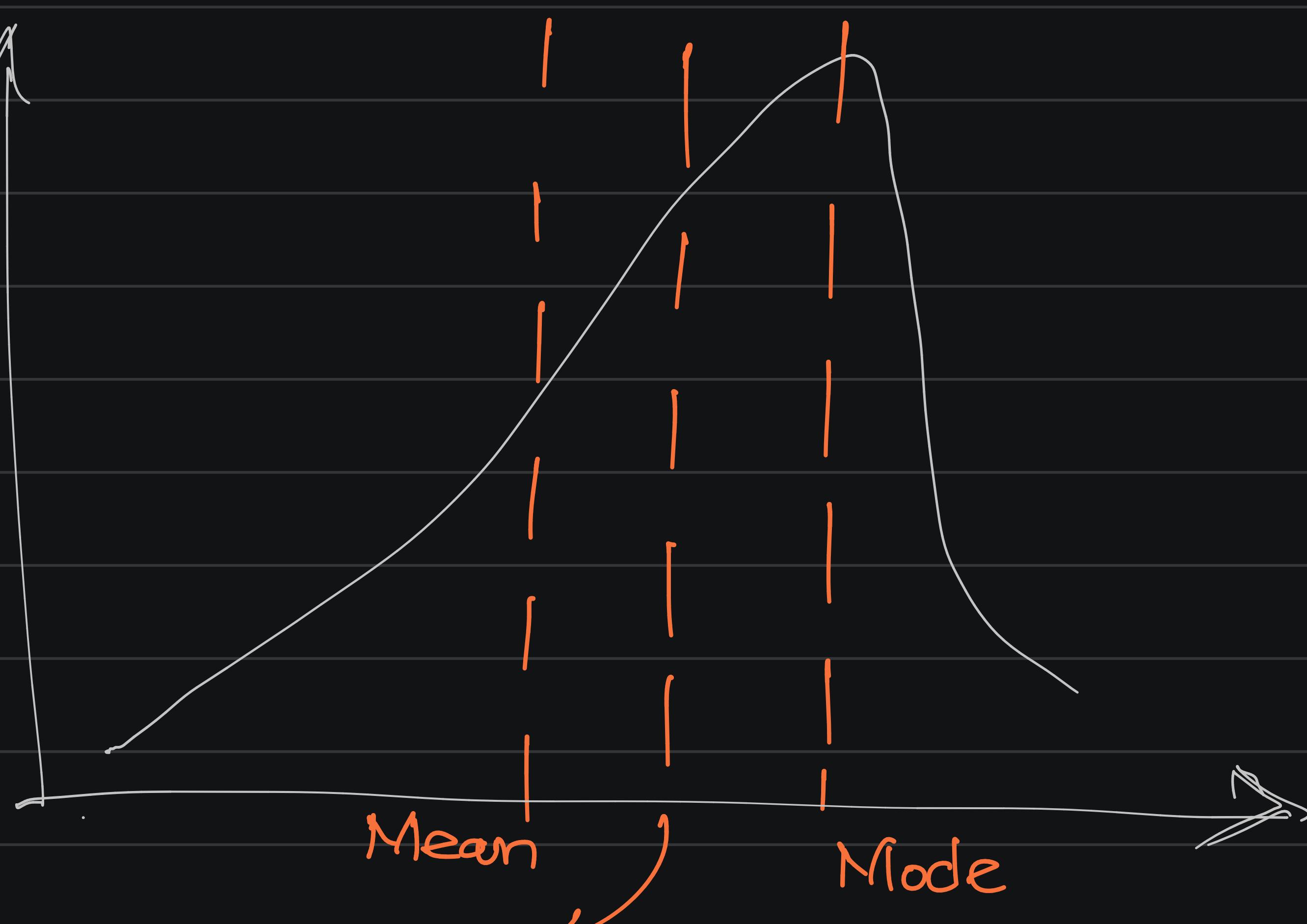
Q = What is the relationship between mean, median and mode for below distributions

(1)



Mean > Median > Mode.

(2)



Mean < Median < Mode.

(2.) Use Z-table to find area under curve.

z	0	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
+0	.50000	.50399	.50798	.51197	.51595	.51994	.52392	.52790	.53188	.53586
+0.1	.53983	.54380	.54776	.55172	.55567	.55966	.56360	.56749	.57142	.57535
+0.2	.57926	.58317	.58706	.59095	.59483	.59871	.60257	.60642	.61026	.61409
+0.3	.61791	.62172	.62552	.62930	.63307	.63683	.64058	.64431	.64803	.65173
+0.4	.65542	.65910	.66276	.66640	.67003	.67364	.67724	.68082	.68439	.68793
+0.5	.69146	.69497	.69847	.70194	.70540	.70884	.71226	.71566	.71904	.72240
+0.6	.72575	.72907	.73237	.73565	.73891	.74215	.74537	.74857	.75175	.75490
+0.7	.75804	.76115	.76424	.76730	.77035	.77337	.77637	.77935	.78230	.78524
+0.8	.78814	.79103	.79389	.79673	.79955	.80234	.80511	.80785	.81057	.81327
+0.9	.81594	.81859	.82121	.82381	.82639	.82894	.83147	.83398	.83646	.83891
+1	.84134	.84375	.84614	.84849	.85083	.85314	.85543	.85769	.85993	.86214
+1.1	.86433	.86650	.86864	.87076	.87286	.87493	.87698	.87900	.88100	.88298
+1.2	.88493	.88686	.88877	.89065	.89251	.89435	.89617	.89796	.89973	.90147
+1.3	.90320	.90490	.90658	.90824	.90988	.91149	.91308	.91466	.91621	.91774
+1.4	.91924	.92073	.92220	.92364	.92507	.92647	.92785	.92922	.93056	.93189

0.59 \Rightarrow 59% → Area under curve for head.

$$\text{Area under curve} = (1 - 0.59) \\ (\text{tail}) = 41\%$$

Day 4.

Agenda :

- (1) Central Limit theorem .
- (2) Probability .
- (3) Permutation & Combination .
- (4) Covariance , Pearson Correlation ,
Spearman Rank Correlation .

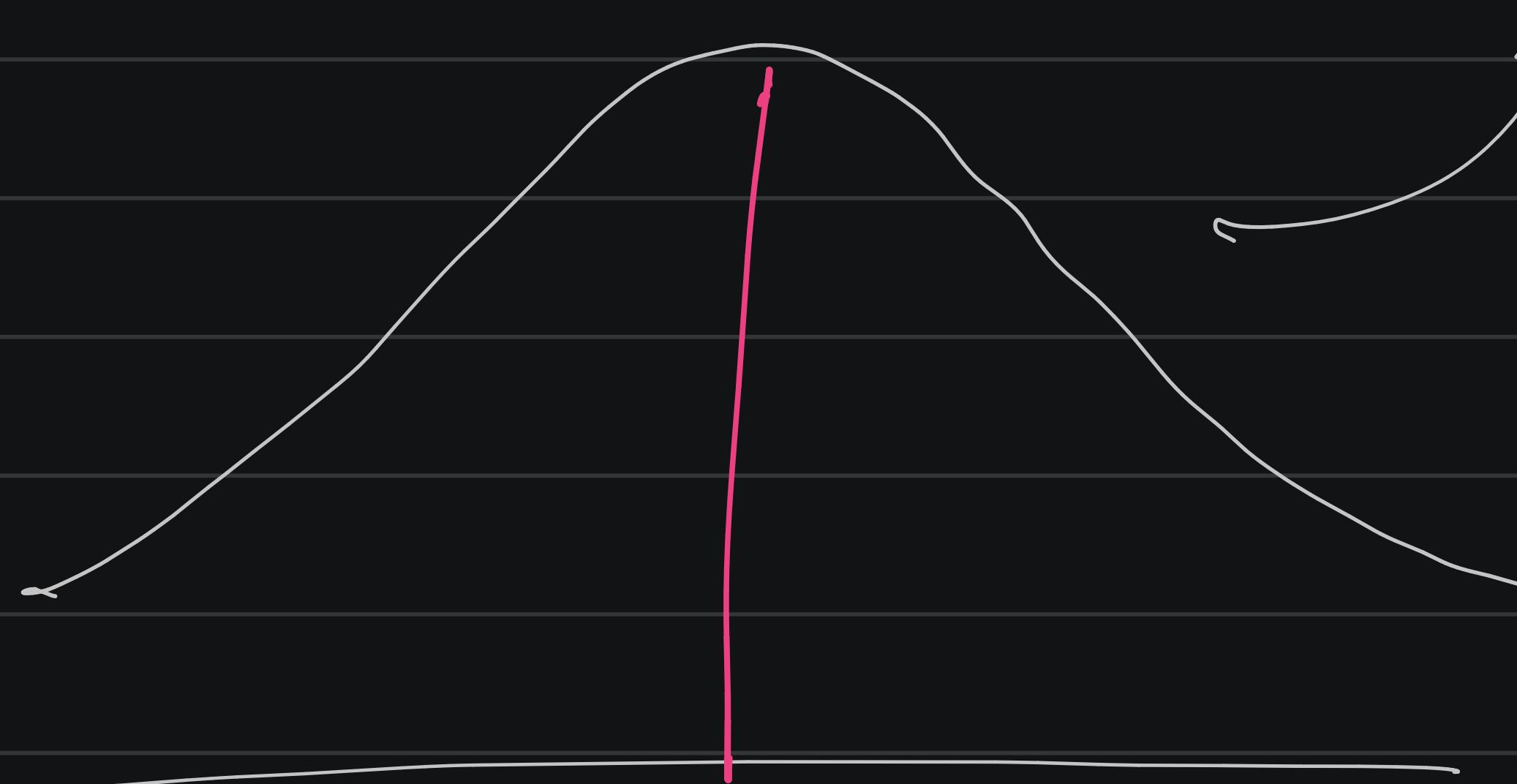
- { (5) Benoulli's distribution .
- (6) Binomial distribution .
- (7) Power Law. (Pareto distribution)

→ Not covered on Day 4.

A. Central Limit theorem.

Or.

The means will follow a gaussian distribution :-



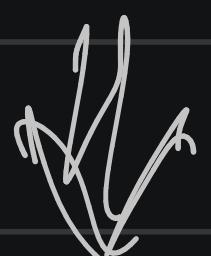
As the distribution is now a gaussian distribution

All the assumptions about gaussian distribution
is valid now.

2. Probability — Probability is a measure

of the likelihood of an event

e.g. Tossing of coin (fair) $P(H) = 0.5$ → Head.



$$P(T) = 0.5$$

→ Tail.

SCHOOL

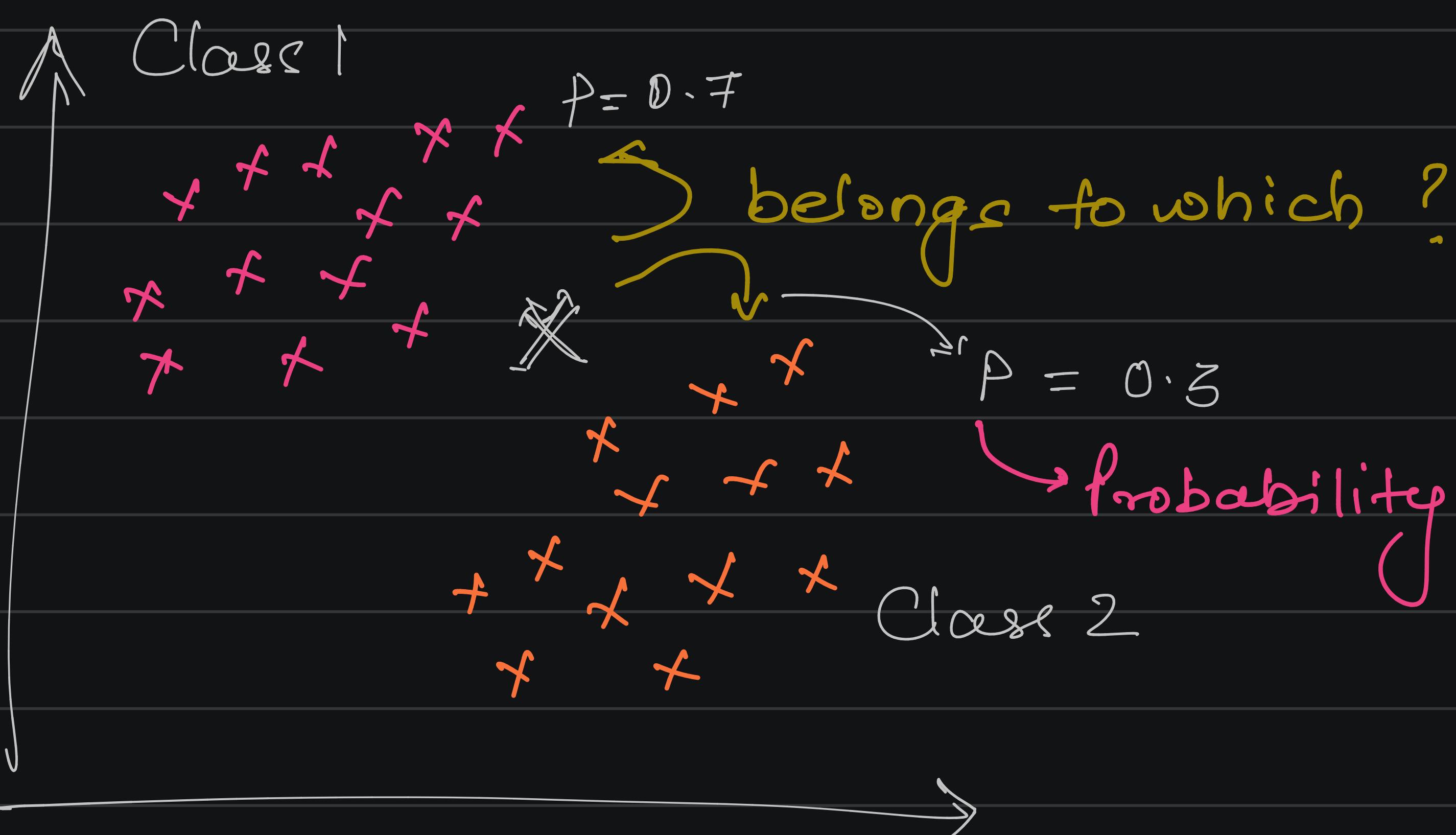


$$COIN \rightarrow P(H) = 1.$$

↓
Unfair coin.

$$\text{Rolling a dice} \quad P(1-\delta) = \frac{1}{\delta}$$

Use in ML : Classification Problem.



Mutual Exclusive event .

Two events are mutually exclusive if they cannot occur at the same time .

e.g. Rolling of dice , Tossing a coin .



Non-mutually Exclusive event .

Two events can occur at the same time .

e.g. Picking randomly a card from a deck of cards , two events "heart" and "king" can be selected .

Mutually exclusive

Q: What is the probability of coin landing on heads or tails ?



Addition rule for mutual exclusive events ,

$$P(H) = 0.5.$$

$$P(T) = 0.5.$$

$$\left[P(H \text{ or } T) = P(H) + P(T) \right]$$
$$= 0.5 + 0.5 = 1.$$

Q: What is the probability of getting 1 or 3 or 6 in a dice roll ?

$$P(1 - 6) = \frac{1}{6}$$

$$P(1 \text{ or } 3 \text{ or } 6) = P(1) + P(3) + P(6)$$

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6}$$

$$= \underbrace{\frac{3}{6}}_{\frac{1}{2}} = \frac{1}{2}$$

Non-mutual Exclusive event

Problem 1. = Bag of marbles : 10 Red, 6 Green
3 (R & G)

When picking randomly from a bag of marbles what is the probability of choosing a marble that is red or green?

Addition rule for non-mutual exclusive

Ans = Event 1.

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

$$P(R \text{ or } G) = \frac{13}{19} + \frac{9}{19} - \frac{3}{19}$$

$$= \frac{19}{19} = 1.$$

Independent Event

e.g. Q What is the probability of rolling a 5 & then a 3 with a normal six sided dice?

$$\text{Ans.} = P(1 \text{ to } 6) = \frac{1}{6}$$

Multiplication rule for independent events.

$$P(A \text{ and } B) = P(A) * P(B)$$

$$\Rightarrow \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

Dependent event

e.g. Q What is probability of drawing a "orange" and then drawing a "yellow" marble from the bag?

$$P(O) = \frac{4}{7}$$

Diagram showing 7 marbles: 4 orange (O) and 3 yellow (Y).

$$P(Y|O) = \frac{3}{6} = \frac{1}{2}$$

Conditional probability of yellow given that orange has happened.

$$P(O \text{ and } Y) = P(O) * P(Y/O)$$

$$= \frac{4}{7} \times \frac{3}{6} = \frac{2}{7}$$

3.

(i) Permutation.

e.g.

Chocolates.

{ dairy milk, kit kat, milky bar . }
 { sneakers, 5-star. }

Arrange these
in the available space.

5 4 3

Possibilities.

Total no. of ways = $5 * 4 * 3$.

= 60 ways.



With permutation, order matters.

Formula.:

n = total no. of objects.

r = # of selection

$$P_r = \frac{n!}{(n-r)!}$$

(from previous example) $P_r = \frac{5!}{(5-3)!} = \frac{5 \times 4 \times 3 \times 2!}{2!}$

$= 60.$

(ii) Combination.

Repetition will not occur. (Only unique combinations allowed)

Formula.:

$$C_r = \frac{n!}{r!(n-r)!}$$

(from previous example.)

$$C = \frac{5!}{3! \times 2!} = \frac{5 \times 4 \times 3!}{3! \times 2 \times 1}$$

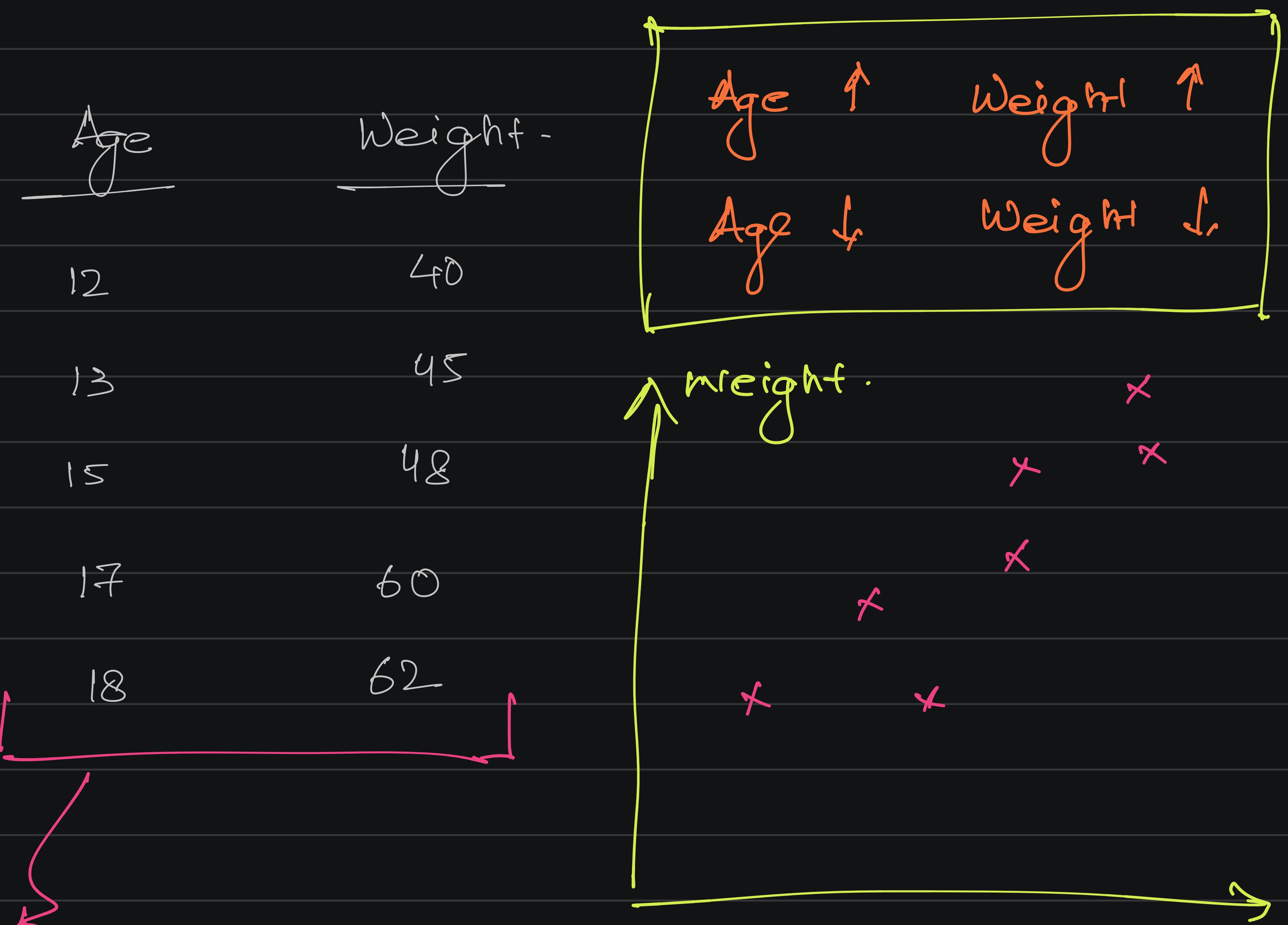
$$= 10.$$

Used in Dream 11. A lot.

4

(i) Covariance

→ important for
feature selection.



Let's quantify the relationship $x \leftrightarrow y$ Age.

using mathematical questions.

Formula:

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x}) * \sum (y_i - \bar{y})}{n - 1}$$

Note. :

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

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

So,

$$\text{Cov}(x, x) = \sigma^2(x)$$

Interview question.

For previous example,

$$\bar{x} = 15 \quad \bar{y} = 51$$

$$\text{Cov}(x, y) = \underline{\underline{24.}}$$

Positive value.

