

STATS

Statistics and its applications -

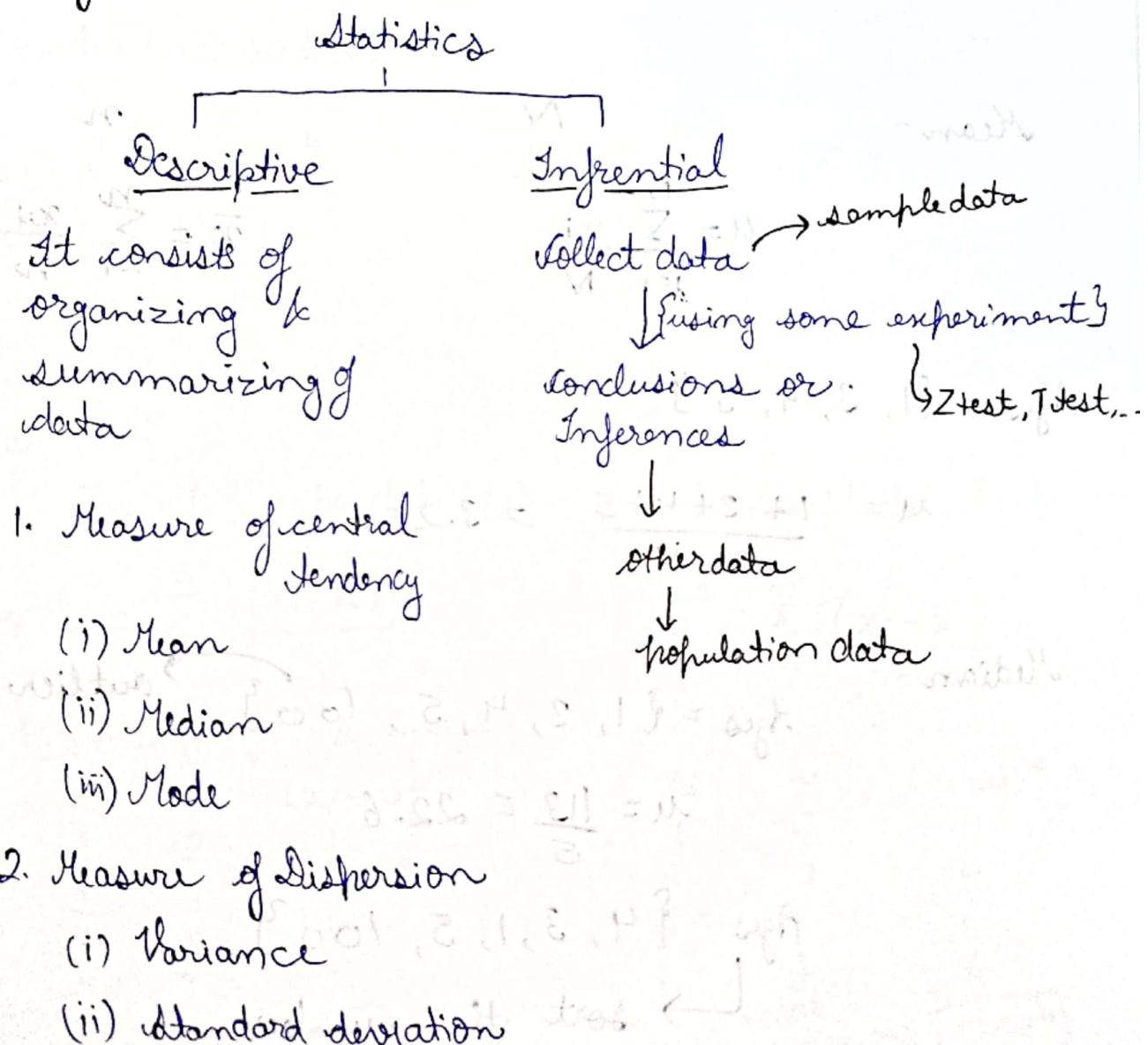
Statistics is a field that deals with collection, organization, analysis, interpretation and presentation of the data.

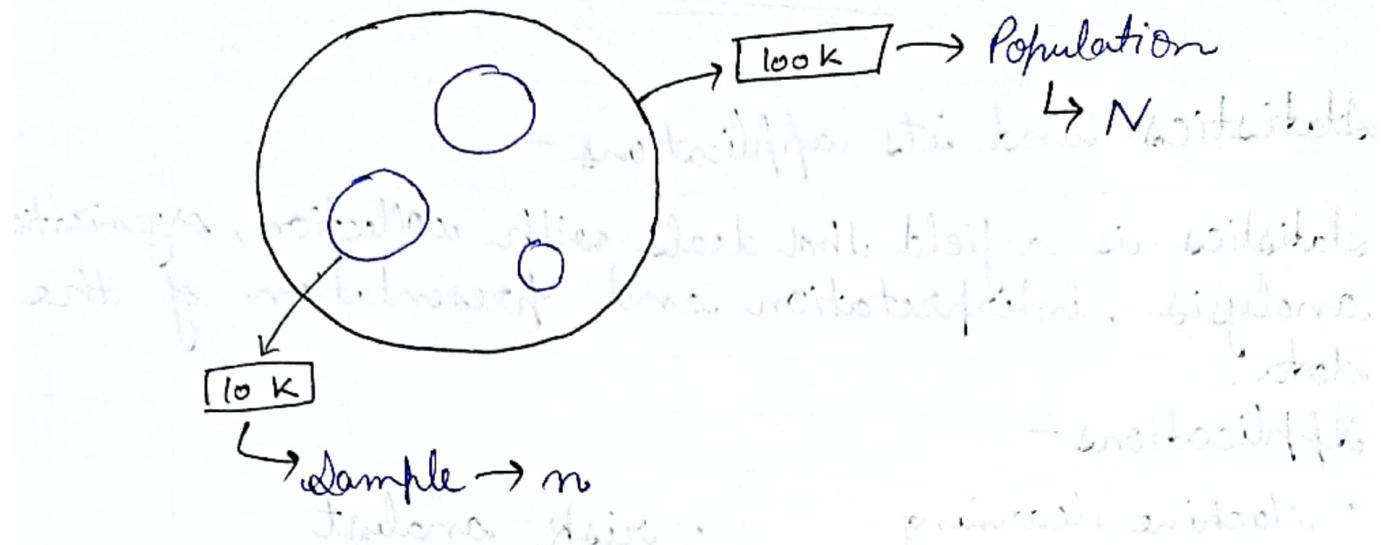
Applications -

- Machine learning
- Data analyst
- BA

- risk analyst
- house wife
- Each & every domain

Types of statistics -





Measure of central tendency -

- ① Mean
- ② Median
- ③ Mode

Mean -

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

(Sum of all values / no. of values)

Ages = {1, 3, 4, 5}

$$\mu = \frac{1+3+4+5}{4} = 3.25$$

Median -

Ages = {1, 3, 4, 5, 100} → outlier

$$\mu = \frac{1+3+4+5}{5} = 22.6$$

Ages = {4, 3, 1, 5, 100}

→ sort the numbers

$$\text{Ages} = \{1, 3, 4, 5, 100\}$$

* Instead of finding total average pick up the central element.

Median - 4

$$\text{Ages} = \{1, 3, \underline{4}, 5, 100, 200\}$$

$$\hookrightarrow \text{Medians} \rightarrow \frac{4+5}{2} = 4.5$$

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

* Select the element having maximum frequency.

$$\text{Mode} = 4$$

Measure of Dispersion

① Variance

② Standard deviation

$$\text{Ages 1} = \{2, 2, 4, 4\}$$

$$\mu = \frac{2+2+4+4}{4} \\ = 3$$

$$\text{Ages 2} = \{1, 1, 5, 5\}$$

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

Variance -

For population data $\{N\}$

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

For sample data $\{n\}$

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

Eg - Using Ages 1 & Ages 2

Index	x_i	μ	$(x_i - \mu)^2$
	2	3	1
	2	3	1
	4	3	1
	4	3	1

$$\boxed{\sum (x_i - \mu)^2}$$

x_i	μ	$(x_i - \mu)^2$
1	3	4
1	3	4
5	3	4
5	3	4

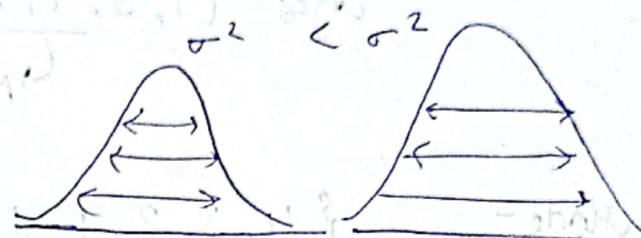
Scanned with CamScanner

Standard deviation - measure of spread

Variance ↑↑ spread ↑↑

$\sqrt{\sigma^2}$ → how far a data pt. is away from mean

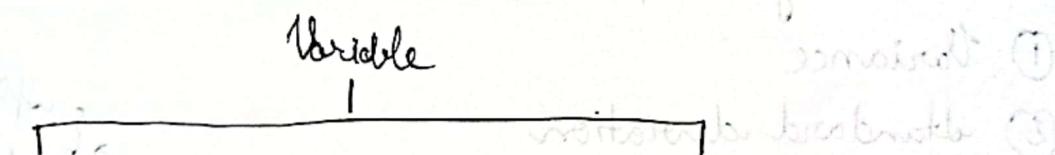
Population standard deviation



Variable - a variable is a property that can take up any value.