Age \propto weight

If Negative value,

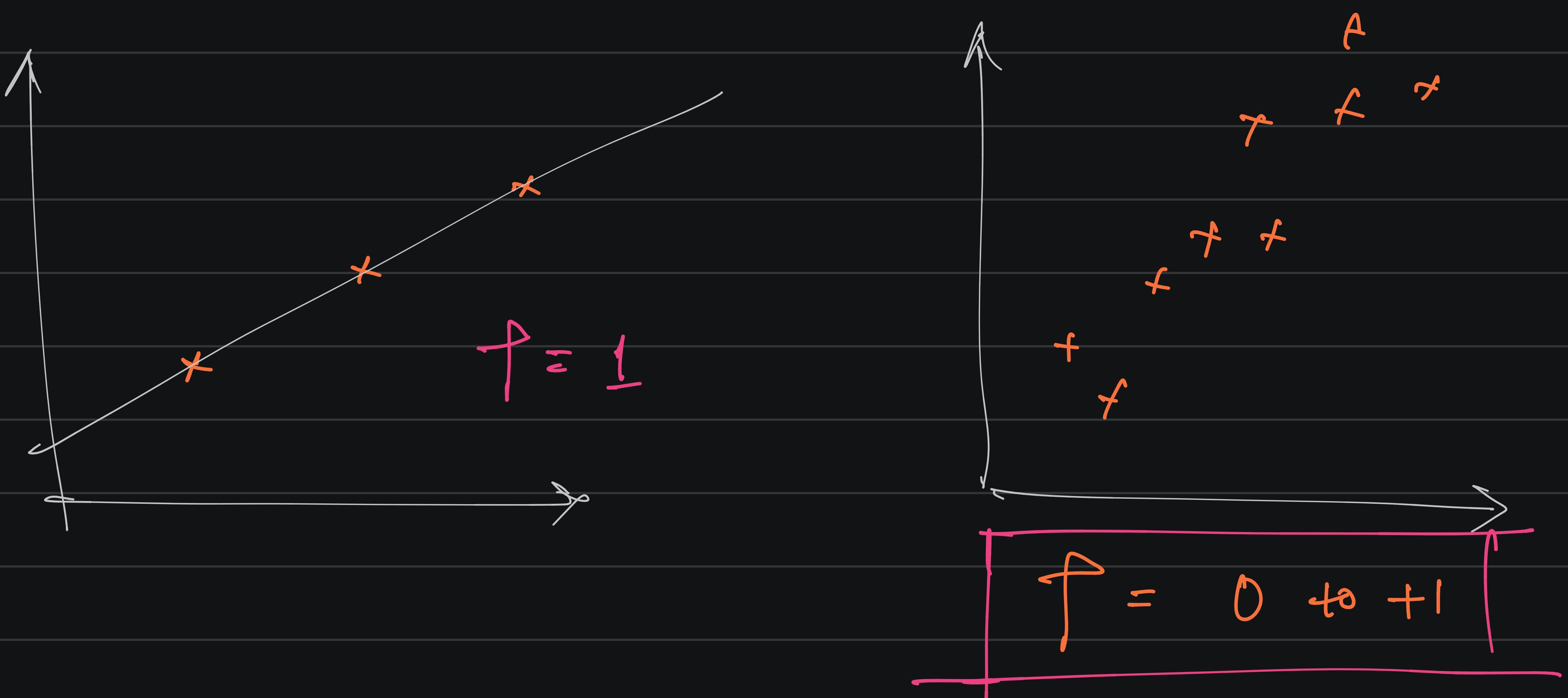
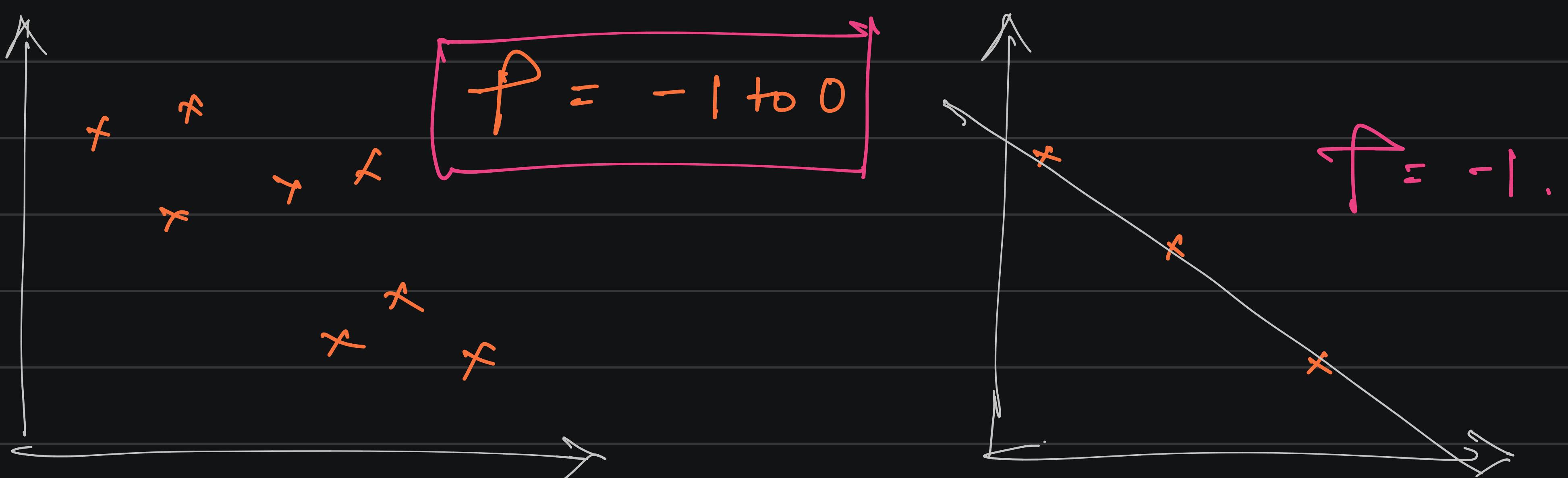
$$x \propto \frac{1}{y}$$

\rightarrow proportionality sign

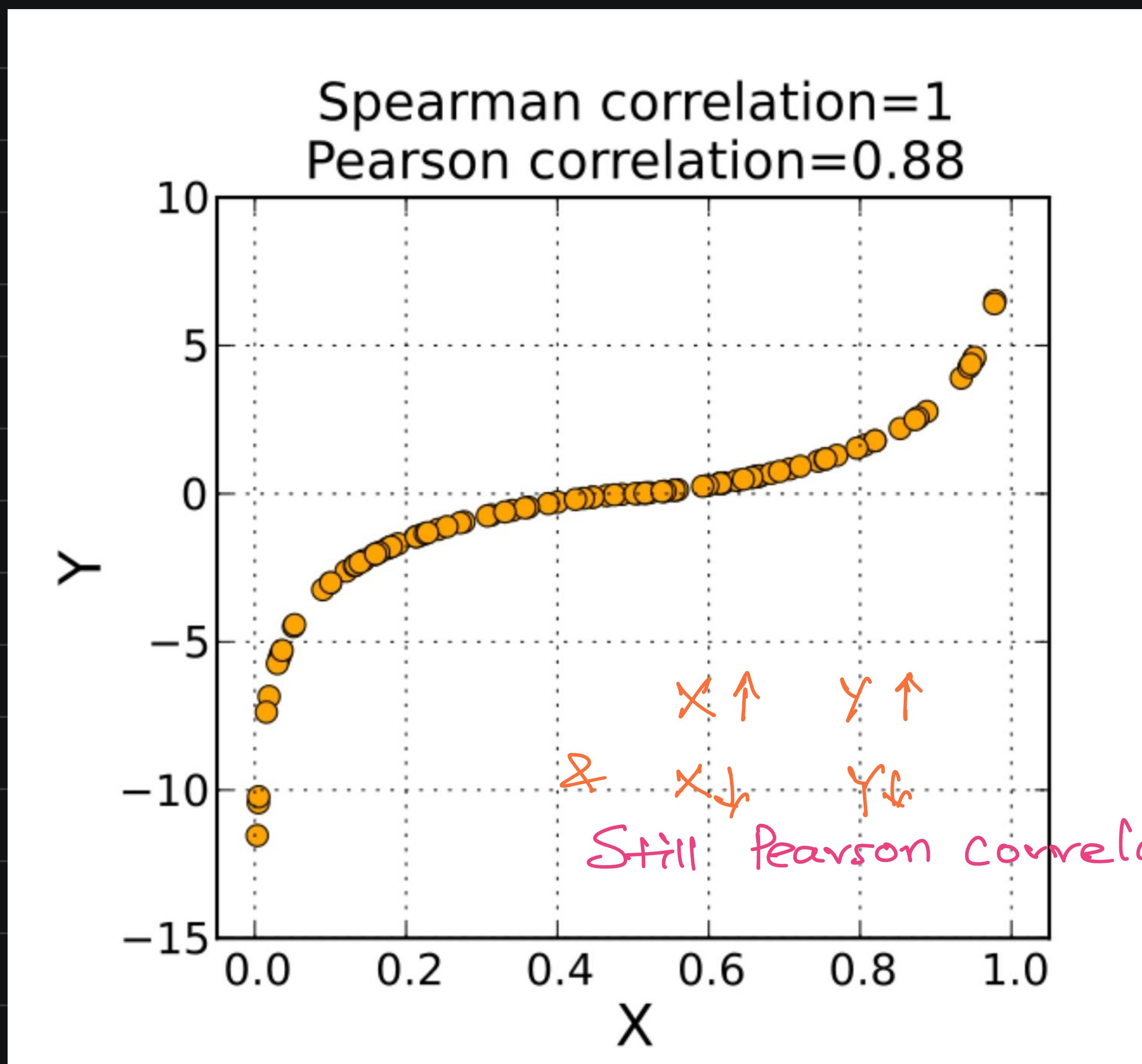
(ii) Pearson correlation coefficient. \rightarrow tries to restrict the correlation between -1 and 1.

$$\rho = \frac{\text{Cor}(X, Y)}{\sigma_X * \sigma_Y}$$

More value towards +1 \rightarrow \uparrow positive correlation.
 More value towards -1 \rightarrow \uparrow negative correlation.



(iii) Spearman's rank correlation.



* Pearson correlation \rightarrow does not work.

on linear data.



So, spearman correlation

is used.



Assigns ranks.

Formula.,

$$\rho_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}$$

X	Y
10	4
8	6
7	8
6	10.

R(X)	R(Y)
4	1
3	2
2	3
1	4.

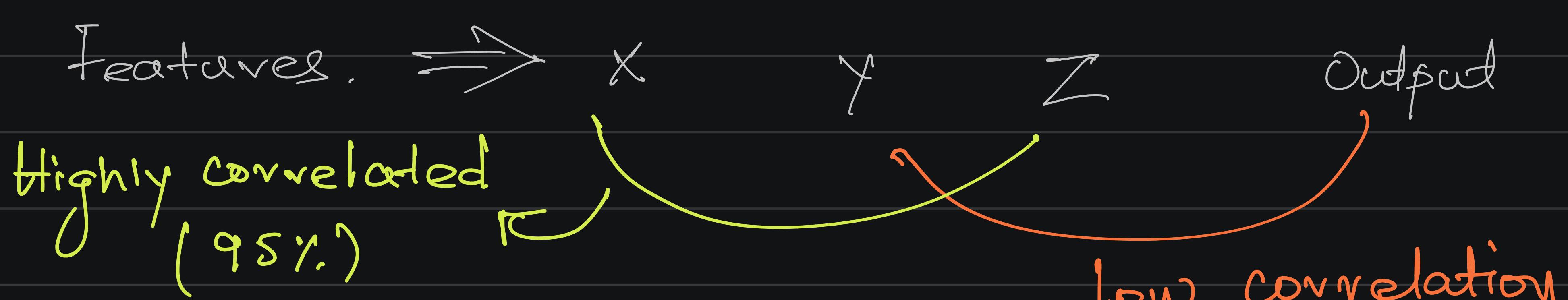
Use ranks to calculate Spearman Correlation

Keep.

Important feature if highly correlated



Why correlation is used?



One of the features can be dropped as both have similar impact.

Not important features

Can be dropped.

Day 5.

Inferential Statistics

(1) Hypothesis testing.

(2) p-value.

(3) Confidence interval

(4) Significance value.

Z-test

t-test

Chi square test

Anova test

(F-test)

3 distributions

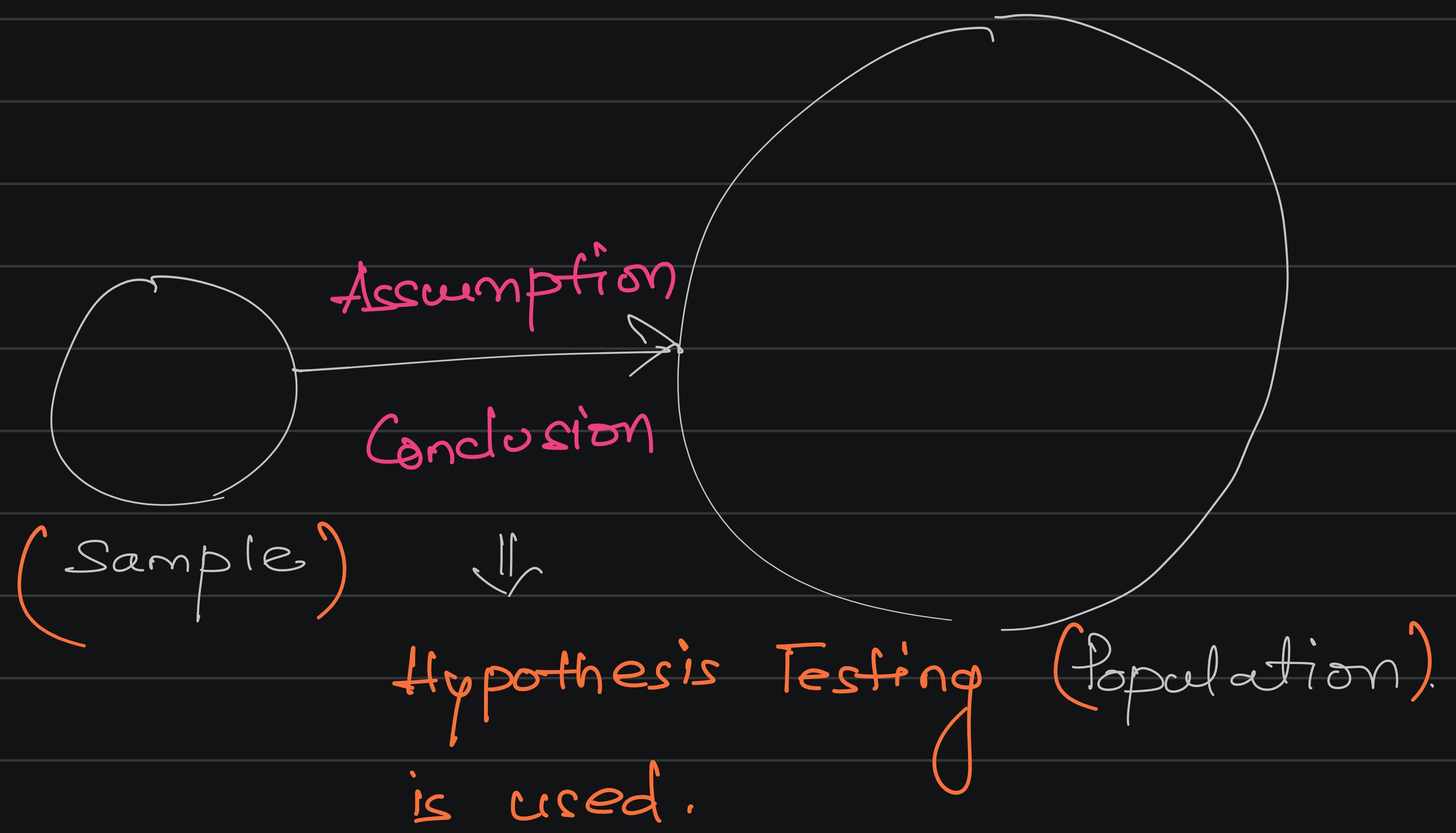
i Bernoulli's

ii Binomial.

iii Power Law.

transformation

Inferential Statistics



Steps of Hypothesis testing

Define
1. Null hypothesis: e.g. The person is not a criminal.
by default → Null hypothesis → TRUE

Define
2. Alternative hypothesis:— Opposite to null hypothesis

e.g. The person is a criminal.

(3) Perform experiments

Exp. = "Coin is fair or not"

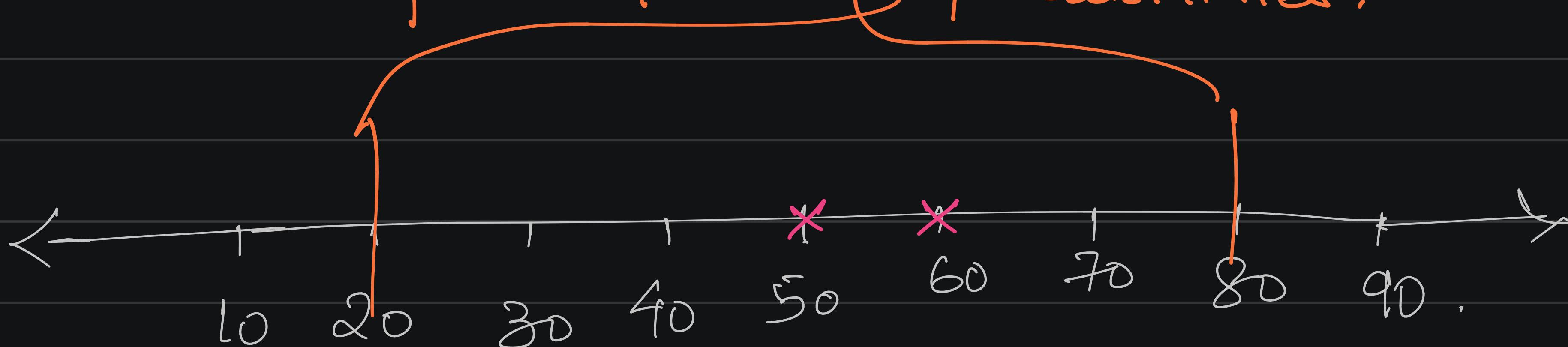
Null hypothesis \rightarrow Coin is fair.

Alternate hypothesis \rightarrow Coin is not-fair.

Experiment: Coin tossed 100 times.

$$\mu = 50, \quad SD(\sigma) = 10,$$

Define a range (confidence interval)
of acceptable probabilities.



Result: $\xrightarrow{\text{rounds}}$
(1) 50 times head. \rightarrow Fair coin

(2) 60 times head. \rightarrow Fair coin.

 if probability goes beyond Confidence
interval \longrightarrow Unfair Coin

~~e.g 2~~ Person is Criminal or not (murder).

(1.) Null hypothesis — Person is not criminal.

(2.) Alternative hypothesis — Person is criminal.

(3.) Experiment / Proofs — DNA, finger print, weapons, footage.

if

Conducted by judge.

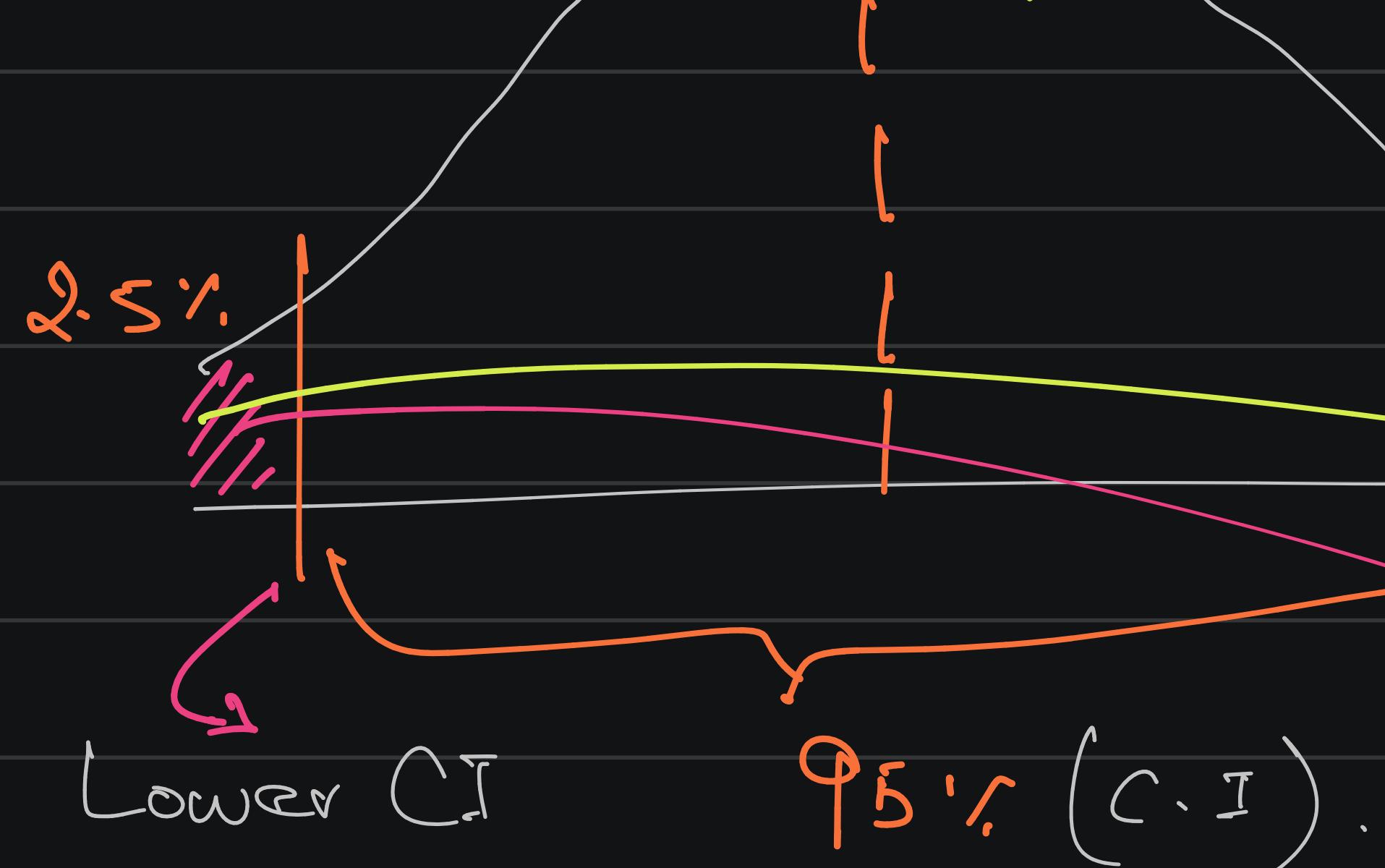
(4.) Conclusions — Accept / Reject null hypothesis

Confidence Interval (C.I.).