Age = 25 → variable

Ages = [11, 20, 25, 30] → not a variable



Quantitative variable

Discrete

(whole no.)

Eg - Age

Continuous

(can be decimal)

Eg - Weight

Qualitative

Categorical

(e.g. - Gender [male, female], colors)

[not is numerical values]

Random variables -

$y = 5x + 2$ → here x is a random variable

$y = 1$, $y = 2$, $y = 3$ etc. which can take diff values.

X is a function

values

are derived from process k experiments

Tossing a coin

$$X = \begin{cases} 0 & - \text{ if head} \\ 1 & - \text{ if tail} \end{cases}$$

Random variables

Discrete

Eg - tossing a coin
rolling a dice

Continuous

Eg - tomorrow how many inches it is going to rain

Percentiles & Quartiles -

① Percentage -

How many % of numbers are odd - {2, 4, 7, 5, 8, 9, 10}

$$\text{Percentage} = \frac{\# \text{ No.'s are odd}}{\# \text{ of obs.}} \times 100 = \frac{3}{7} \times 100$$

② Percentile - A percentile is a value below which a certain percentage of obs. lie.

Eg - 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

* What is the percentile of ranking 10%?

$$\begin{aligned} \text{Percentile rank of } x &= \frac{\# \text{ of values below } x}{n} \times 100 \\ &= \frac{16 \times 100}{20} = 80\% \end{aligned}$$

* What value exists at the percentile ranking of 25% ??

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

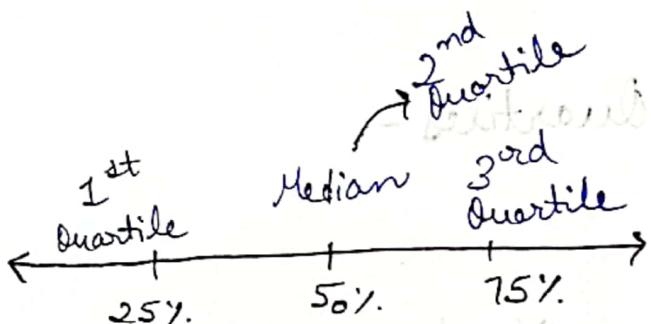
$$= \frac{25}{100} \times (20+1) = 5.25$$

5.25 index indication -

$$\text{Avg}(5+6) = \frac{5+5}{2} = 5 \rightarrow 25\%$$

↓
indexes

③ Quartiles -



Five number summary -

1. Minimum

2. First Quartile

3. Median

4. Third Quartile

5. Maximum

Dataset -

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27

{Remove the outlier}

[Lower fence \rightarrow Higher fence]

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

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

$$\text{Inter Quartile range (IQR)} = Q_3 - Q_1 = 7 - 3 = 4$$

$$Q_1(25\%) = \frac{\text{percentile}}{100} \times (n+1) = \frac{25}{100}(20) \\ = 5^{\text{th}} \text{ index} = 3$$

$$Q_3(75\%) = \frac{75}{100} \times 20 = 15^{\text{th}} \text{ index} = 7$$

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

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

$\therefore 27$ is an outlier

① Minimum = 1

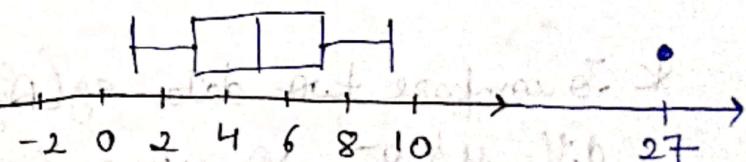
② Q₁ = 3

③ Median = 5

④ Q₃ = 7

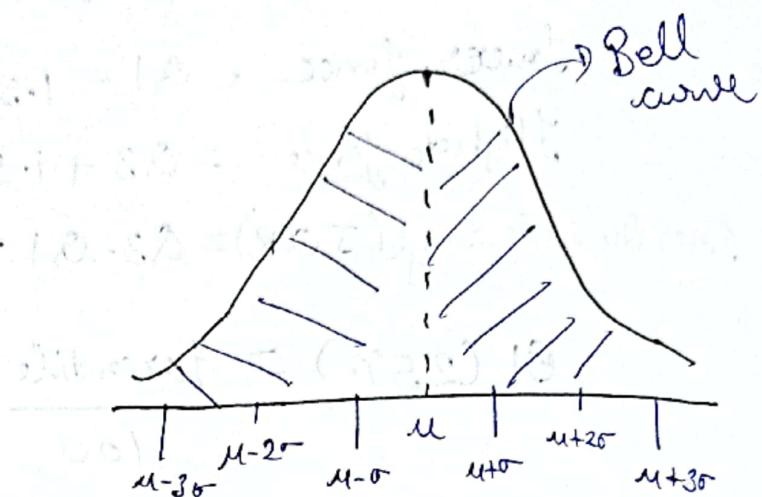
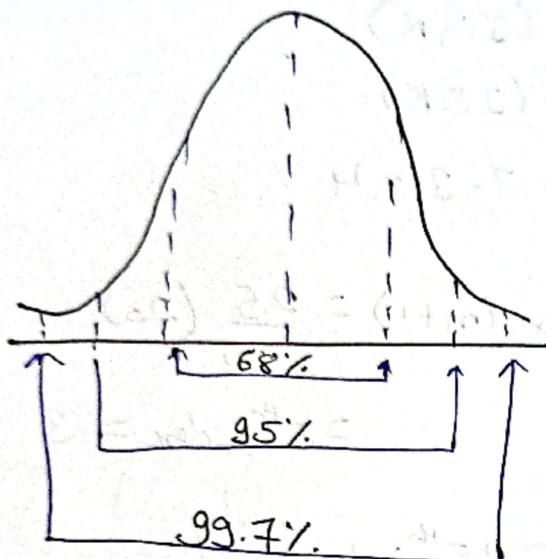
⑤ Maximum = 9

Box plot



Normal Distribution / Gaussian distribution

Empirical formula
[68 - 95 - 99.7]



Standard normal distribution

$$X \approx N(\mu, \sigma)$$

$$\Downarrow$$

$$Y \approx SND(\mu=0, \sigma=1)$$

Z-score

$$Z = \frac{X_i - \mu}{\sigma}$$

Eg - 1, 2, 3, 4, 5

$$\mu = 3, \sigma = 1$$

$$Z = \frac{x_i - \mu}{\sigma} = \{-2, -1, 0, 1, 2\}$$

$$\hookrightarrow SND(\mu=0, \sigma=1)$$

* To compare two data eg (Age, Weight) both will have diff. units so using Z-score we find their SND.

Application of Z-score -

Z-score → these are standardized values that can be used to compare scores in diff. distribution.

Eg → Cricket - 2020 2021

India vs Australia test match

2020 (Stats - 3 test)

Test average = 181

standard avg. = 12

Rishabh final score = 187

2021 (stats)

Rishabh test avg. = 182

standard deviation = 5

final score = 185

Compared to the rest of the test scores, in which year was Rishabh's score in final game