Significance value

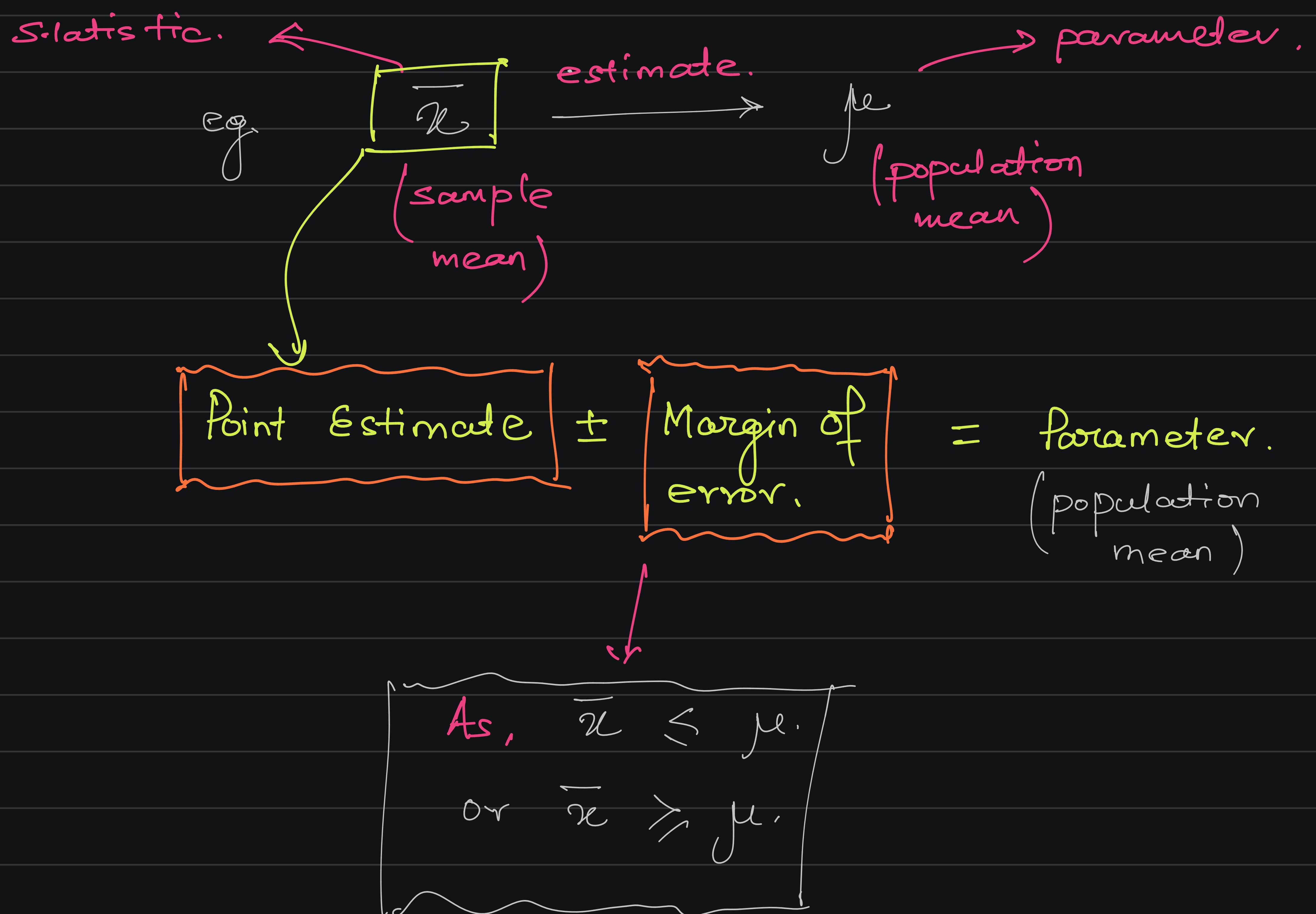
95% C.I. Accept null hypothesis.

if C.I. = 95%



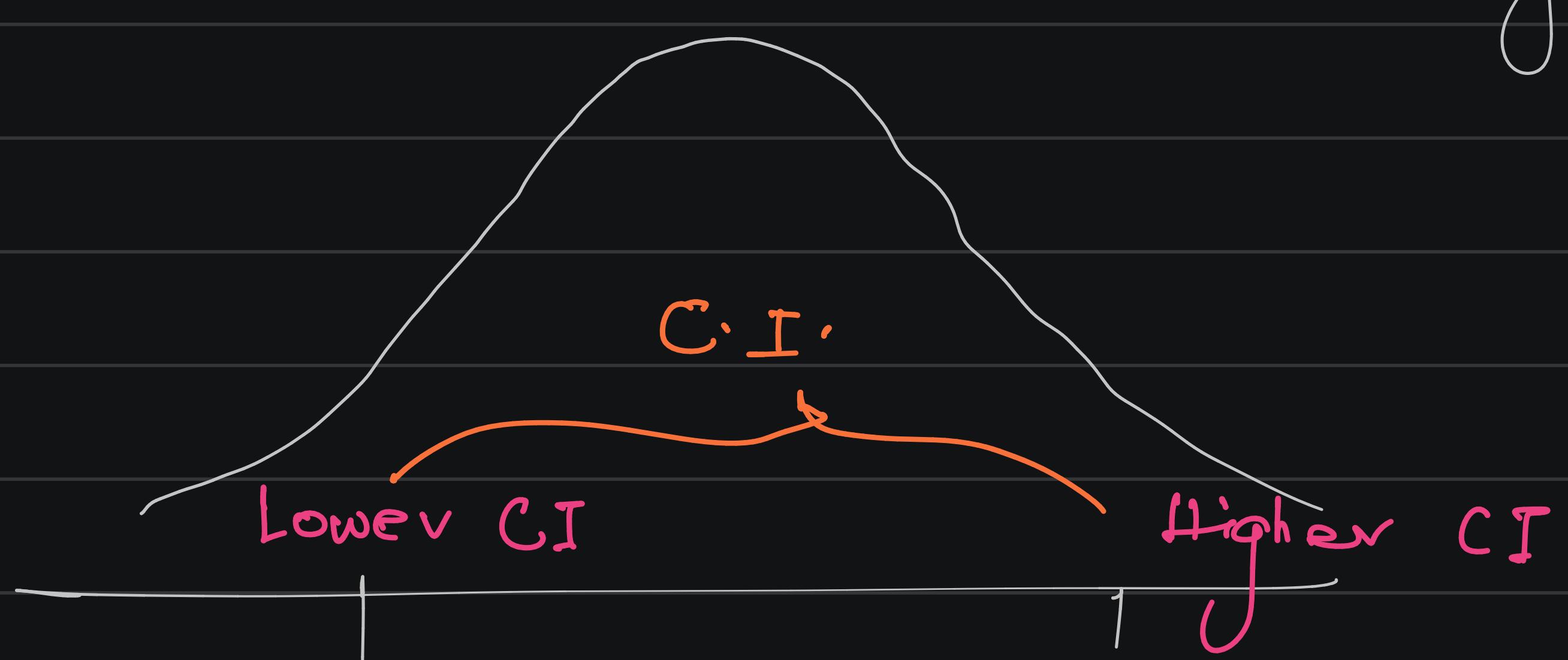
Significance value.

 Point Estimate — the value of any statistic that estimates the value of a parameter is called point estimate.



$$\text{Lower CI} = \text{Point Estimate} - \text{Margin of error}$$

$$\text{Higher CI.} = \text{Point estimate} + \text{Margin of error}$$

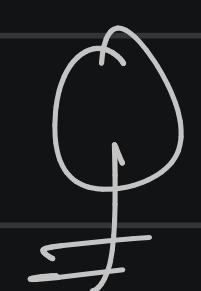




Margin of error =

$$Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

population SD
⇒ Standard error.
significance value.

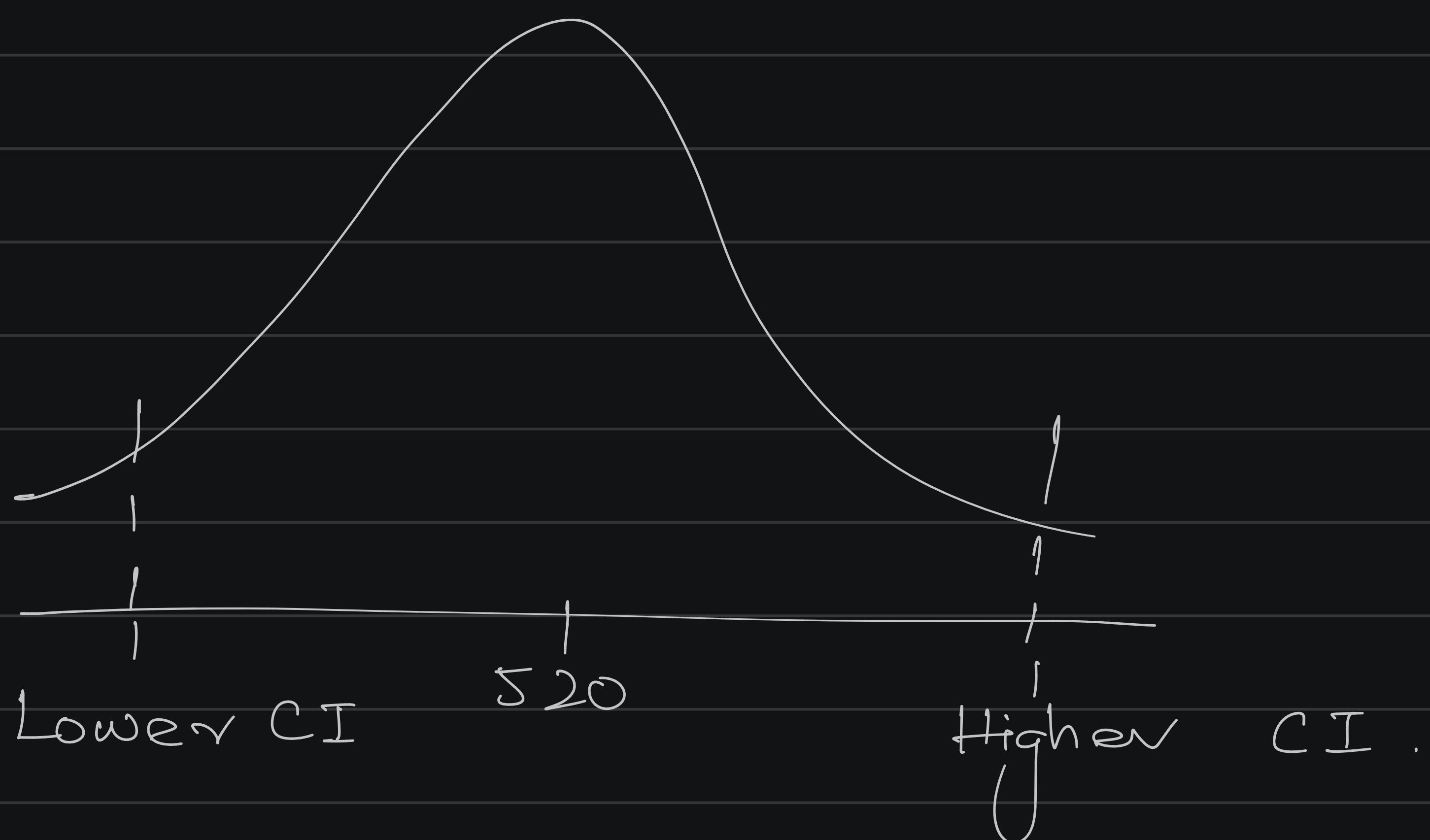


On the quant test of CAT exam, a sample of 25 test takers has a mean of 520 with a population SD of 100. Construct a 95% CI about the mean?

Ans

$$n = 25, \bar{x} = 520, \sigma = 100$$

$$CI = 95\%$$

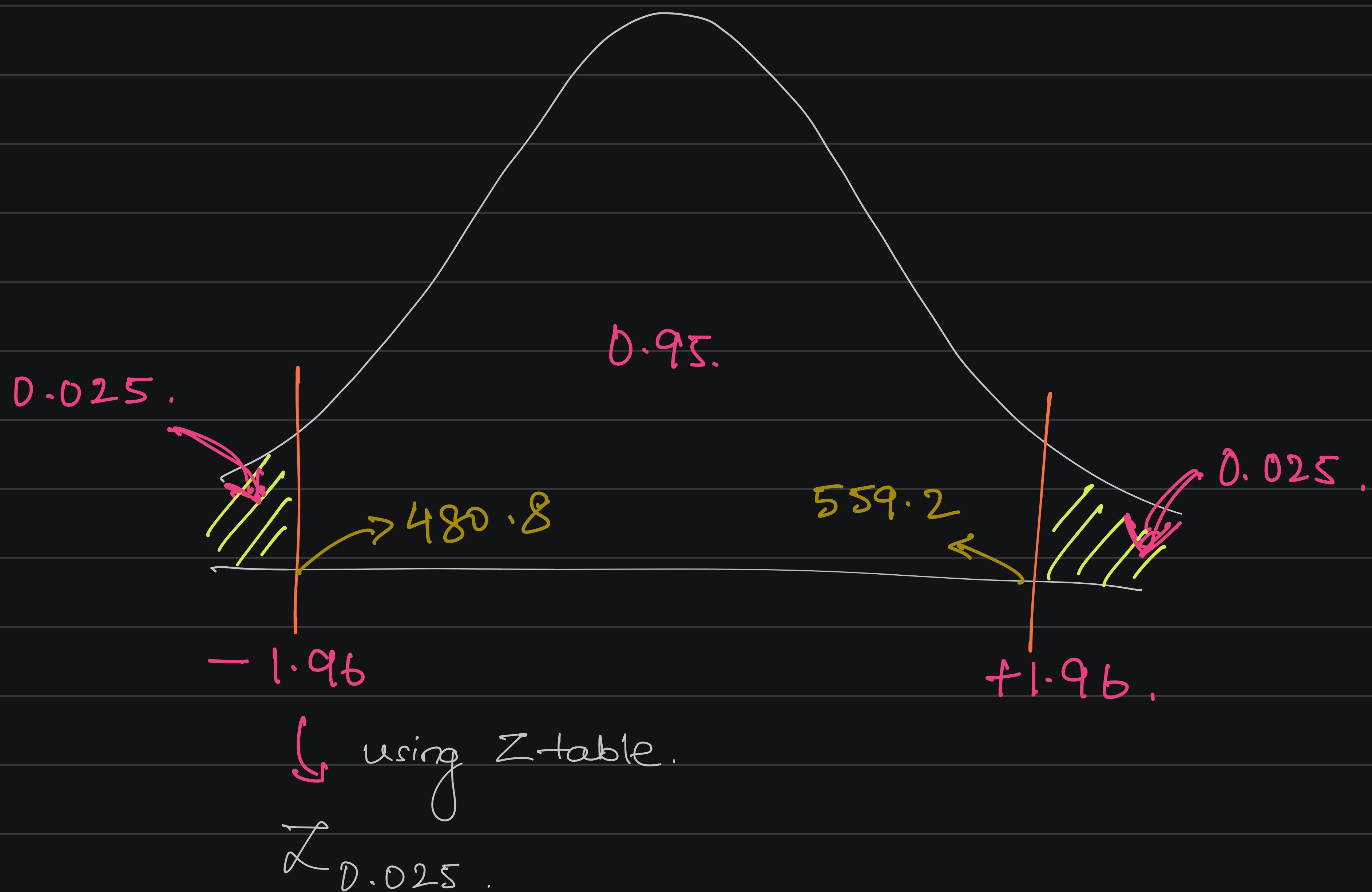


Lower CI = Point estimate - Margin of error.

$$= 520 - Z_{0.05/2} \frac{\sigma}{\sqrt{n}}$$

$$= 520 - Z_{0.025} \frac{100}{\sqrt{25}}$$

$$= 520 - (1.96 \times 20) = 480.8$$



$$\text{Higher CI} = 520 + 1.96 \times 20 = 559.2$$

∴ C.I. = { 480.8 to 559.2 }

$$Q \bar{x} = 480, \sigma = 85, n = 25.$$

$$CI = 0.90, (90\%)$$

Find the C.I range.

$$\text{Lower C.I} = \text{Point estimate} - \text{Margin of error}$$

$$= 480 - Z_{0.5/2} \left(\frac{85}{\sqrt{25}} \right)$$

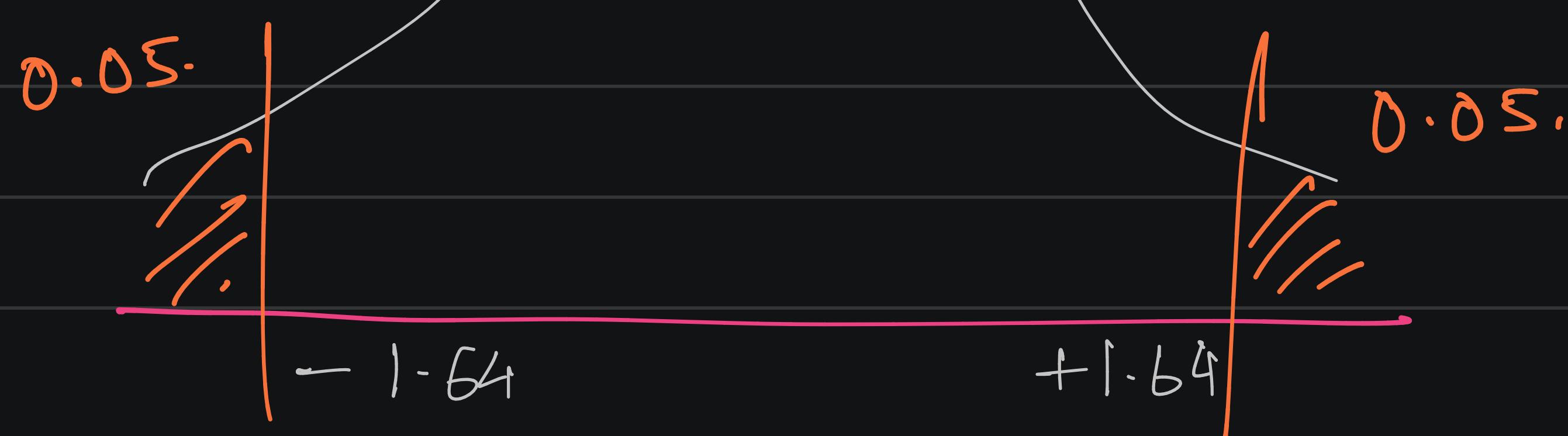
$$= 480 - Z_{0.05} * \frac{85}{5}$$

$$= 452.12$$

$$\text{Higher C.I.} = 480 + Z_{0.1/2} \left(\frac{85}{\sqrt{25}} \right)$$

$$= 507.8$$

$$\text{C.I range} = \{ 452.12 \text{ to } 507.8 \}$$



Q On the quant's test of CAT exam, a sample of 25 test takers has a mean of 520, with a sample standard deviation of 80. So. Construct 95% C.I. about the mean?

Ans. = $\bar{x} = 520$, $s = 80$, C.I. = 95%

$$S.V. = 1 - 0.95 = 0.05$$

≡

$$\left[\bar{x} \pm t_{\alpha/2} \left(\frac{s}{\sqrt{n}} \right) \right]$$

t-test

$$\text{Degree of freedom} = n - 1 = 25 - 1 \\ = 24.$$

$$\text{Lower C.I.} = 520 - \frac{t_{0.05}}{2} \left(\frac{80}{\sqrt{25}} \right)$$

$$= 520 - t_{0.025} \left(\frac{80}{5} \right)^{1/2}$$

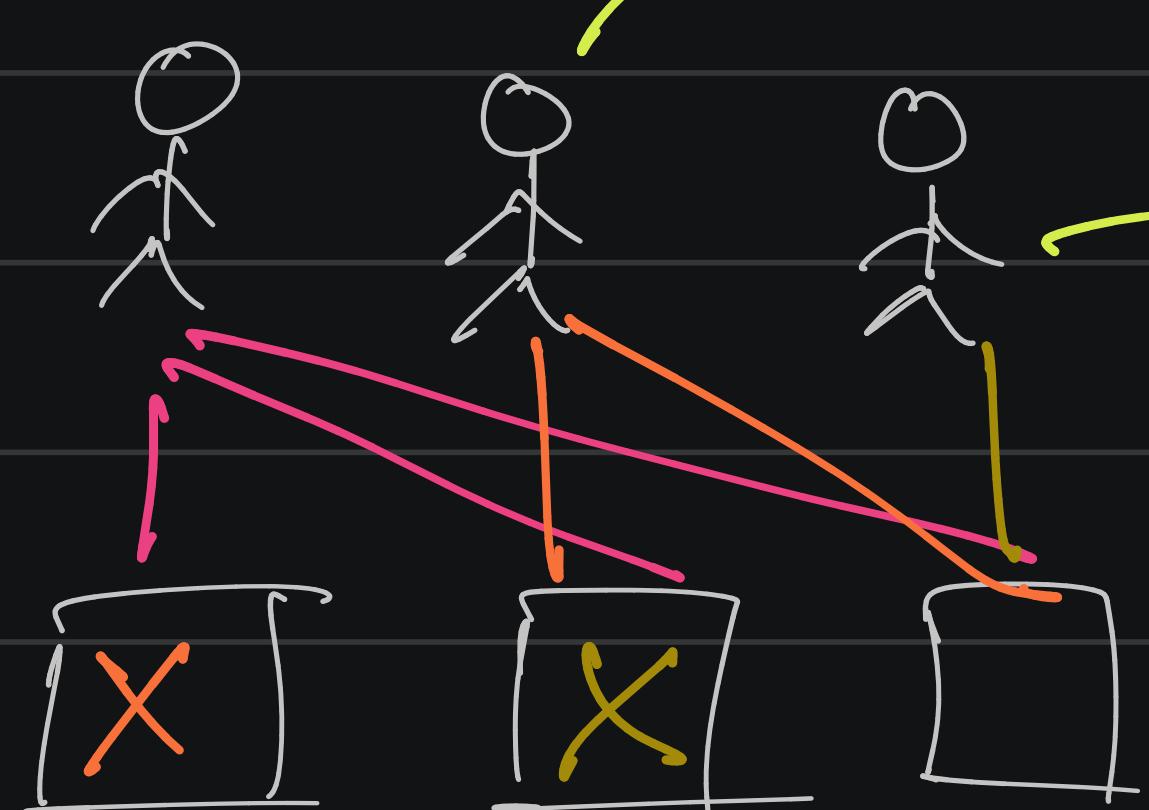
$$= 520 - \left(\frac{2.064 * 16}{5} \right)$$

$$\text{Lower CI} = 486.97 \text{ from t-table}$$



Degree of Freedom.

Choice No. 1.



Choice No. 2

→ No. choice .

→ Chars.

∴ In a sample of 'n' elements ,

$$\boxed{\text{Degree of Freedom} = n - 1}$$

$$\text{Higher CI} = 520 + 2.064 * 16.$$

$$= 553.024$$

(1) 1-tail and 2-tail test .

Q = Colleges. in Town A has 85% placement

rate. A new college was recently

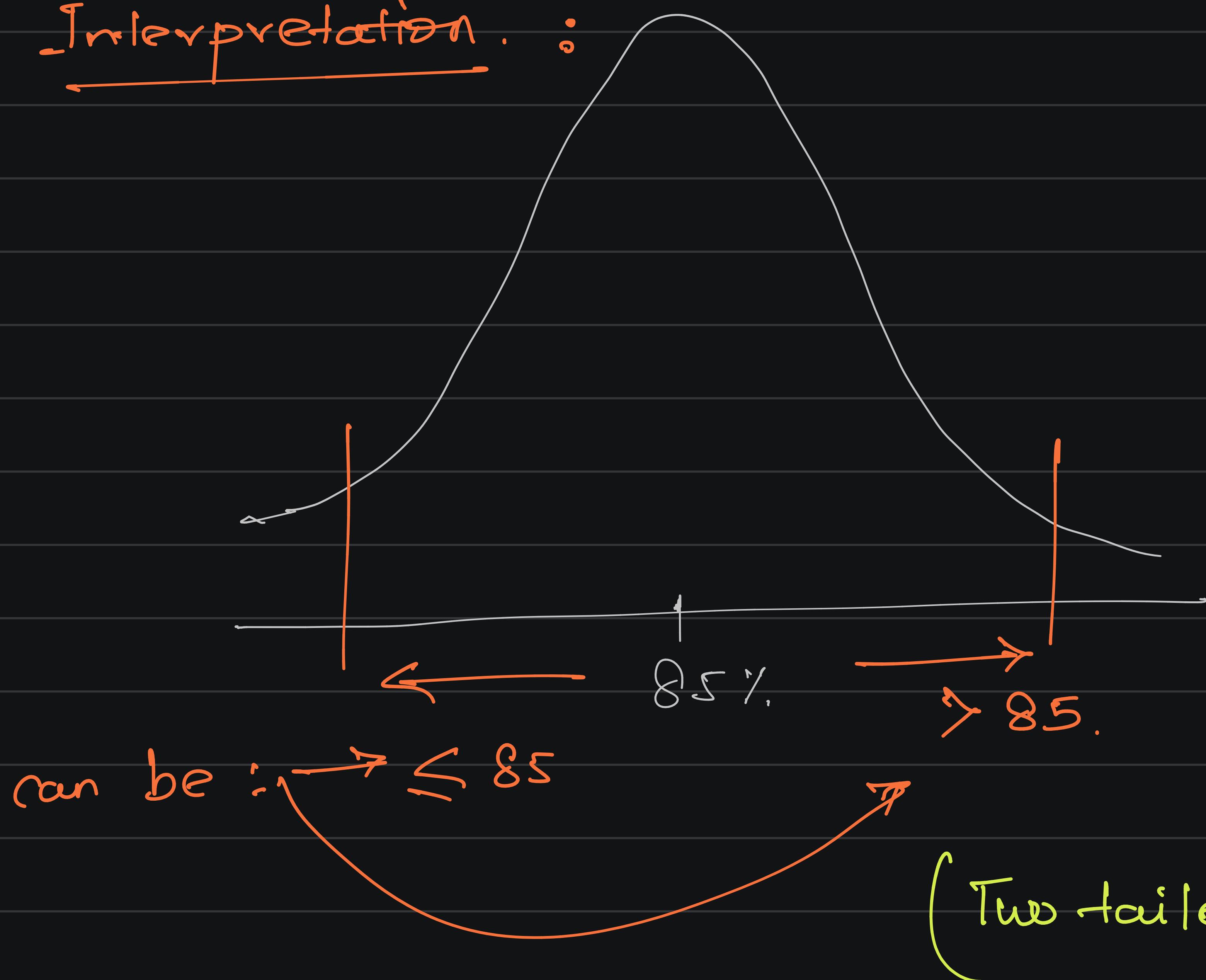
opened and it was found. that a

sample of 150 students had a placement
rate of 88% with a standard deviation .

deviation of 4%. Does this college has

- a different placement rate with CT = 95%?

Interpretation :



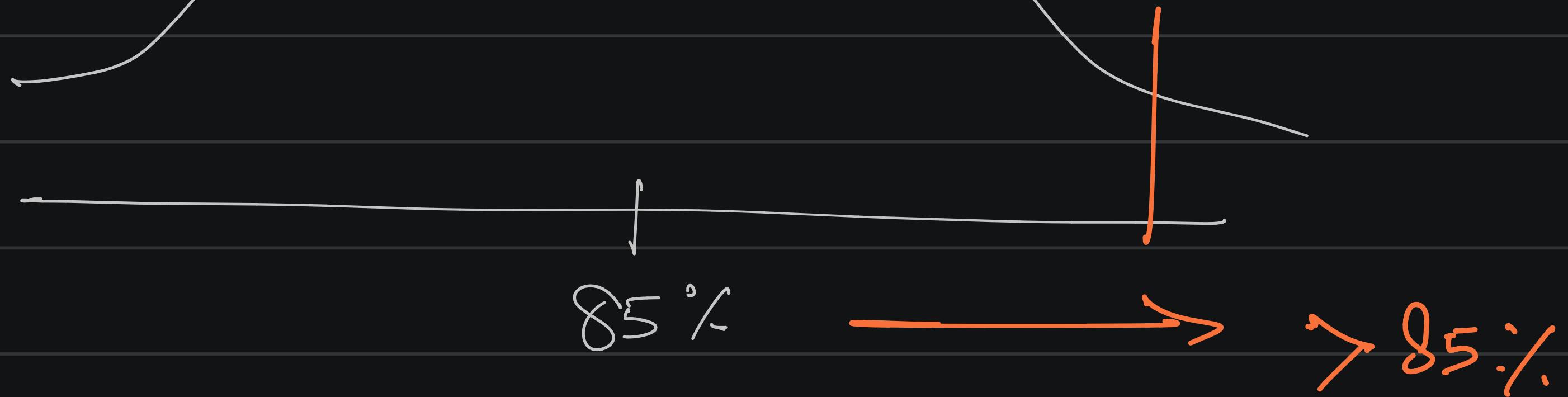
can be : ≤ 85

(Two tailed test).

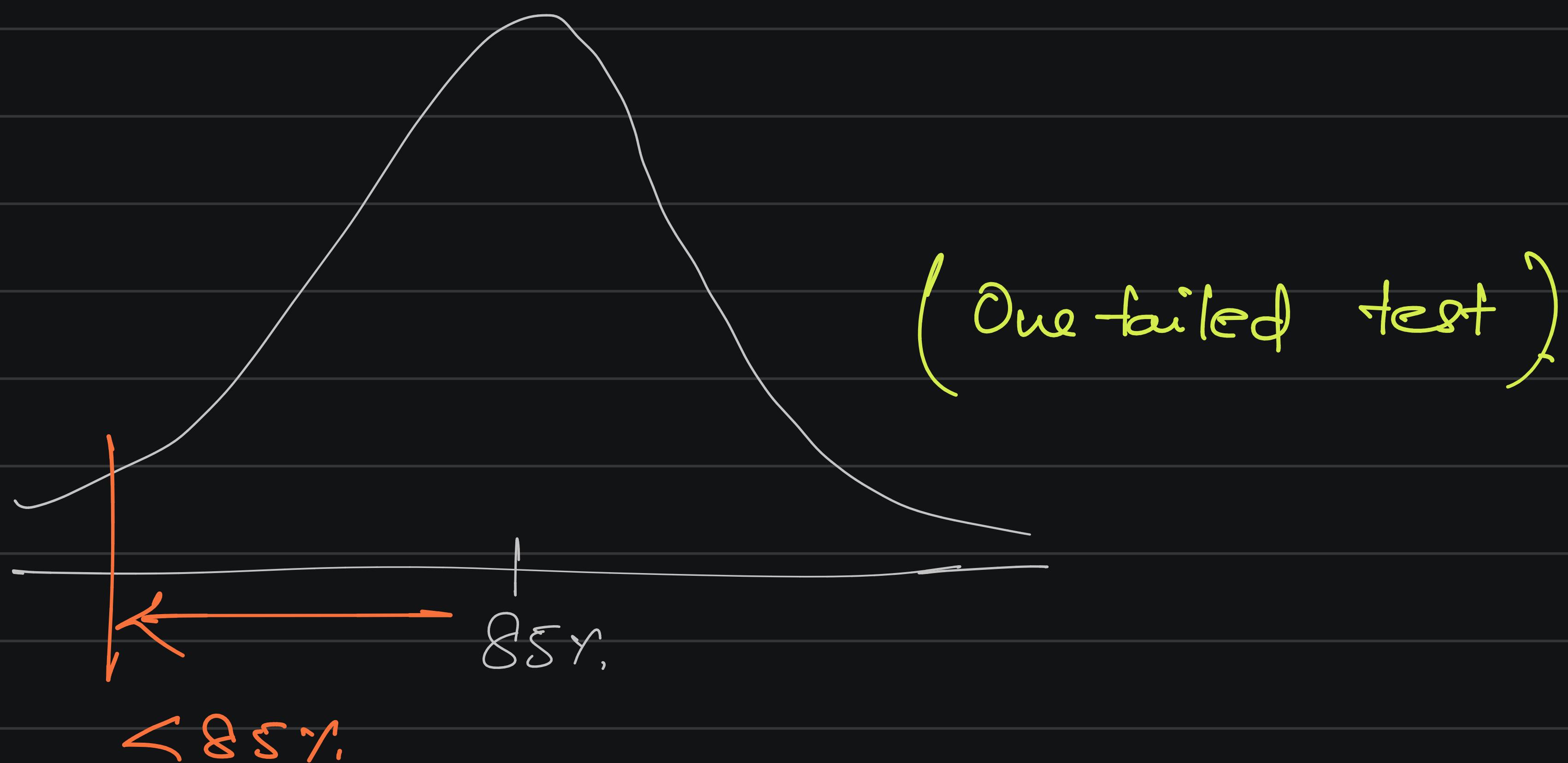
Change in question : Does this college has

a placement rate $> 85\%$.

(One-tailed test)



(if) Can placement rate be $< 85\%$.



Hypothesis testing

(1.) Z-test.

(2.) t-test.

Q = A factory has a machine that fills 80ml of baby medicine in a bottle. An employee believes, the av. amount of baby medicine is not 80ml. Using 40 samples he measured the average amount dispersed by the machine to be 78ml with a SD = 2.5.

(a.) State null and alternate hypothesis.

(b.) At $CI = 95\%$, is there enough evidence to support Machine is working properly or not.

Ans.

Null Hypothesis : $\mu = 80 \text{ ml}$.

→ machine working properly.

Alternate hypothesis : $\mu \neq 80 \text{ ml}$.

→ machine not working properly.



Step-2 : $CI = 0.95$, $\alpha = 0.05$

Step 3 : $n = 40$,

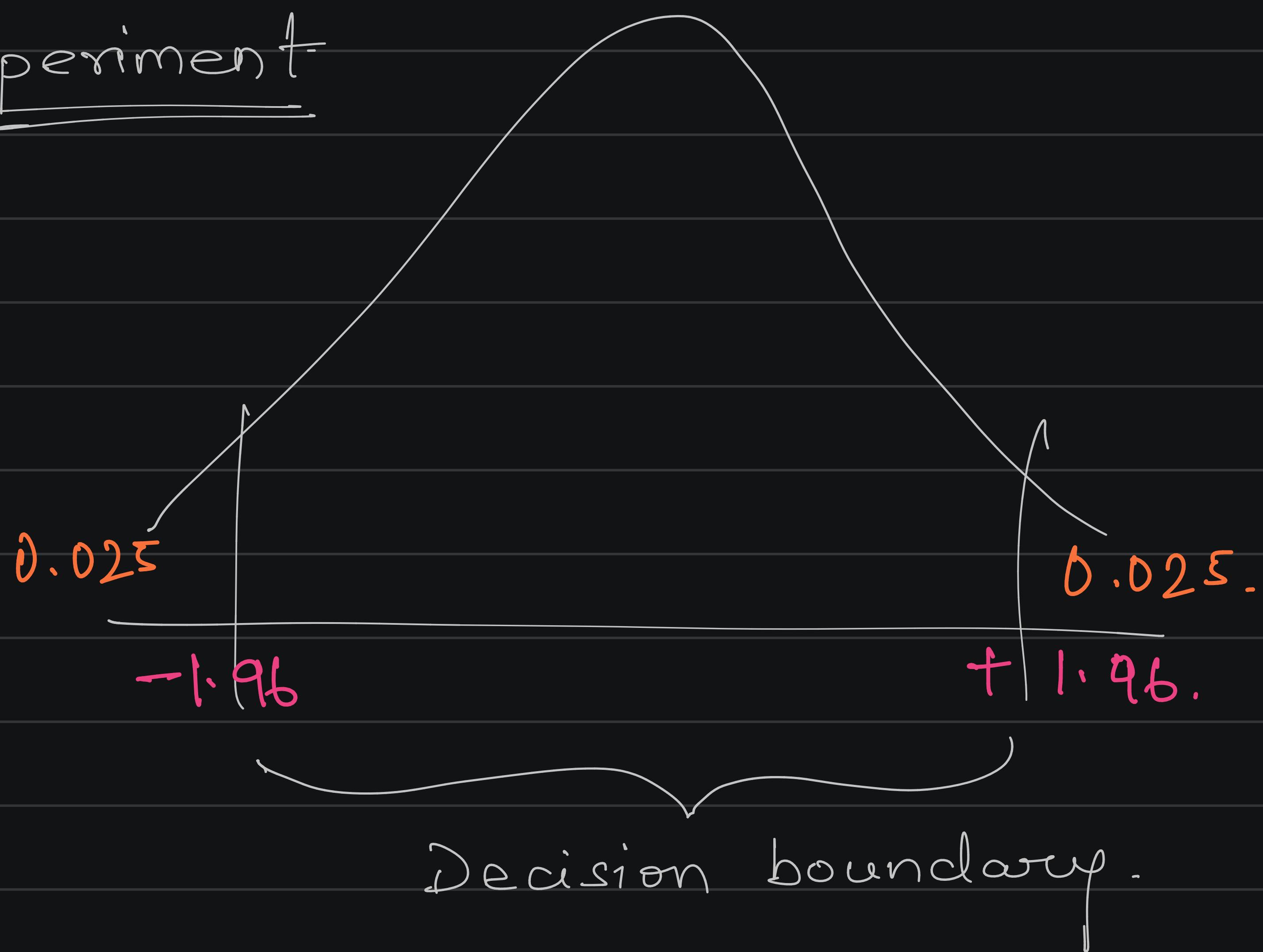
$$S = 2.5.$$

* $n \geq 30$ or population $SD (\sigma)$ given. $\left\{ \begin{array}{l} n \geq 30 \text{ or population } \\ SD (\sigma) \text{ given.} \end{array} \right\} Z\text{-test}$

$Z\text{-test reqd.}$

* $n < 30$ and Sample $SD (s)$ given. $\left\{ \begin{array}{l} n < 30 \text{ and} \\ \text{Sample } SD (s) \text{ given} \end{array} \right\} t\text{-test}$

Experiment



Calculate test statistics. (Z-test)

$$Z = \frac{\bar{x} - \mu}{\sqrt{s/\sqrt{n}}} \rightarrow \text{Standard error.}$$

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

$$= \underline{-5.05}$$

Conclusion: —

Decision rule: If $Z = -5.05$ is less than -1.96 or greater than $+1.96$, reject null hypothesis.

So, we rejected the null hypothesis.

Conclusion : Machine is not working properly.

Q) A complain was registered, the boys in a Govt. School are underfed. Average weight of the boys of age 10 is 32 kgs with $SD = 9 \text{ kgs}$. A sample of 25 boys were selected from the school, and the average weight was found to be 29.7 kgs.

With $C.I = 95\%$, check if it is true or false.

Ans. $\sigma = 9 \text{ kg}$, $n = 25$.

$\mu = 32 \text{ kg}$.

$\rightarrow Z\text{-test}$.

Step 1. Null Hypothesis: Boys are not underfed.

Alternate hypothesis: Boys are underfed.

Step 2 : C.I. = 0.95, $\alpha = 0.05$.

Step 3 :



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

$$= \frac{29.7 - 32}{9 / \sqrt{25}}$$

$$= - \frac{2.3 \times 5}{9}$$

within decision

boundary.

∴ Null hypothesis is accepted

Conclusion. — Boys in the government school are not under-fed.

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

It consists of organizing and summarizing of data.

② Inferential Stats

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.

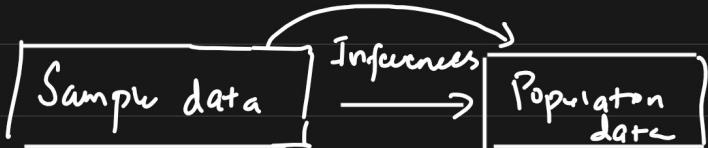
Ages {21, 20, 18, 34, 17, 22, 24, 25, 26, 23, 22}

^{min}
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?

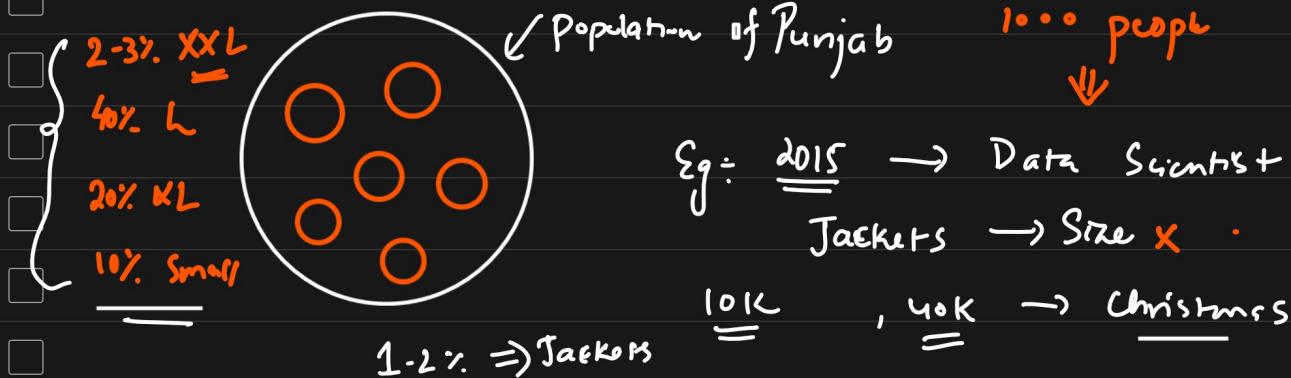


Population And Sample Data → Inferrential Statistics Results

Elections : Punjab

{ AAP, Congress }

Exit Polls ←

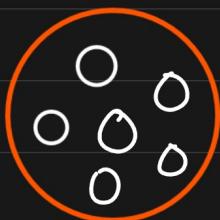


Population (N) ✓

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]

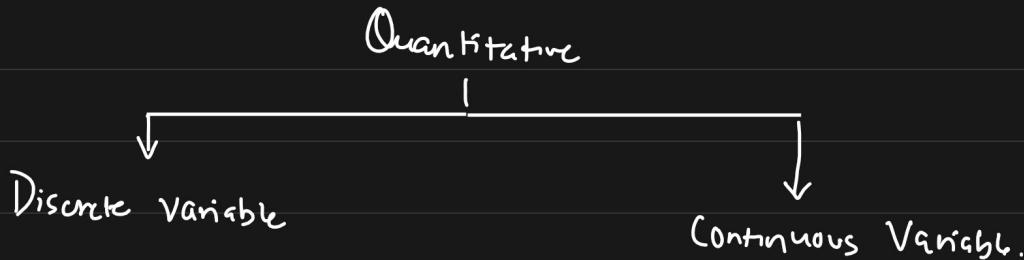
{ Quantitative → Qualitative Variable. }

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

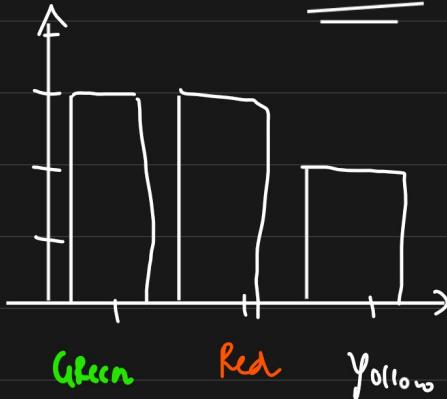
Sample Data : Green, Red, Yellow, Green, Red, Yellow, Green, Red

↓

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 } ↓

Degree

PHD

1

Salary

✓

B.G

3

✓

Master

2

✓

BCA

4

✓

12

5

✓

Assignment

↓

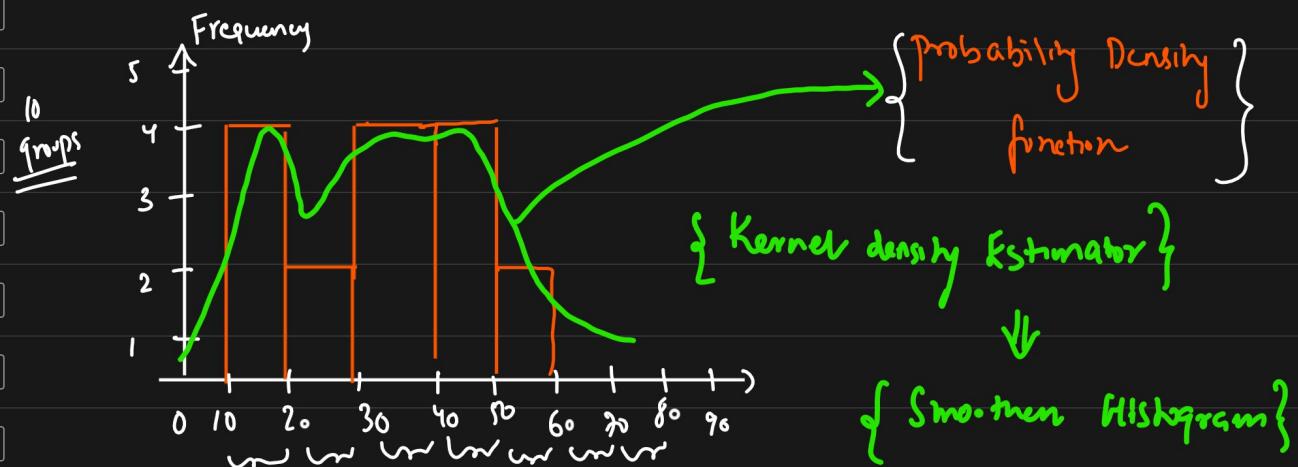
④ Ratio data ✓

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

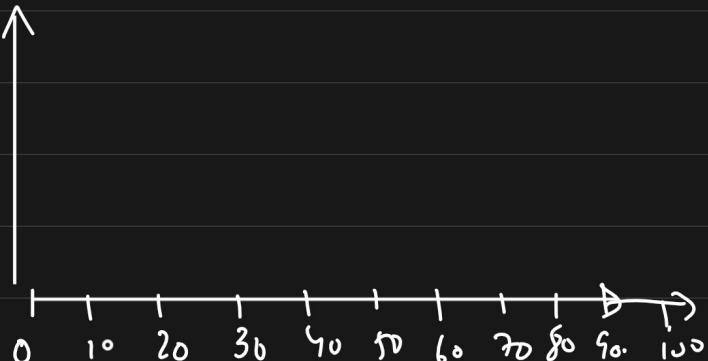
Eg: 10, 13, 18, 22, 27, 32, 38, 40, 45, 51, 56, 57, 88, 90, 92, 94, 99

bins \downarrow
10

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

40-50 50-60 60-70

70-80 80-90 90-100

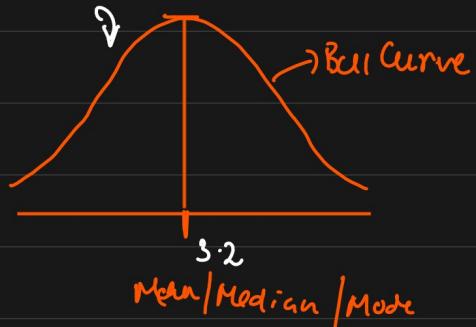


Intermediate Stats

- 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

{ }
↓

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

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

↑

$n=6$

$n-1$

$$\mu = 2.83$$

{let consider

$$10.84$$

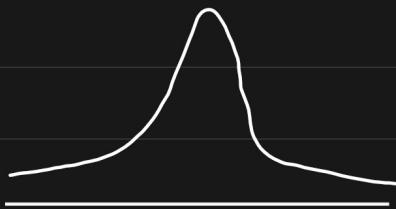
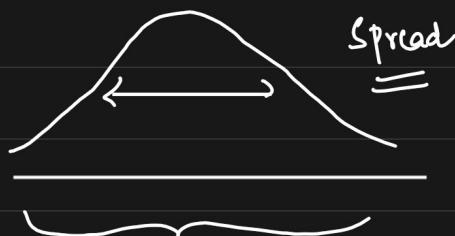
$$\frac{\sigma^2}{\text{Variance}} = \frac{6.42}{\text{as an example}}$$

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

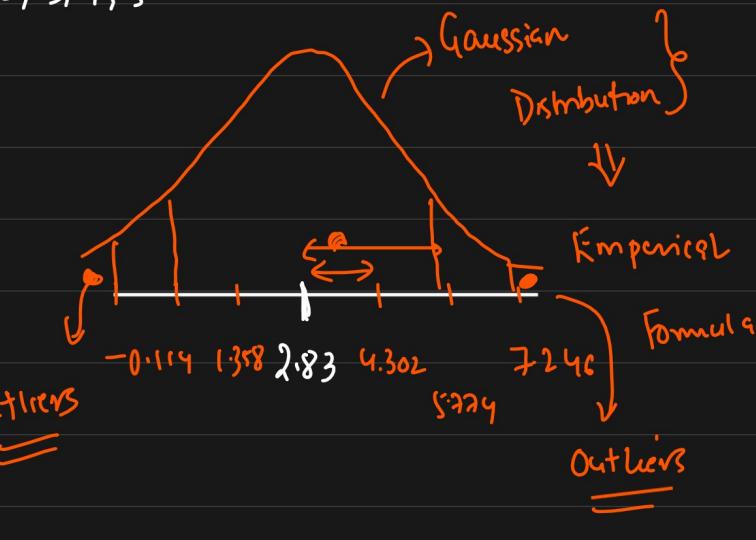
$\frac{1.358}{1.358}$

1.358

1.472

$\frac{1.358}{1.358}$

0.114



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

$$\begin{array}{r} 1.472 \\ \hline 5.274 \end{array}$$

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

(*) 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? $n=20$

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

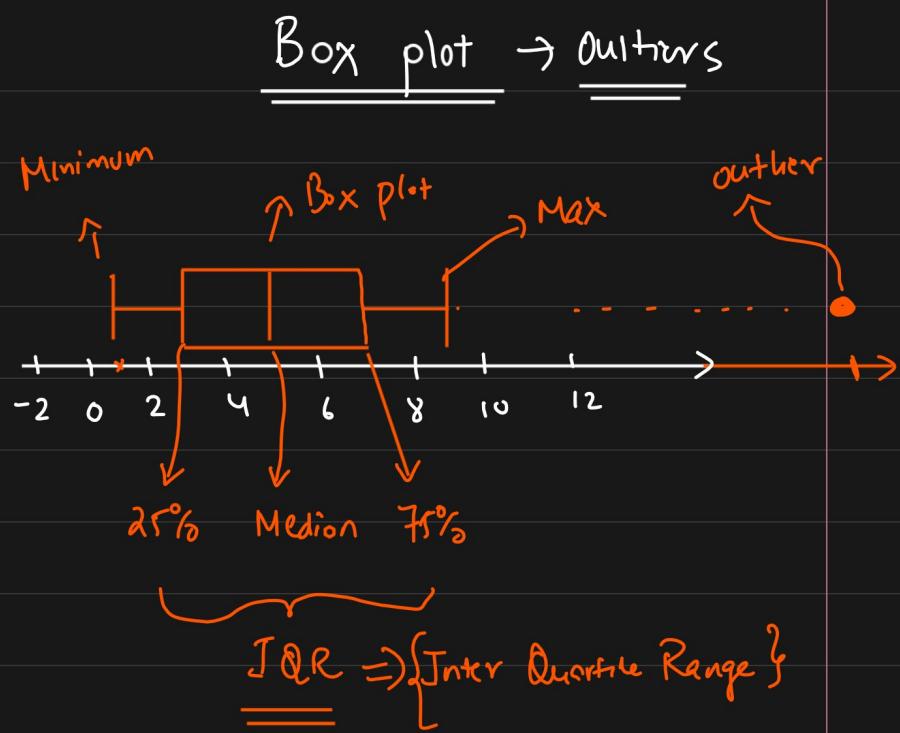
Q₁ = 3

Median = 5

Q₃ = 7

Max = 9

Use of Box plot



① Distributions

① Normal / Gaussian Distribution ✓

② Standard Normal Distribution ✓

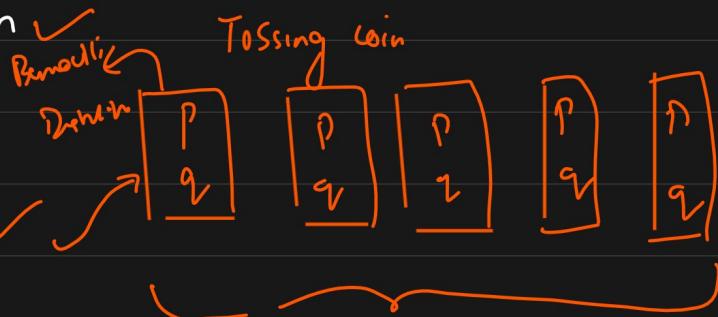
③ Z-Score ✓

④ Log Normal Distr ✓

⑤ Bernoulli's Distribution ✓

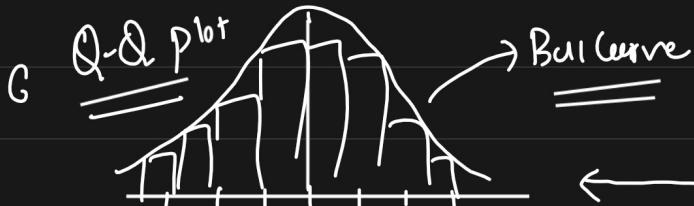
⑥ Binomial Distribution }

① Gaussian / Normal Distribution



Properties

{ Power Law }



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

↳ DATASET → IRIS Dataset } → Petal, Sepal length
domain (exposure)

② Weight of human brain

③ Height → Doctor

$$68.2, -95.4, -99.7$$

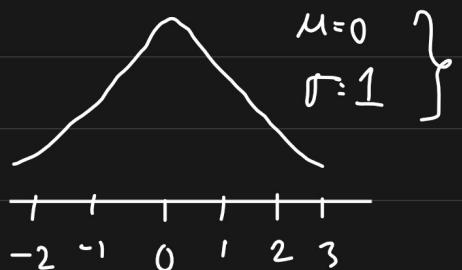
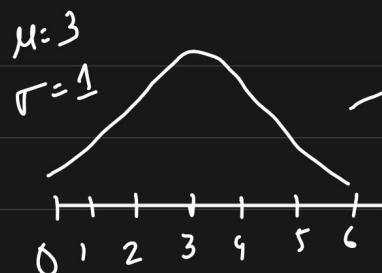
Outliers

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

$$\begin{array}{ll} \text{Age} & \text{Weight} \\ 25 & 75 \\ 26 & 80 \end{array}$$

$$28 \quad \checkmark \quad 85$$

$$30 \quad 60$$

$$32 \quad 70$$

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

$$\begin{array}{ll} \text{Salary} & \text{INR} \\ 25K & 20K \end{array}$$

$$40K$$

$$80K$$

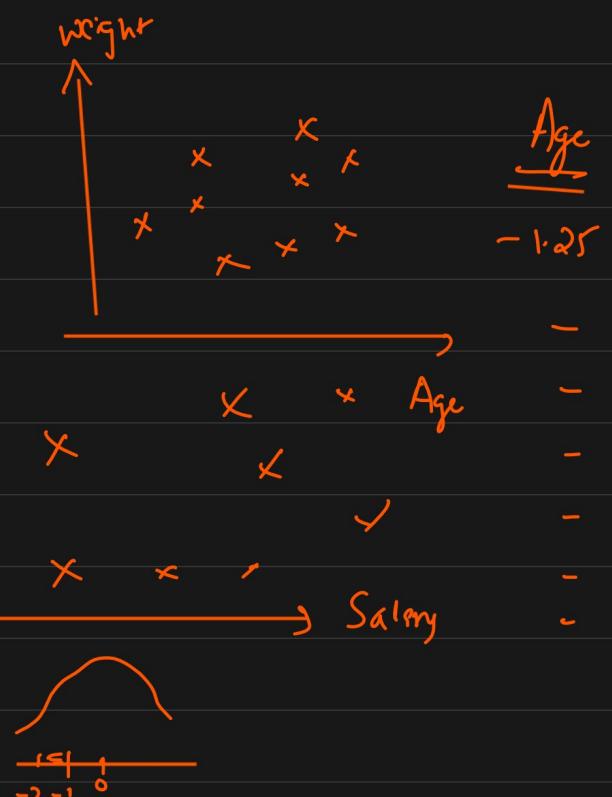
$$\text{un}$$

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

Same unit scale ??

Maths → Scale

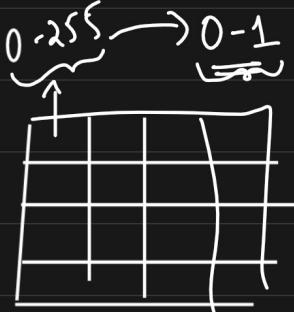
Standardization



Normalization

[Min Max Scalar]

↓
0 to 4
0 to 1



Convolutional

Neural Netw.

ML Disease ✓
Standardization

g

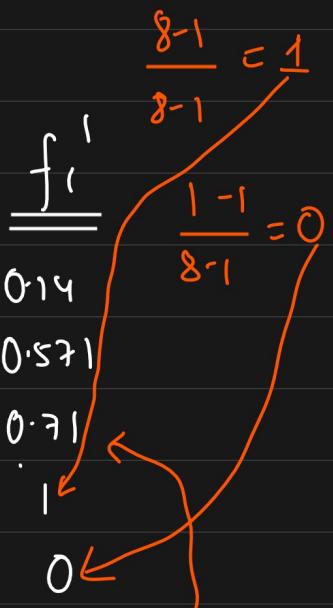
Normalization

↓
← CNN

f1

Normalization

(0 - 1)



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

↓

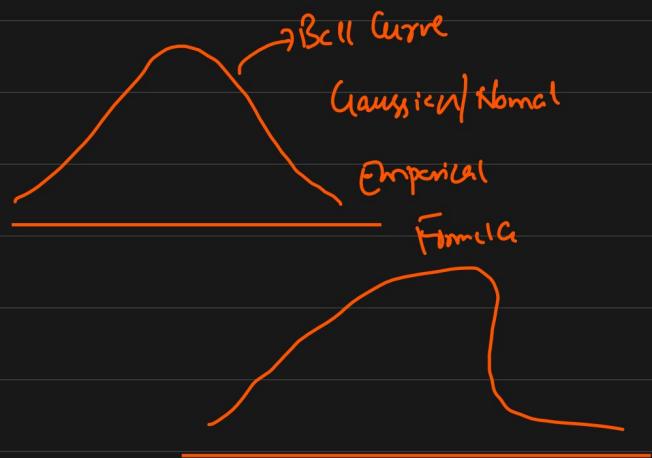
Min Max Scalar

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

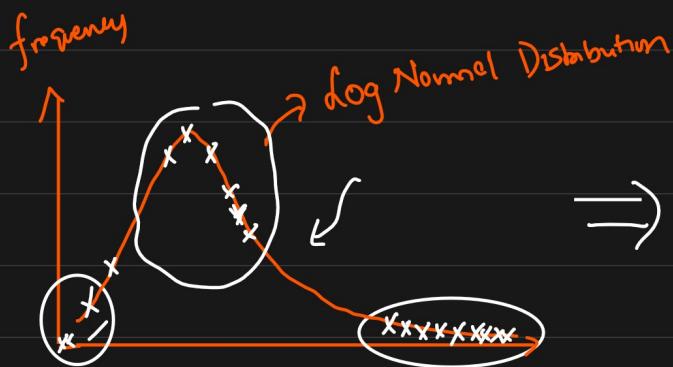
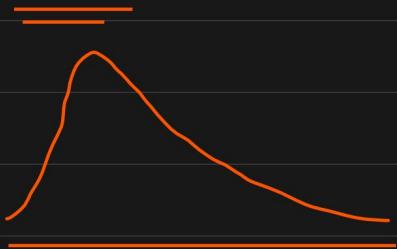
$$\frac{6 - 1}{8 - 1} = \frac{5}{7}$$

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

Log Normal Distribution

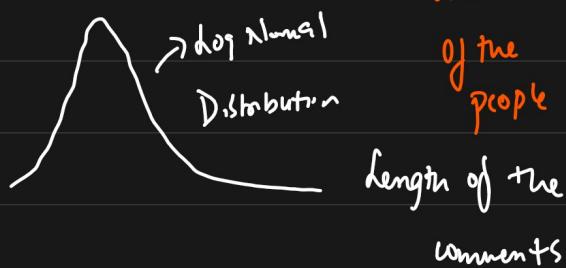
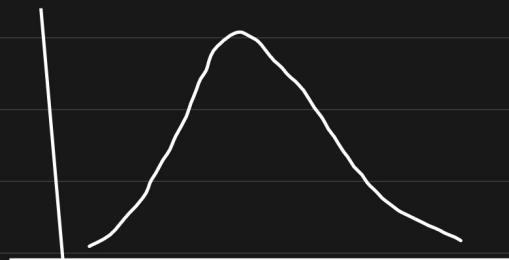


Skewed Curve



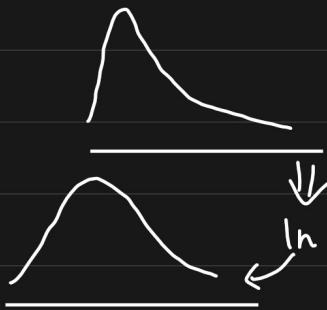
Gaussian Distribution

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

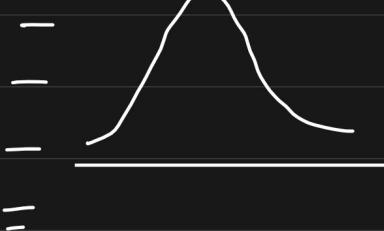
$\mathcal{Y} = \ln(x)$

25

30

40

45

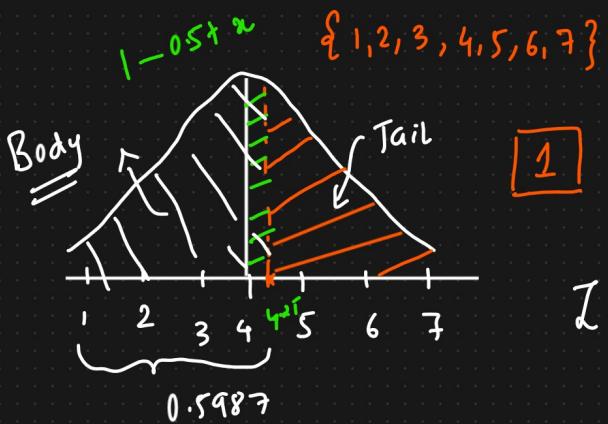


① Bernoulli's Distribution

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??

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

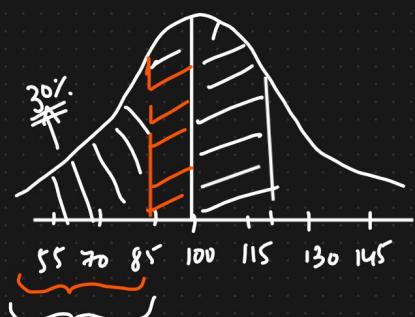
Question : What percentage of scores fall above 4.25?

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

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?

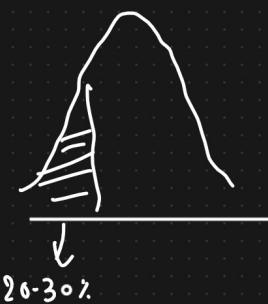
Ans)



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

① Area under this curve

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



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

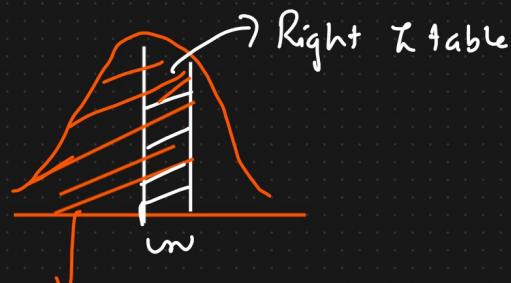


$$Z\text{score} = \frac{125 - 100}{15} = \frac{25}{15} = 1.667$$

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

1.667

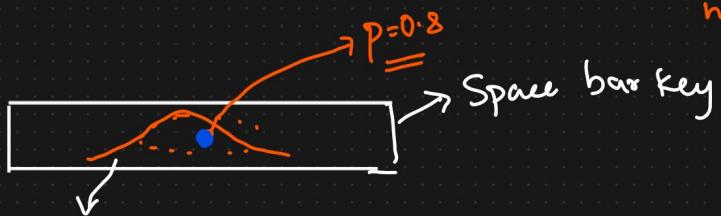
$$\underline{0.5 - 0.4515 = 0.0485} \Rightarrow 4.8\%$$



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 \rightarrow Court

SMOLAY

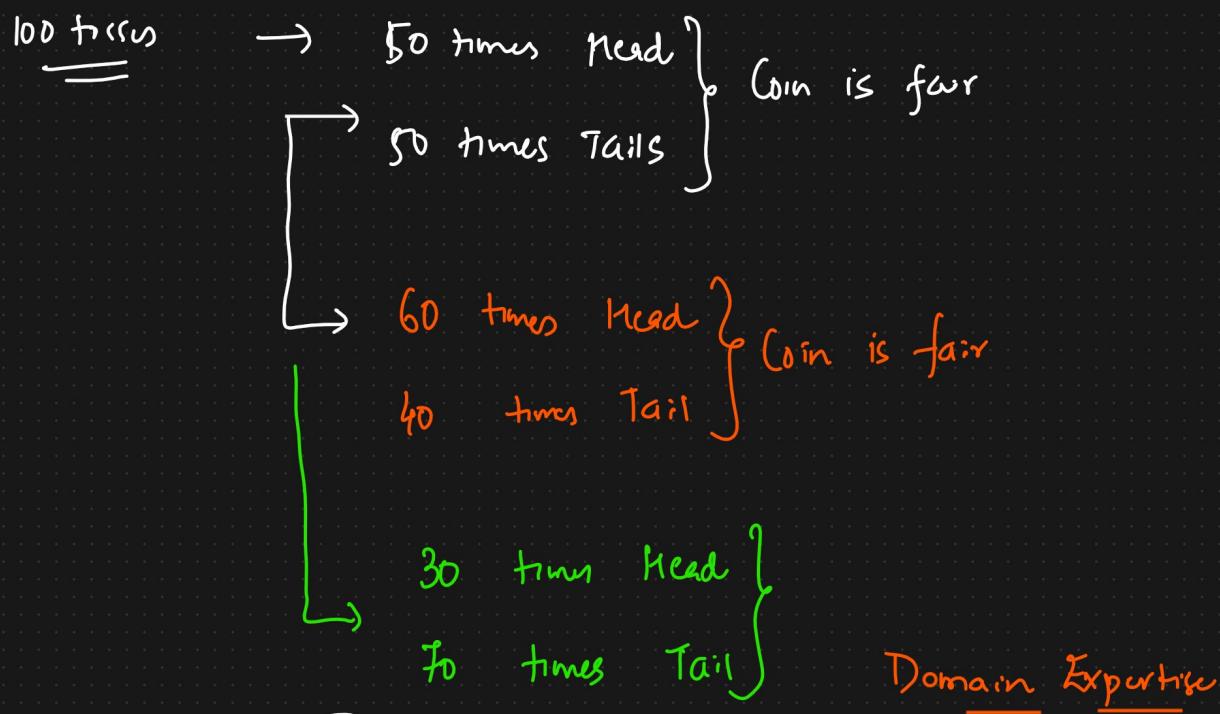
$$P(H) = 100\% \quad P(T) = 0\%$$

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

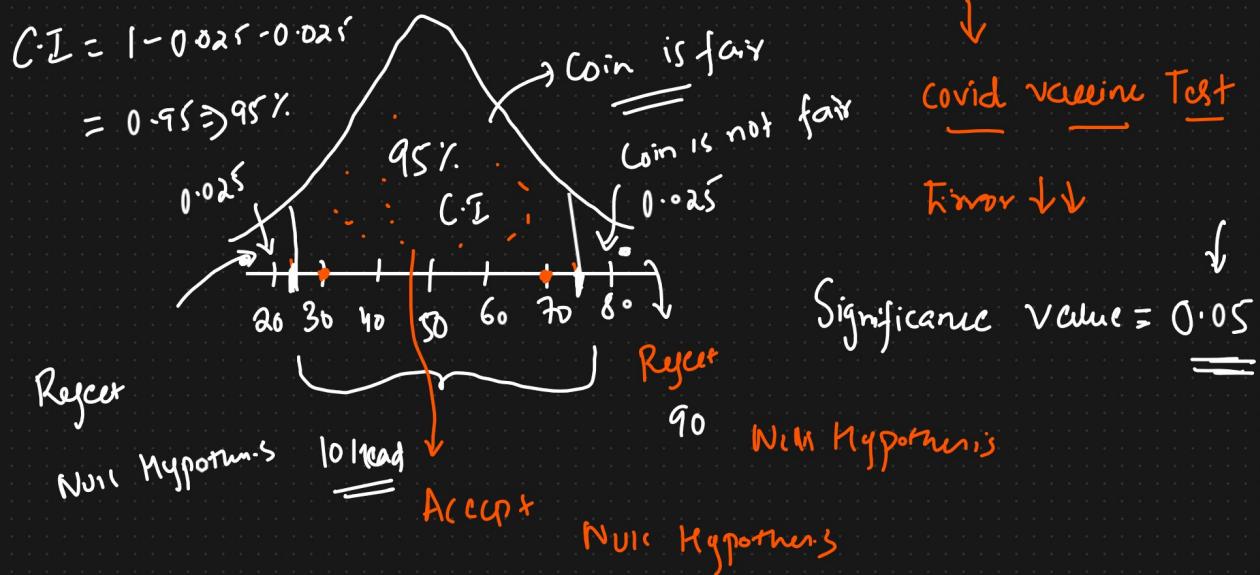
② Alternative Hypothesis — Coin is not fair $\rightarrow (H_1)$

③ Experiments

④ Reject or Accept the Null Hypothesis



Confidence Interval, Significance Values

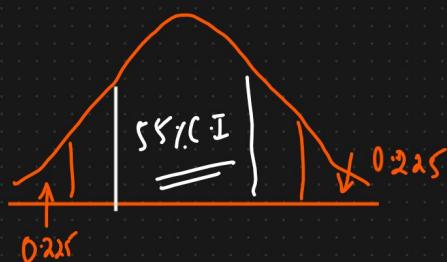


$$\lambda = 0.45$$

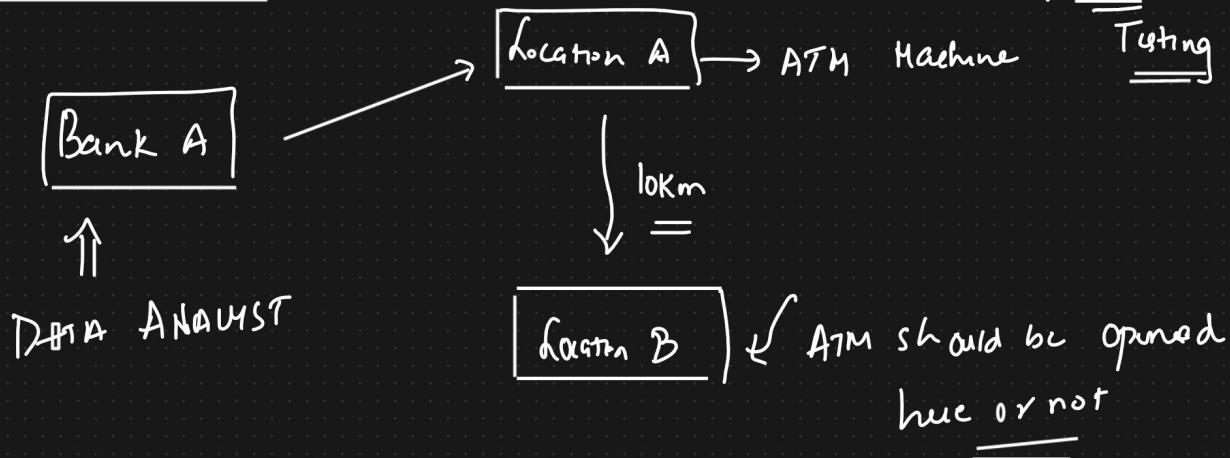
Medical

$f \uparrow \uparrow$

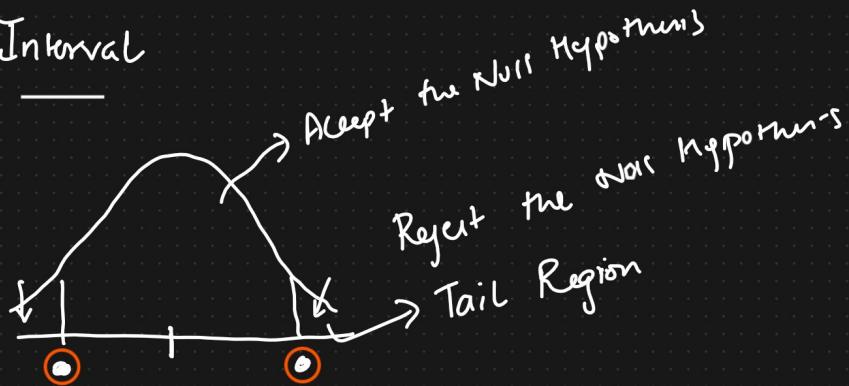
$$\frac{0.45}{2} = 0.225$$



Real World Project



① Confidence Interval



Point Estimate

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

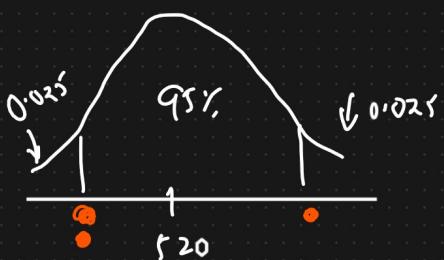


Confidence Interval

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

- Q) On the quant test of CAT Exam, the population 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. =

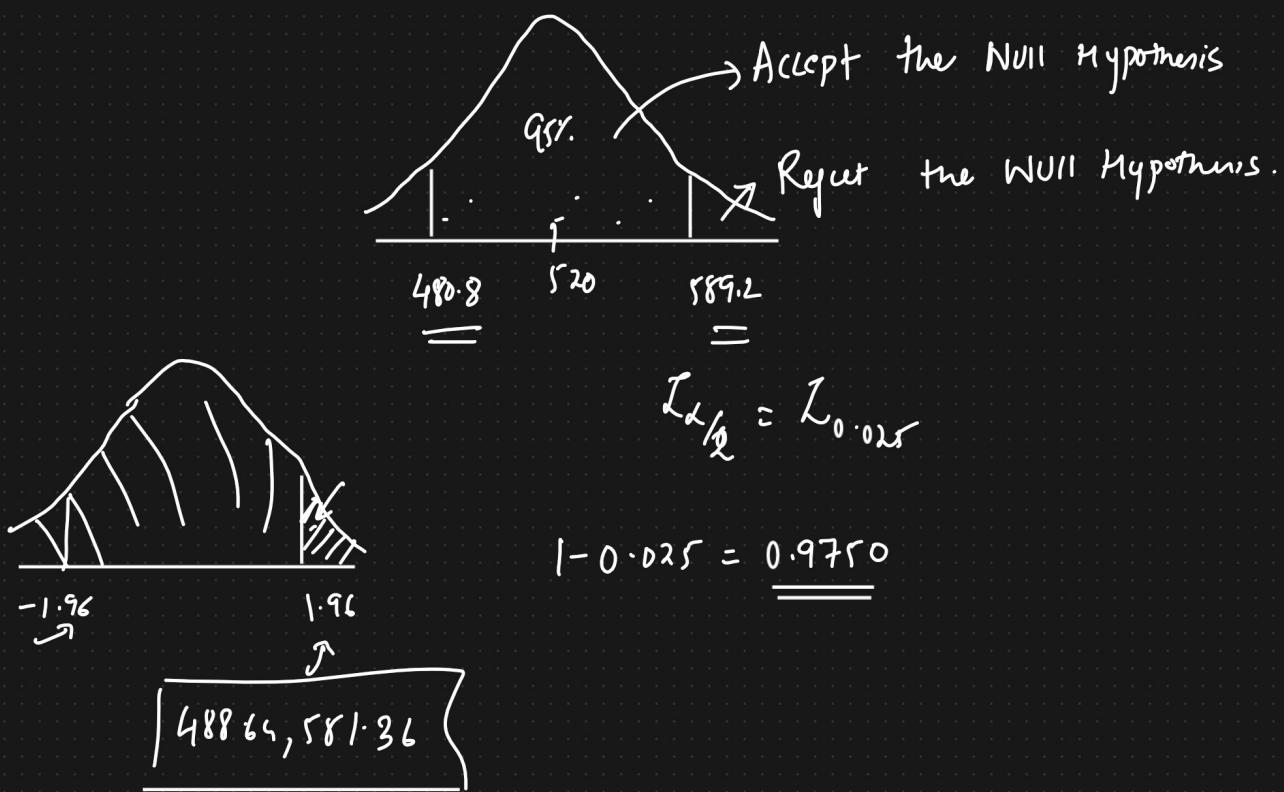
$$\bar{x} \pm Z_{\alpha/2} \left[\frac{\sigma}{\sqrt{n}} \right] \rightarrow \text{Standard Error}$$

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

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

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

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



- 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

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

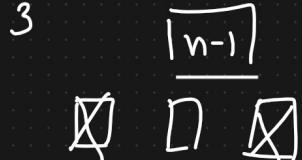
t -test $\Rightarrow t$ - table { Because population Sd is not given }

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

$$t_{0.025}$$

t -test

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



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 ✓ Null Hypothesis is True or False.

in reality if it is false → 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?

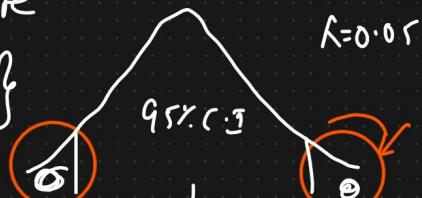
$$\alpha = 0.05 \Rightarrow 95\% \text{ C.I} \rightarrow 85\%$$

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]

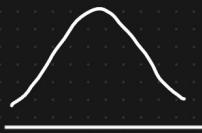
b) t test {t table}

c) Z test proportion population.

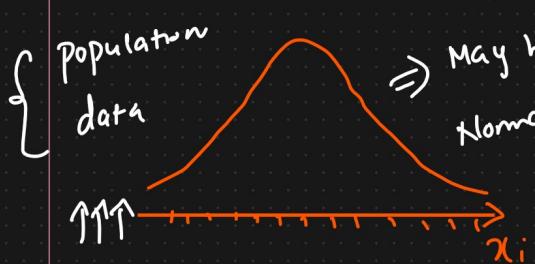
d) Chi Square (Categorical Test)

e) ANNOVA (F Test)

Influential



① Central limit Theorem



\Rightarrow May be Gaussian

$n > 30$

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

$\rightarrow \bar{x}_m$

\Rightarrow It may not

Sample m

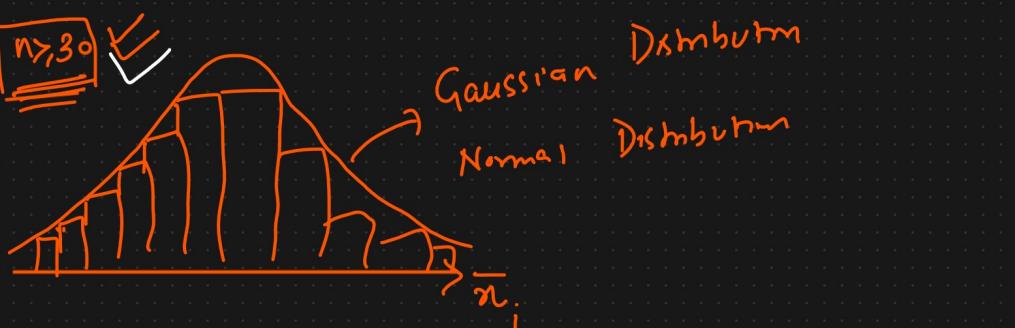
x_i

10,

$n > 30$

Sample mean

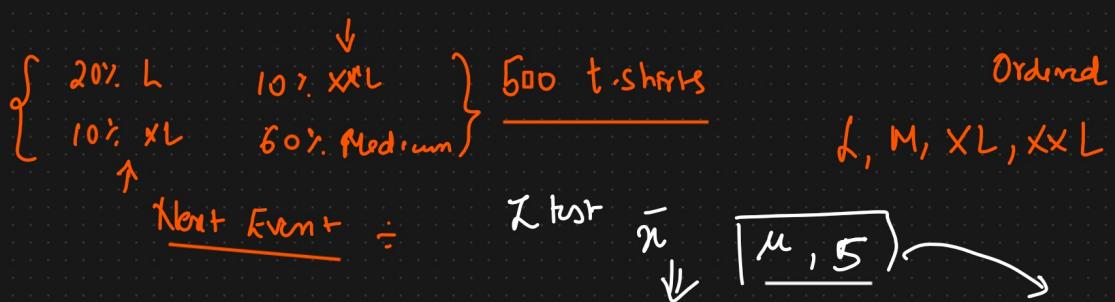
distribution



(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?

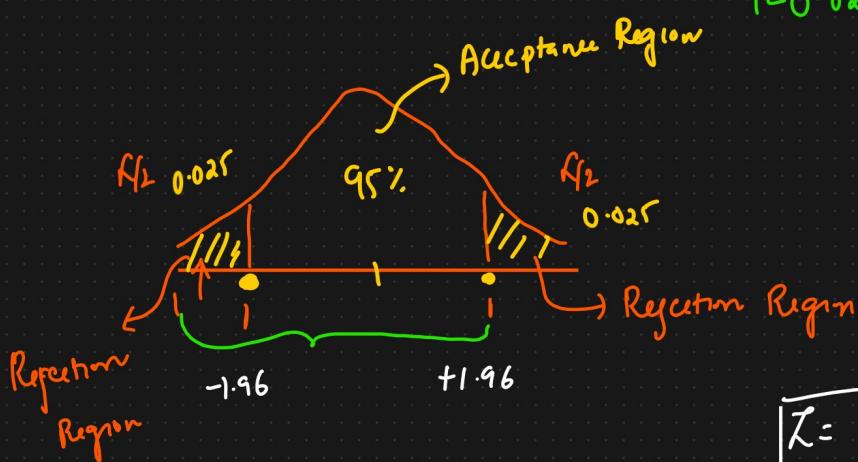
Step 2:

$$\alpha = 0.05$$

$$C.I = 95\%$$

$n > 30$ $n \leq 30$
 (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.

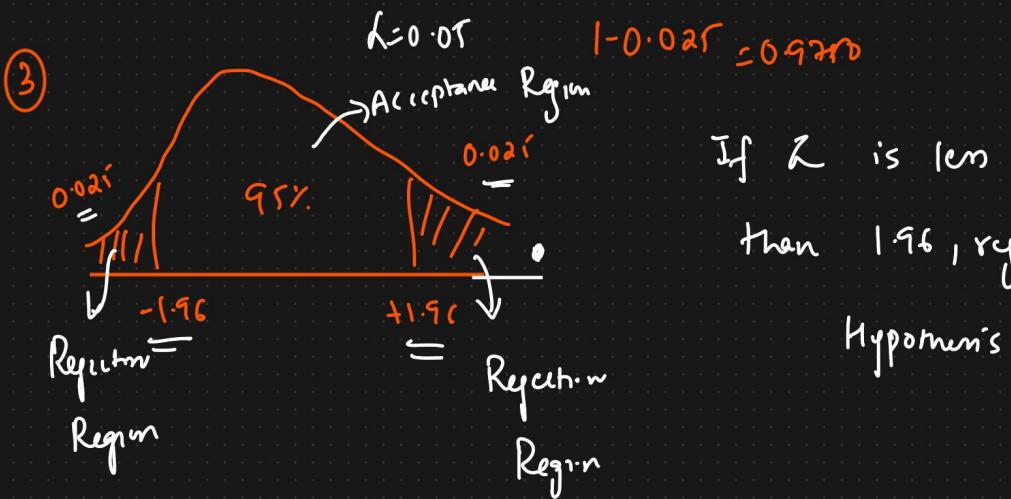
A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect Intelligence? $\left\{ \begin{array}{c} 95\% \\ \hline \downarrow \\ C.I. \end{array} \right\}$

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

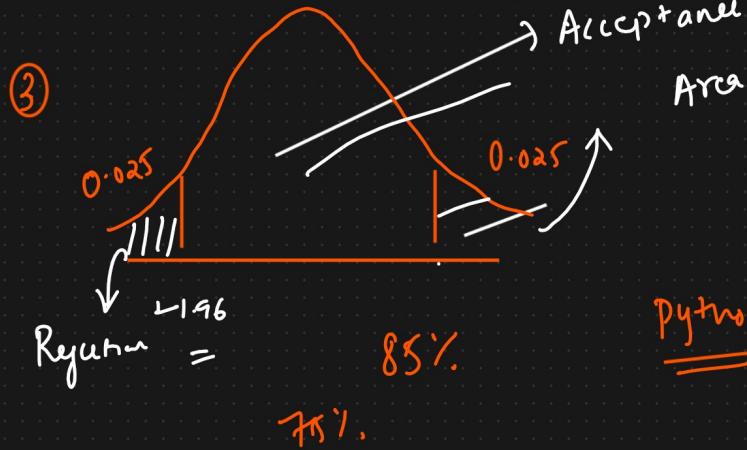
(*) A Complain was registered , the boys in the Municipal Primary School are underfed. Average weight of boys of age 10 is 32kgs with $S.D = 9\text{ kgs}$. A sample of 25 boys was selected from the municipal school and the average weight was found to be 29.5 kgs ? With $C.I = 95\%$ Check whether it is 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 $\alpha = 0.05$ $1 - 0.95 = 0.05$

$$H_1 = \mu < 32$$



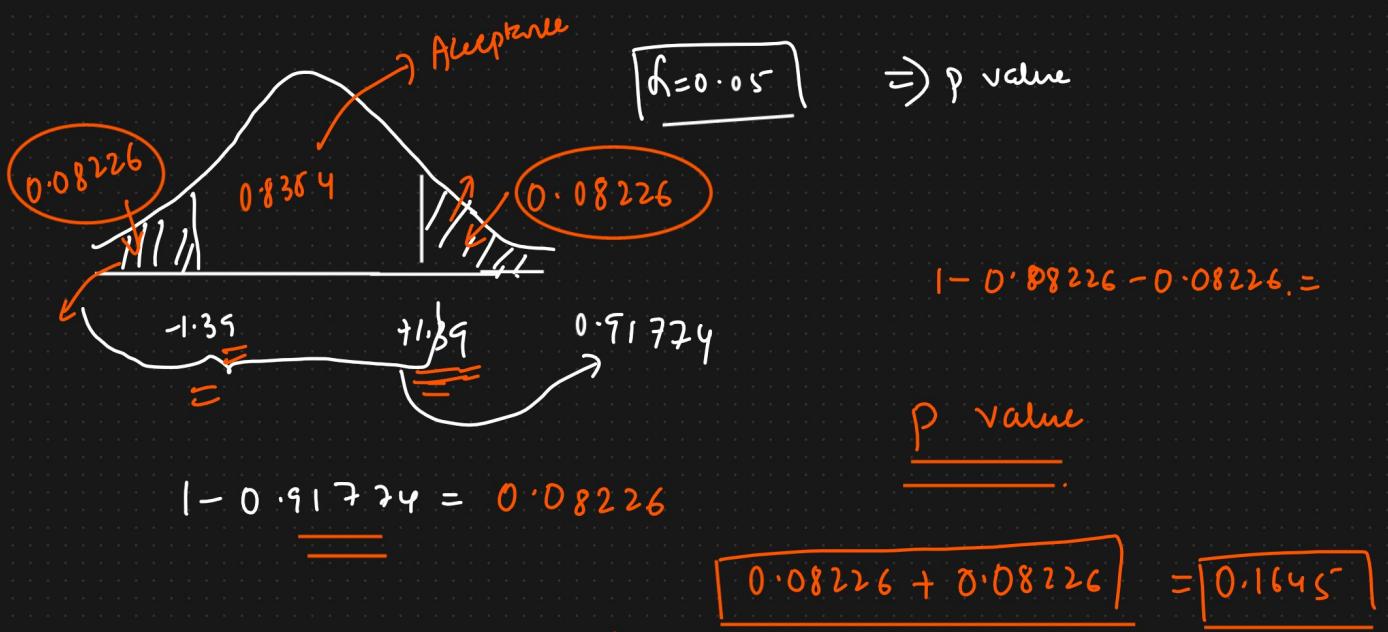
$$\textcircled{4} \quad Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

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

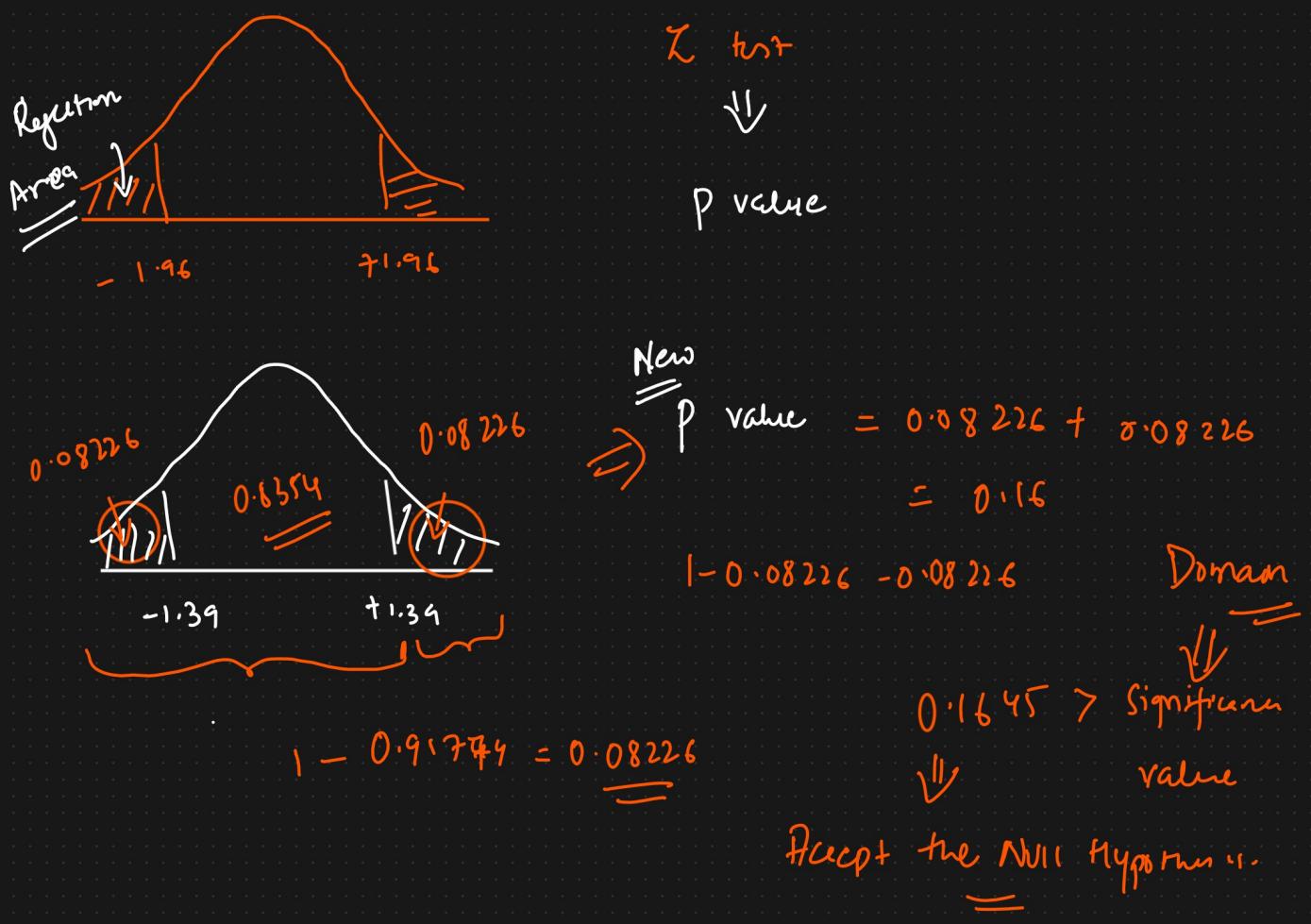
Conclusion : $-1.39 > -1.92$ therefore we accept the Null Hypothesis

So, the boys are not underfed.

So, the boys are not underfed.



Significance value
 $0.1645 > 0.05$
 \downarrow
 $p \geq 0.05 \Rightarrow \text{Accept the Null Hypothesis}$

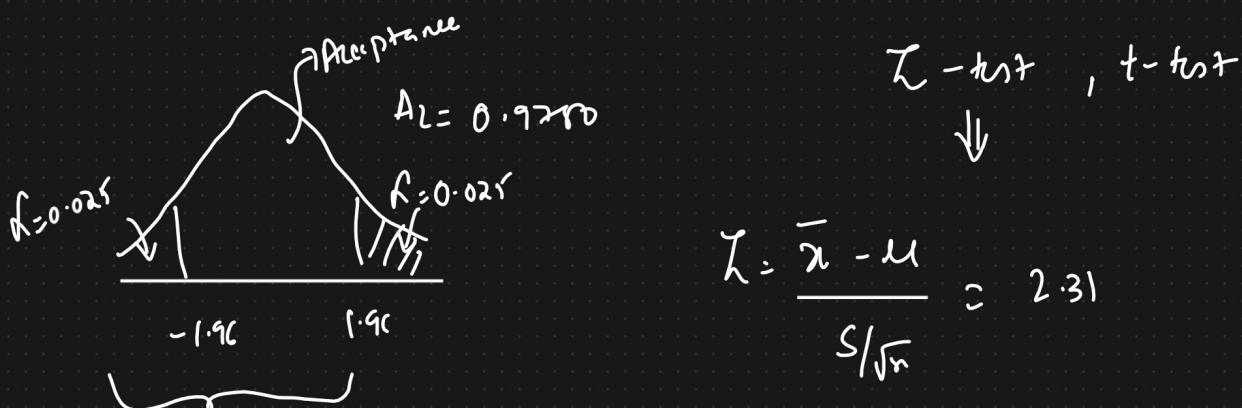


④ The average weight of all residents in town XYZ is 168 lbs. A nutritionist believes the true mean to be different. She measured the weight of 36 individuals and found the mean to be 169.5 lbs with a standard deviation of 3.9.

(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

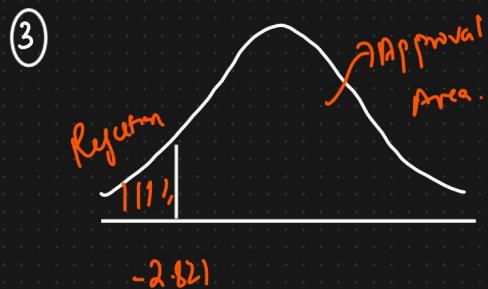
⑤ A company manufactures bike batteries with an average life span of 2 or more years. An engineer believes this value to be less. Using 10 samples, he measures the average life span to be 1.8 years with a standard deviation of 0.15.

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
 $H_1 : \mu < 2$ ≤ 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) $H_0: p_0 = 0.70.$

$H_1: p_0 \neq 0.70$

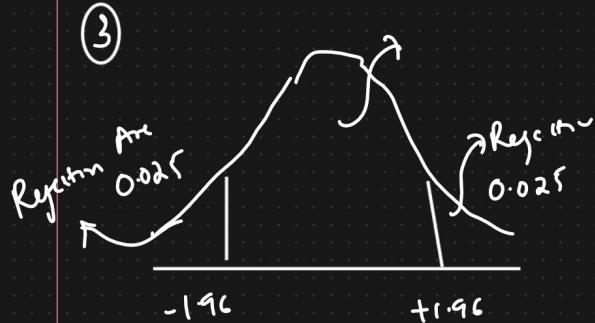
$$n = 200 \quad X = 130 \\ \hat{P} = \frac{X}{n} = \frac{130}{200} = \frac{13}{20} = 0.65$$

$$q_0 = 1 - p_0$$

② $\alpha = 0.05 \quad C.I = 95\%$

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

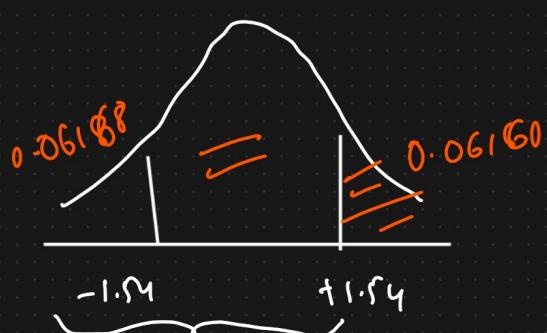
$$= \frac{0.65 - 0.70}{\sqrt{\frac{0.7 \times 0.3}{200}}} \approx -1.54$$



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

$x \uparrow \quad y \uparrow$

$\downarrow \quad \quad \quad \downarrow$
 $x =$

$y =$

{ quantity the relationship

between $x \& y$ }

$x \uparrow \quad y \downarrow$

-

-

$x \downarrow \quad y \uparrow$

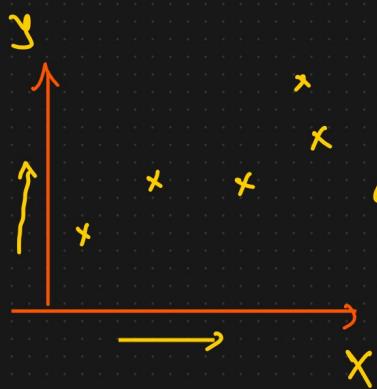
-

-

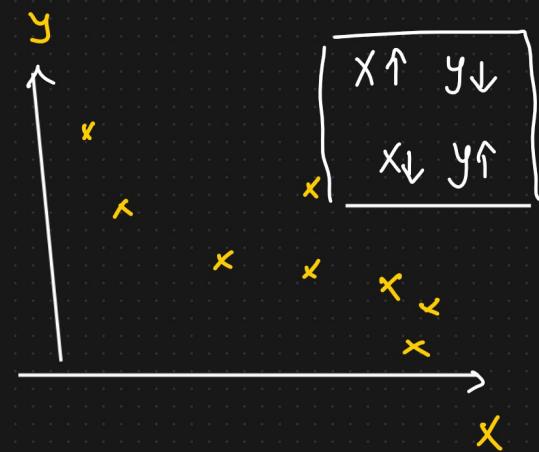
$x \downarrow \quad y \downarrow$

-

-



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



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

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

\Downarrow

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

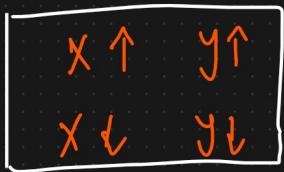
$$= \frac{\sum (x_i - \bar{x}) \times (x_i - \bar{x})}{N-1}$$

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

↓

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

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

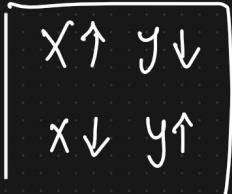


$$\text{Cov}(x, y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1} \quad \left. \right\}$$

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

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 $\left[\begin{matrix} -1 & \rightarrow & 1 \end{matrix} \right]$

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

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

Limit	-400
f 500	-300
-400	f 1000

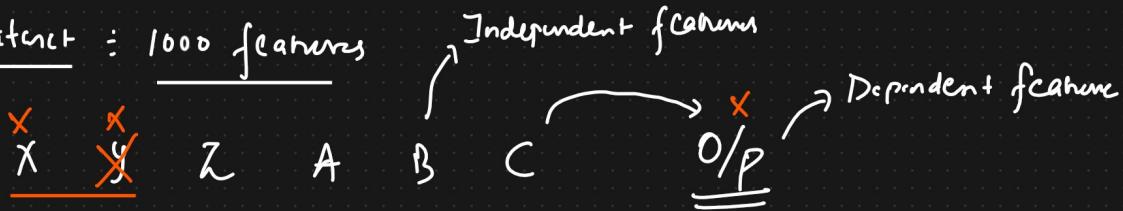
$\boxed{\infty}$

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 correlated

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

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

-ve correlation

↓
Keep it

③ Spearman Rank Correlation

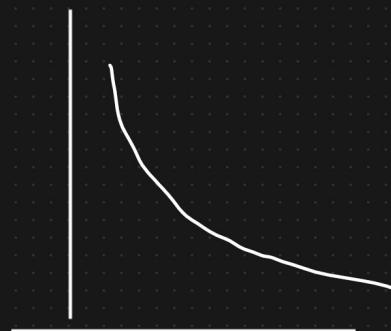
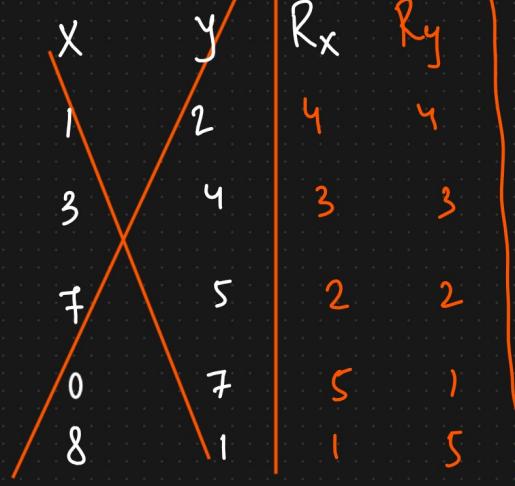
$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sqrt{R(x)} \sqrt{R(y)}}$$

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.

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

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

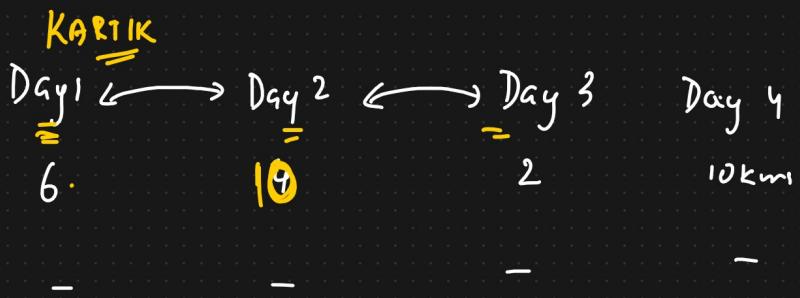


Types of ANOVA

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	Running Kms
Levles	1
	2

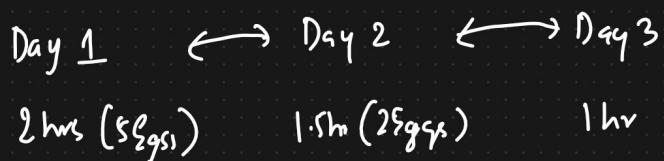


Ques. Study hours of KARTIK

④ Factorial ANOVA



Gym



⑥ Factorial ANOVA : Two or more factor (each of which with atleast 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

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

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

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

Statistics

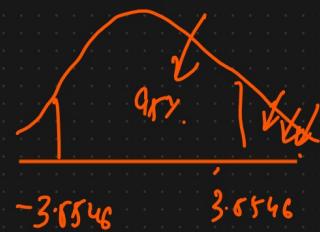
$$N = 21 \quad n = 7 \\ =$$

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

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		

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

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

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

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

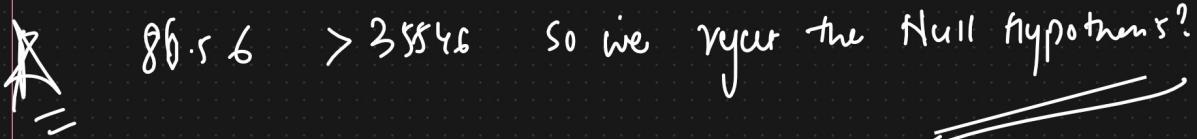
$$\left. \begin{array}{l} P=0.48 \\ d.f.=0.05 \end{array} \right\} = \sum y^2 - \left[\frac{57^2 + 47^2 + 21^2}{7} \right] = 10.29$$

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

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

Final Conclusion

Accept

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

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

Virginia

=

=

Pctz1 width

-

-

-

-

-

-

-

-

-

-

-

=

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

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

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

$$0.0118$$

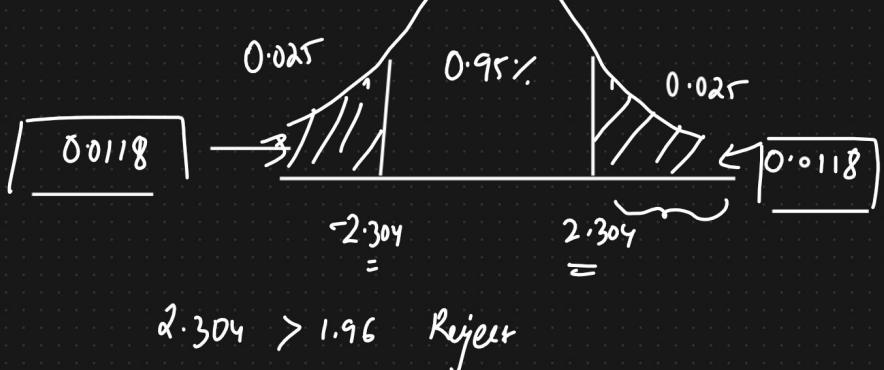
$$0.0228$$

$$0.0228$$

$$d =$$

$$Z_{\text{test statistic}}$$

Z_{test}



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

$$\sigma / \sqrt{n}$$

$$Z = 2.304$$

$2.304 > 1.96 \text{ Reject}$

Stats Interviews Questions

1. Question: What is the Central Limit Theorem and why is it important?

Answer: The Central Limit Theorem (CLT) states that the distribution of the sum (or average) of a large number of independent, identically distributed random variables approaches a normal (or Gaussian) distribution, regardless of the original distribution of the variables. It's crucial in statistics because it allows us to make inferences about populations using the normal distribution, which has well-understood properties.

2. Question: Explain Type I and Type II errors.

Answer:

- **Type I Error (False Positive, or Alpha error):** It's when you incorrectly reject a true null hypothesis.
- **Type II Error (False Negative, or Beta error):** It's when you fail to reject a false null hypothesis. The significance level (usually denoted by α) is the probability of making a Type I error. The power of a test is 1 minus the probability of making a Type II error (β).

3. Question: What is R-squared in linear regression?

Answer: R-squared, also known as the coefficient of determination, measures the proportion of the variance in the dependent variable that can be explained by the independent variables in a regression model. An R-squared value of 1 indicates that the regression predictions perfectly fit the data. Values of R-squared close to 1 indicate a large proportion of the variance in the dependent variable is explained by the regression model, while values close to 0 indicate the opposite.

4. Question: What is the difference between correlation and causation?

Answer: Correlation indicates a mutual relationship or association between two variables. When one variable changes, there's a tendency for the other variable to change in a specific direction. However, correlation does not imply causation. Causation means that a change in one variable is responsible for a change in another. For example, even if there's a strong correlation between ice cream sales and the number of drowning incidents, this doesn't mean buying more ice cream causes more drownings. A lurking variable, like temperature, can be influencing both.

5. Question: What is the difference between a parametric and a non-parametric test?

Answer: Parametric tests make assumptions about the parameters of the population distribution from which the sample is drawn, such as assuming that the population has a normal distribution. Examples include t-tests and ANOVA. Non-parametric tests, on the other hand, do not make strong assumptions about the population's distribution. Examples include the Mann-Whitney U test and Kruskal-Wallis test.

6. Question: Explain p-value.

Answer: The p-value is a measure used to help determine the significance of the results in hypothesis testing. It represents the probability of observing the current data, or something more extreme, given that the null hypothesis is true. A small p-value (typically ≤ 0.05) indicates strong evidence against the null hypothesis, so we reject it. A larger p-value suggests weaker evidence against the null hypothesis, so we fail to reject it.

7. Question: Describe the difference between cross-validation and bootstrapping.

Answer: Cross-validation is a technique for evaluating the performance of a statistical model by partitioning the original sample into a training set to train the model, and a test set to evaluate it. One common method is k-fold cross-validation. Bootstrapping, on the other hand, is a resampling technique used to estimate the distribution of a statistic (like the mean or variance) by sampling with replacement from the data. It helps assess the variability of a sample statistic and construct confidence intervals.

These are just a few potential questions. Depending on the role, interviewers might go deeper into specific topics or might also incorporate more practical, hands-on problems.

Some More Difficult Questions

1. Question: Can you explain the different measures of central tendency?

Answer: The three main measures of central tendency are the mean, median, and mode:

- **Mean:** It is the average of all the numbers in a dataset.
- **Median:** It is the middle value in a dataset when the numbers are arranged in order.
- **Mode:** It is the number that appears most frequently in a dataset.

2. Question: What is the difference between population and sample?

Answer: A population includes all members of a specified group, while a sample is a subset of the population. Statistics calculated on a population are called parameters, while those calculated on a sample are called statistics.

3. Question: How do you handle missing data?

Answer: Handling missing data can involve various techniques:

- **Deletion:** Remove records with missing values.
- **Imputation:** Fill missing values with estimated ones, e.g., using the mean, median, or mode of the known values, or using more complex algorithms or models to predict the missing value.
- **Analysis:** Use statistical techniques designed to handle missing values, such as multiple imputation or full information maximum likelihood estimation.

4. Question: What is the interquartile range (IQR) and why is it useful?

Answer: The IQR is a measure of statistical dispersion and is calculated as the difference between the upper (Q3) and lower (Q1) quartiles in a dataset. It is useful for understanding the spread of the data and for identifying outliers, as it is not affected by extremely large or small values.

5. Question: Explain the concept of skewness in statistics.

Answer: Skewness is a measure of the asymmetry of the probability distribution of a real-valued random variable. A negative skew indicates that the left tail of the distribution is longer, while a positive skew indicates that the right tail is longer. A skewness of zero indicates a perfectly symmetrical distribution.

6. Question: Can you describe what a box plot represents?

Answer: A box plot, or box-and-whisker plot, visually displays the distribution of a dataset, including its central tendency and variability. The box represents the interquartile range (IQR, Q3-Q1), the line inside the box shows the median, and the whiskers extend to the smallest and largest observations in the dataset.

7. Question: What is the difference between variance and standard deviation?

Answer: Variance and standard deviation are both measures of dispersion or spread in a dataset. Variance is the average of the squared differences from the mean, while the standard deviation is the square root of the variance. The standard deviation is more commonly used because it is in the same units as the data.

8. Question: What is a z-score and what is it used for?

Answer: A z-score is a statistical measurement that describes a value's relation to the mean of a group of values. It is measured in terms of standard deviations from the mean. A z-score is used to determine how unusual a value is, and it's commonly used for hypothesis testing, outlier detection, and comparison of scores from different datasets.

9. Question: Can you explain what covariance and correlation are?

Answer:

- **Covariance:** It is a measure of the joint variability of two random variables. A positive covariance indicates that the variables tend to increase and decrease together, whereas a negative covariance indicates that as one variable increases, the other tends to decrease.
- **Correlation:** It is the normalization of covariance to have values between -1 and 1, providing a measure of the strength and direction of the linear relationship between the two variables. A correlation of 1 indicates a perfect positive linear relationship, -1 indicates a perfect negative linear relationship, and 0 indicates no linear correlation.

These questions can help interviewers evaluate a candidate's understanding and knowledge of descriptive statistics concepts.

Lets Increase The Complexity

1. Question: How does the presence of outliers affect the mean and median of a dataset?

Answer: Outliers can greatly affect the mean because the mean considers all values in its calculation. An extreme outlier can pull the mean up or down, making it less representative of the central location of the data. The median, however, is more resistant to outliers since it depends only on the middle value(s) of an ordered dataset. In datasets with outliers, the median can often be a better representation of central tendency.

2. Question: Describe the concept of kurtosis. How is it different from skewness?

Answer: Kurtosis measures the "tailedness" of a probability distribution. High kurtosis indicates a distribution with tails heavier or more extreme than the normal distribution, and low kurtosis indicates a distribution with tails lighter than the normal distribution. While skewness deals with the asymmetry and direction of skew (left or right), kurtosis deals with the extremities (or outliers) in the distribution tails.

3. Question: How do you interpret the value of a Pearson correlation coefficient?

Answer: The Pearson correlation coefficient, often denoted as r , measures the strength and direction of a linear relationship between two variables. Its values range between -1 and 1.

- $r = 1$: Perfect positive linear relationship.
- $r = -1$: Perfect negative linear relationship.
- $r = 0$: No linear correlation. The closer r is to 1 or -1, the stronger the linear relationship. However, a strong correlation does not imply causation.

4. Question: Explain Simpson's Paradox and its implications in descriptive statistics.

Answer: Simpson's Paradox occurs when a trend or relationship between two variables reverses or disappears when they are examined in the context of a third variable. This can happen due to confounding factors. It emphasizes the importance of considering all relevant factors when interpreting statistical relationships.

5. Question: In a given dataset, what are the differences and relationships between the range, variance, and standard deviation?

Answer:

- **Range:** The difference between the maximum and minimum values in the dataset.
- **Variance:** The average of the squared differences from the mean.
- **Standard Deviation:** The square root of the variance.

The range provides a sense of the full spread of the data but is sensitive to outliers. The variance gives a measure of how data points differ from the mean, but it's in squared units of the data. Standard deviation, being the square root of variance, gives dispersion in the original units of the data and is commonly used because of this.

6. Question: How would you decide between using the mean vs. median as a measure of central tendency?

Answer: The decision often depends on the shape of the data distribution and the presence of outliers:

- For a symmetric distribution without outliers, the mean and median will be close, and either could be used.
- For skewed distributions or distributions with outliers, the median is usually a better representation because it is less affected by extreme values.

7. Question: Why might standard deviation be a misleading measure of spread in some situations?

Answer: Standard deviation can be misleading, especially when the data contains outliers, since it considers all deviations from the mean in its calculation. Extreme values can inflate the standard deviation, making it seem as though the data is more spread out than it actually is. In such cases, other measures like the interquartile range might be more appropriate.

Lets Try Some UseCases On Descriptive Stats

UseCase: A company sells products in three regions: North, South, and West. The sales team wants to understand the sales performance across these regions to allocate resources more efficiently.

Dataset:

Region	Monthly Sales (in thousands)
North	12, 15, 14, 13, 17, 19, 20
South	22, 21, 20, 23, 25, 26, 28
West	32, 30, 31, 29, 30, 33, 35

Question 1: Which region has the highest average monthly sales?

Process:

1. Calculate the mean (average) for each region.

$$\text{Formula for Mean: } \mu = \frac{\sum x}{n}$$

- μ = mean
- $\sum x$ = sum of all observations
- n = number of observations

2. Compare the means to determine the region with the highest average sales.

Answer:

- **North:** $\mu = \frac{12 + 15 + 14 + 13 + 17 + 19 + 20}{7} = 15.7$
- **South:** $\mu = \frac{22 + 21 + 20 + 23 + 25 + 26 + 28}{7} = 23.6$
- **West:** $\mu = \frac{32 + 30 + 31 + 29 + 30 + 33 + 35}{7} = 31.4$

The **West** region has the highest average monthly sales.

Question 2: Which region has the most consistent monthly sales (lowest variability)?

Process:

1. Calculate the standard deviation for each region to measure the spread of sales.

$$\text{Formula for Variance: } \sigma^2 = \frac{\sum (x - \mu)^2}{n}$$

$$\text{Formula for Standard Deviation: } \sigma = \sqrt{\sigma^2} = \sigma$$

- σ = standard deviation
- μ = mean
- \sum = sum of squared differences from the mean
- n = number of observations

2. Compare the standard deviations. The region with the lowest standard deviation is the most consistent.

Answer:

$$\bullet \text{ North Variance: } \sigma^2 = \frac{(12 - 15.7)^2 + \dots + (20 - 15.7)^2}{7} = 8.96$$

$$\text{North Standard Deviation: } \sigma = \sqrt{8.96} = 2.99 = 8.96$$

$$\bullet \text{ South Variance: } \sigma^2 = \frac{(22 - 23.6)^2 + \dots + (28 - 23.6)^2}{7} = 6.8$$

$$\text{South Standard Deviation: } \sigma = \sqrt{6.8} = 2.61 = 6.8$$

$$\bullet \text{ West Variance: } \sigma^2 = \frac{(32 - 31.4)^2 + \dots + (35 - 31.4)^2}{7} = 4.67$$

$$\text{West Standard Deviation: } \sigma = \sqrt{4.67} = 2.16 = 4.67$$

The **West** region has the most consistent monthly sales due to the lowest standard deviation.

Using these processes, you can provide insights into various data-driven questions by employing descriptive statistics. The core concept involves using measures of central tendency (like mean) and measures of spread (like standard deviation) to glean insights from datasets.

Lets Increase the Complexity On Usecases

Usecase: A hospital wants to analyze the recovery times of patients undergoing a specific surgery. The data for recovery times (in days) over a month is as follows:

Patient Group	Recovery Times (days)
A	5, 6, 4, 5, 7, 5, 6
B	7, 8, 7, 9, 8, 7, 9
C	5, 7, 6, 5, 6, 6, 5

Question 1: Which patient group has the quickest median recovery time?

Process:

1. Sort the recovery times for each group in ascending order.
2. Find the median (middle value) for each group.
3. Compare the medians.

Answer:

- **Group A Median:** The middle value of the sorted list (4, 5, 5, 5, 6, 6, 7) is 5.
- **Group B Median:** The middle value of the sorted list (7, 7, 7, 8, 8, 9, 9) is 8.
- **Group C Median:** The middle value of the sorted list (5, 5, 5, 6, 6, 6, 7) is 6.

Patient Group A has the quickest median recovery time of 5 days.

Question 2: Which patient group has the least variation in recovery times?

Process:

1. Calculate the range (difference between maximum and minimum values) for each group.
2. The group with the smallest range has the least variation.

Answer:

- **Group A Range:** $7 - 4 = 3$
- **Group B Range:** $9 - 7 = 2$
- **Group C Range:** $7 - 5 = 2$

Patient Groups B and C both have the least variation in recovery times with a range of 2 days.

Question 3: How do the interquartile ranges (IQR) of the groups compare?

Process:

1. Calculate the first quartile (Q1) and third quartile (Q3) for each group.
2. Subtract Q1 from Q3 to get the IQR for each group.
3. Compare the IQRs.

Answer:

- **Group A IQR:** For the sorted list (4, 5, 5, 5, 6, 6, 7), Q1 = 5 and Q3 = 6. $IQR = 6 - 5 = 1$.
- **Group B IQR:** For the sorted list (7, 7, 7, 8, 8, 9, 9), Q1 = 7 and Q3 = 9. $IQR = 9 - 7 = 2$.
- **Group C IQR:** For the sorted list (5, 5, 5, 6, 6, 7), Q1 = 5 and Q3 = 6. $IQR = 6 - 5 = 1$.

Patient Groups A and C have the same IQR of 1 day, which is less than Group B's IQR.

These use cases illustrate how to utilize various descriptive statistics measures to analyze and interpret real-world data. By understanding the distributions, central tendencies, and variations of datasets, decisions can be more data-driven and informed.

Use cases on Different type Of Distributions

Use case: An e-commerce company analyzes its website's page load times in seconds over a month to optimize user experience. The data includes:

Day	Load Times (seconds)
1	3, 2.5, 2.8, 3.1, 15 (Outlier due to a server glitch)

2	2.6, 2.5, 2.7, 2.9, 2.8
3	2.7, 2.8, 2.6, 2.5, 3
...	...

Question 1: What impact do outliers have on the average load time?

Process:

1. Calculate the mean with and without outliers.
2. Compare both means to gauge the effect of outliers.

Answer:

With the outlier: Mean = $(3 + 2.5 + 2.8 + 3.1 + 15) / 5 = 5.28$

Without the outlier: Mean = $(3 + 2.5 + 2.8 + 3.1) / 4 = 2.85$

The outlier significantly increases the average page load time by 2.43 seconds.

Question 2: How can we transform load times to normalize the data?

Process:

1. Use logarithmic transformation.
2. Compute the logarithm (base 10 or natural logarithm) of all page load times.

Answer: Log-transforming the data can help in dealing with skewed data or data with outliers. If the original load time was 3 seconds, the transformed value using a natural log would be $\ln(3) \approx 1.0986$.

Question 3: Describe the distribution of load times using histograms.

Process:

1. Divide the data into bins (e.g., 2-2.5 seconds, 2.5-3 seconds).
2. Count the number of observations within each bin.
3. Plot the frequency of observations vs. bins.

Answer: Using the histogram, you might find, for instance, that most page load times cluster around 2.5-3 seconds, indicating the mode of the distribution. Peaks would represent common load times, while troughs would show less frequent load times.

Question 4: What is the Probability Density Function (PDF) for day 2's load times?

Process:

1. Estimate the PDF from the data (often using kernel density estimation).
2. Plot the continuous curve, showing how densities of load times vary.

Answer: The PDF will be a continuous curve indicating the probability of the page taking a specific time to load. For instance, the peak around 2.7 seconds might have a higher value, indicating it's the most common load time for day 2.

Question 5: What is the Probability Mass Function (PMF) for load times on day 3?

Process:

1. For discrete data, compute the proportion of each unique load time.
2. Plot these proportions.

Answer: The PMF might show, for instance, that the probability of the page taking exactly 2.7 seconds to load is 0.2 (or 20%). It gives probabilities for discrete outcomes.

These analyses can be deepened using more data and more advanced statistical methods, but the use case provides an insight into how different techniques in descriptive statistics can be used in a practical scenario.

Try this By your Own

Use Case 1: A pharmaceutical company has developed a new drug. During clinical trials, they measured the time (in hours) it took for patients to show symptom relief. They're particularly interested in how quickly the drug works.

Dataset Sample:

sql	Patient Number	Relief Time (hours)
	1	3.5
	2	2.8
	3	4.1

Question 1: Do the relief times follow a normal distribution?

Process:

1. Plot a histogram of the relief times.
2. Overlay a normal distribution curve on the histogram.

Answer: If the histogram matches closely with the normal distribution curve, then the relief times likely follow a normal distribution.

Question 2: What percentage of patients experienced relief within 3 hours, assuming the data follows a normal distribution?

Process:

1. Calculate the z-score for 3 hours: $z = \frac{x - \mu}{\sigma}$
2. Look up this z-score in a z-table to find the percentage of patients.

Answer: If the z-score is, for example, -0.5 and corresponds to 30% on the z-table, then 30% of patients experienced relief within 3 hours.

UseCase 2: A factory produces light bulbs. They have a dataset of the number of bulbs produced each day and the percentage of defective bulbs. They want to improve the quality control process.

Dataset Sample:

sql	Day	Defective Bulbs (%)
	1	2
	2	1.5
	3	3

Question 1: Do the percentages of defective bulbs follow a Poisson distribution?

Process:

1. If the occurrence of defects is rare and random, the distribution might follow a Poisson distribution.
2. Plot the PMF of the observed defects and compare with the PMF of a Poisson distribution with the same mean.

Answer: If the observed PMF aligns closely with the Poisson PMF, it's likely that the defect rates follow a Poisson distribution.

Question 2: If the data follows a binomial distribution, what is the probability that more than 5% of the bulbs are defective on any given day?

Process:

1. Use the binomial probability formula: $P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$ Where:
 - n is the total number of trials (bulbs produced)
 - k is the number of successes (defective bulbs)
 - p is the probability of success on a single trial.
2. Calculate the probability for 5%, 6%, 7%,... and sum these probabilities.

Answer: The sum of the probabilities gives the likelihood that more than 5% of the bulbs are defective on any given day.

These use cases illustrate the application of different types of distributions (normal, Poisson, binomial) in real-world scenarios. In practice, determining the fit of a distribution would require more rigorous statistical testing, but this gives an overview of the process.

1. Z-test

Question: A national examination board believes that the students in state X score an average of 52 in mathematics. A state education official disputes this and collects a random sample of 100 student scores from the state. The sample has an average score of 54 with a standard deviation of 10. At the 0.05 significance level, is the official correct?

Solution:

- **Null Hypothesis (H₀):** The students in state X have an average score of 52.
- **Alternative Hypothesis (H_a):** The students in state X do not have an average score of 52.

Step by Step Process:

1. Calculate the z-score: $z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$

Where:
 \bar{X} = sample mean = 54
 μ = population mean = 52
 σ = sample standard deviation = 10
 n = number of samples = 100

2. Compare the z-score to the critical z-value for a 0.05 significance level (two-tailed).

3. If $|z| > z\text{-critical}$, reject the null hypothesis.

```
python
import math
import scipy.stats as stats

X_bar = 54
mu = 52
sigma = 10
n = 100

z = (X_bar - mu) / (sigma/math.sqrt(n))
p = 1 - stats.norm.cdf(abs(z))

alpha = 0.05
if p < alpha:
    print("Reject the null hypothesis")
else:
    print("Do not reject the null hypothesis")
```

2. T-test

Question: A company claims its new energy drink increases stamina. 15 people were tested before and after consuming the drink. Test if the drink has a significant effect on stamina at the 0.05 significance level.

Solution:

Given that the measurements are paired (before and after for the same individual), use a paired t-test.

Step by Step Process:

1. Compute the difference in stamina for each individual.
2. Compute the mean and standard deviation of these differences.
3. Calculate the t-statistic.
4. Compare the t-statistic to the critical t-value for a 0.05 significance level.

```
python
import numpy as np

before = np.array([...]) # insert stamina values before drinking
after = np.array([...]) # insert stamina values after drinking

differences = after - before
t_stat, p_value = stats.ttest_rel(after, before)

alpha = 0.05
if p_value < alpha:
    print("Reject the null hypothesis")
else:
    print("Do not reject the null hypothesis")
```

3. ANOVA

Question: A farmer tests three types of fertilizers to see which one produces the highest crop yield. Is there a significant difference in yield across the fertilizers?

Solution:**Step by Step Process:**

1. Use one-way ANOVA to compare the means of crop yields from the three fertilizers.
2. If the p-value is below the significance level, there is a significant difference.

```
python
fertilizerA = np.array([...]) # insert yields for fertilizer A
fertilizerB = np.array([...]) # insert yields for fertilizer B
fertilizerC = np.array([...]) # insert yields for fertilizer C

f_stat, p_value = stats.f_oneway(fertilizerA, fertilizerB, fertilizerC)

alpha = 0.05
if p_value < alpha:
    print("Reject the null hypothesis")
else:
    print("Do not reject the null hypothesis")
```

4. Chi-Square Test

Question: A company wants to know if there's a relationship between gender (male, female) and product preference (Product A, Product B). They survey 100 customers. Is product preference independent of gender?

Solution:**Step by Step Process:**

1. Construct a contingency table of gender vs. product preference.
2. Compute the chi-square statistic and p-value.
3. If the p-value is below the significance level, they are not independent.

```
python
# Contingency table: rows = gender, columns = product preference
observed = np.array([[30, 20], # males
                     [25, 25]]) # females

chi2_stat, p_value, _, _ = stats.chi2_contingency(observed)

alpha = 0.05
if p_value < alpha:
    print("Reject the null hypothesis")
else:
    print("Do not reject the null hypothesis")
```

5. Regression

Question: An e-commerce website wants to understand if the time spent on the website (in minutes) predicts the total amount spent (in dollars). They gather data from 100 users. Determine if there's a relationship.

Solution:**Step by Step Process:**

1. Run a simple linear regression with time spent as the independent variable and amount spent as the dependent variable.
2. If the p-value for the slope is below the significance level, there's a significant relationship.

```
python
from statsmodels import api as sm

time_spent = np.array([...]) # insert time spent by users
amount_spent = np.array([...]) # insert amount spent by users

X = sm.add_constant(time_spent) # adding a constant
model = sm.OLS(amount_spent, X).fit()

alpha = 0.05
if model.pvalues[1] < alpha:
    print("Reject the null hypothesis")
else:
    print("Do not reject the null hypothesis")
```

Remember, these are simplifications of what you might encounter in the real world. In practice, you would also check assumptions, consider additional covariates, and potentially apply more sophisticated statistical models.

This usecase definitely you should solve

Scenario: Online Retail Store A/B Testing

An online retail store has recently introduced a new webpage design to increase the amount of time users spend on the page and ultimately increase purchases. They have conducted A/B testing, where Group A is exposed to the old design, and Group B to the new design. They've collected data on the time spent on the webpage and whether a purchase was made.

Objective: Determine if the new webpage design leads to a significant increase in both time spent on the webpage and the likelihood of making a purchase.

Steps:

1. Define the Problem:

- **Null Hypothesis (H0):** The new webpage design does not significantly affect the time spent on the webpage and the likelihood of making a purchase.
- **Alternative Hypothesis (HA):** The new webpage design significantly affects the time spent on the webpage and the likelihood of making a purchase.

2. Data Collection:

- Collect data on time spent on the webpage and purchasing behavior for both groups.

3. Data Exploration and Preprocessing:

- Understand the basic statistics of the datasets.
- Handle missing values if any.
- Check and handle outliers.

4. Perform T-Test on Time Spent:

- Conduct an Independent Samples t-test to compare the mean time spent on the webpage by the two groups.

5. Perform Chi-Square Test on Purchase Behavior:

- Construct a contingency table of the groups and purchasing behavior.
- Conduct a Chi-Square test to check the independence of the group and purchasing behavior.

6. Decision Making:

- Based on the p-values from the t-test and Chi-Square test, reject or fail to reject the null hypothesis.
 - Make recommendations for the business.
-

Python Code:

```

python
import numpy as np
import pandas as pd
import scipy.stats as stats

# Suppose `data` is the collected data with 'group', 'time_spent', and 'purchase' columns
# 'group' - 'A' for control group and 'B' for test group
# 'time_spent' - time spent by the user on the webpage
# 'purchase' - 1 if the user made a purchase, 0 otherwise

# Sample data creation
data = pd.DataFrame({
    'group': ['A', 'A', 'B', 'A', 'B', 'A', 'B'],
    'time_spent': [3, 5, 7, 4, 6, 8, 2, 7],
    'purchase': [0, 1, 1, 0, 1, 1, 0, 1]
})

# Step 4: Perform T-Test on Time Spent
group_A_time_spent = data[data['group'] == 'A']['time_spent']
group_B_time_spent = data[data['group'] == 'B']['time_spent']

```

```
t_stat, p_value_time = stats.ttest_ind(group_A_time_spent, group_B_time_spent)

# Step 5: Perform Chi-Square Test on Purchase Behavior
contingency_table = pd.crosstab(data['group'], data['purchase'])
chi2_stat, p_value_purchase, _, _ = stats.chi2_contingency(contingency_table)

# Step 6: Decision Making
alpha = 0.05

if p_value_time < alpha:
    print("There is a significant difference in time spent on the webpage between the two groups.")
else:
    print("There is no significant difference in time spent on the webpage between the two groups.")

if p_value_purchase < alpha:
    print("There is a significant difference in purchasing behavior between the two groups.")
else:
    print("There is no significant difference in purchasing behavior between the two groups.")
```

In this script:

- Data for two groups (A and B) is analyzed.
- A T-test is conducted to compare the mean time spent on the webpage for the two groups.
- A Chi-Square test is conducted to assess the association between the group and purchasing behavior.
- Decisions are made based on the p-values from the tests to determine whether there are significant differences in time spent on the webpage and purchasing behavior between the two groups.

Recommendation:

- If there is a significant increase in both time spent and likelihood of making a purchase, recommend the implementation of the new webpage design.
- If not, further analysis and perhaps more A/B testing may be needed to identify effective strategies for improving webpage performance and sales.

Lets Solve with each step by step

Scenario: Online Retail Store A/B Testing

Hypothetical Data:

Group A (Old Design):

- Time spent (minutes): [3, 5, 4, 6, 5, 5, 6, 4]
- Purchases: [0, 1, 0, 1, 1, 0, 0, 1]

Group B (New Design):

- Time spent (minutes): [6, 7, 7, 7, 8, 6, 7, 8]
- Purchases: [1, 1, 1, 1, 1, 1, 0, 1]

Step-by-Step Solution with Numerical Calculation:

1. Define the Problem:

As stated in the previous answer.

2. Data Collection:

Data has been hypothetically given.

3. Data Exploration and Preprocessing:

- Calculate the means and standard deviations for both groups:

Group A Time spent:

- Mean = $(3 + 5 + 4 + 6 + 5 + 5 + 6 + 4) / 8 = 4.75$ minutes
- Standard Deviation ≈ 1.16 minutes

Group B Time spent:

- Mean = $(6 + 7 + 7 + 7 + 8 + 6 + 7 + 8) / 8 = 7$ minutes
- Standard Deviation ≈ 0.76 minutes

4. Perform T-Test on Time Spent:

- Compute the t-statistic using the formula: $t = \frac{X_1 - X_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{n_1 s_1^2 + n_2 s_2^2}{\sqrt{X_1 - X_2}}$
-

Where:

- $X_1 X_2$ are the sample means of group A and B.
- $s_1 s_2$ are the sample standard deviations of group A and B.
- $n_1 n_2$ are the sample sizes of group A and B.

Using the values: $t = \frac{4.75 - 7}{\sqrt{\frac{1.16^2}{8} + \frac{0.76^2}{8}}} = 81.162 + 80.762$

$$\sqrt{4.75 - 7} \approx -5.91 \quad t \approx -5.91$$

Consulting a t-distribution table for $df=14$ (since $df = n_1 + n_2 - 2$) and $\alpha=0.05$ (two-tailed), the critical t-value is approximately ± 2.145 .

Since $-5.91 < -2.145$, we reject the null hypothesis for time spent.

5. Perform Chi-Square Test on Purchase Behavior:

Construct a 2x2 contingency table:

	Purchase=0	Purchase=1	Total
Group A	4	8	
Group B	7	8	

- Calculate the expected frequencies for each cell:

For example, for Group A and Purchase=0: Expected frequency = $(\text{row total} * \text{column total}) / \text{grand total} = (8 * 5) / 16 = 2.5$

- Compute the chi-square statistic using the formula:

$$\chi^2 = \sum \frac{(observed - expected)^2}{expected} \quad \chi^2 = \sum \text{expected}(observed - expected)^2$$

Using the observed and expected frequencies: $\chi^2 \approx 3.6$

For a 2x2 table with $\alpha=0.05$, the critical value from the chi-square distribution is approximately 3.841.

Since $3.6 < 3.841$, we fail to reject the null hypothesis for purchase behavior.

6. Decision Making:

- The t-test indicates a significant difference in time spent on the webpage between the two designs.
- The chi-square test suggests that purchase behavior isn't significantly different between the groups, although it's close to the threshold.

Recommendation:

- The new webpage design seems to engage users for longer periods.
- While purchase behavior hasn't shown a statistically significant change, it's close, and with a larger sample, it might.

Further A/B testing or possibly combining this new design with other strategies could be beneficial for sales.