2020 1 + better? $2.020 = 0.79$

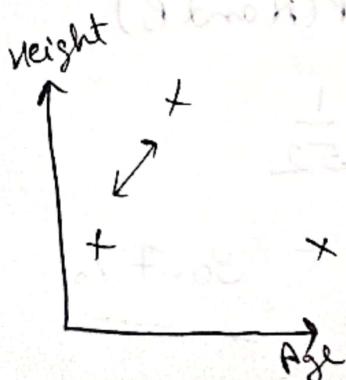
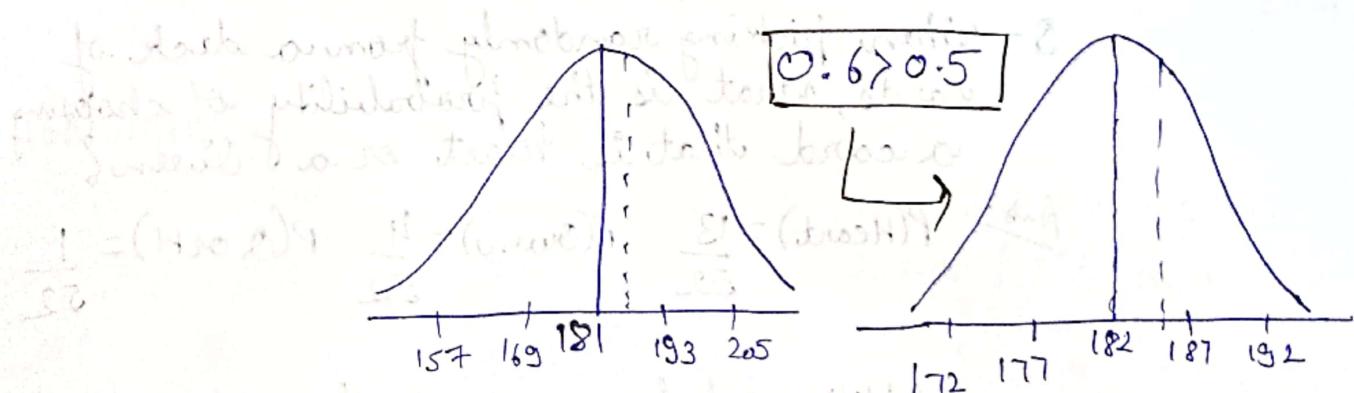
$$\text{Ans: } Z = \frac{x-\mu}{\sigma}$$

$$= \frac{187-181}{12} = 0.5$$

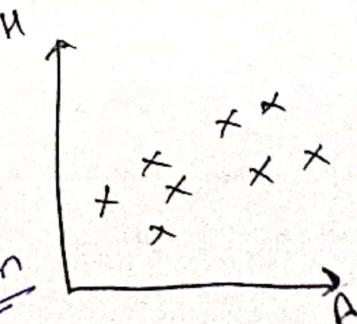
$$Z = \frac{185-182}{5}$$

$$= 0.6$$

→ In 2020 Rishabh's performance will be 2021



$S.D.$
z-score
(same scale)
i.e. Standardization



Log normal distribution -

$$X \approx N(\mu, \sigma^2)$$



$X \approx \text{log normal distribution}$



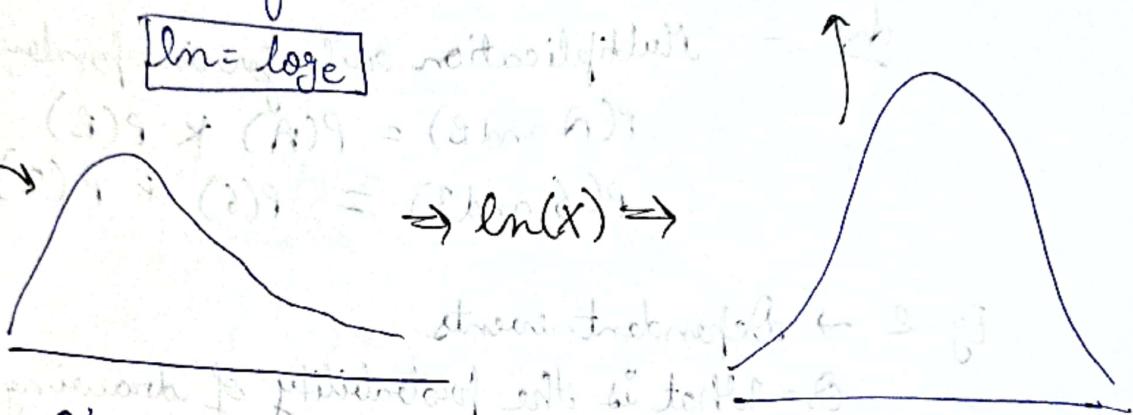
Then $y = \ln(x)$ have a normal distribution

$\ln = \log_e$

$$(1)^n + (1)^n = (2^n)^n$$

$$X =$$

$$\rightarrow \ln(x) \leftarrow$$



So $X \approx \text{log normal distribution}$

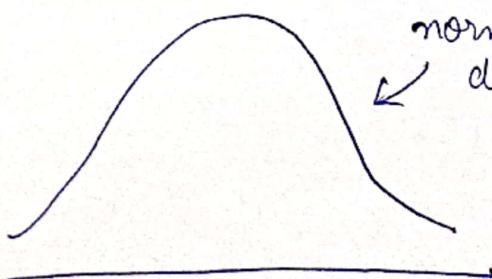
Central limit theorem -

The CLT states that regardless of the shape of population distributions the distribution of sample mean will be approximately normal.

$$\bar{S}_1, \bar{S}_2, \bar{S}_3, \dots, \bar{S}_n$$

All these will be

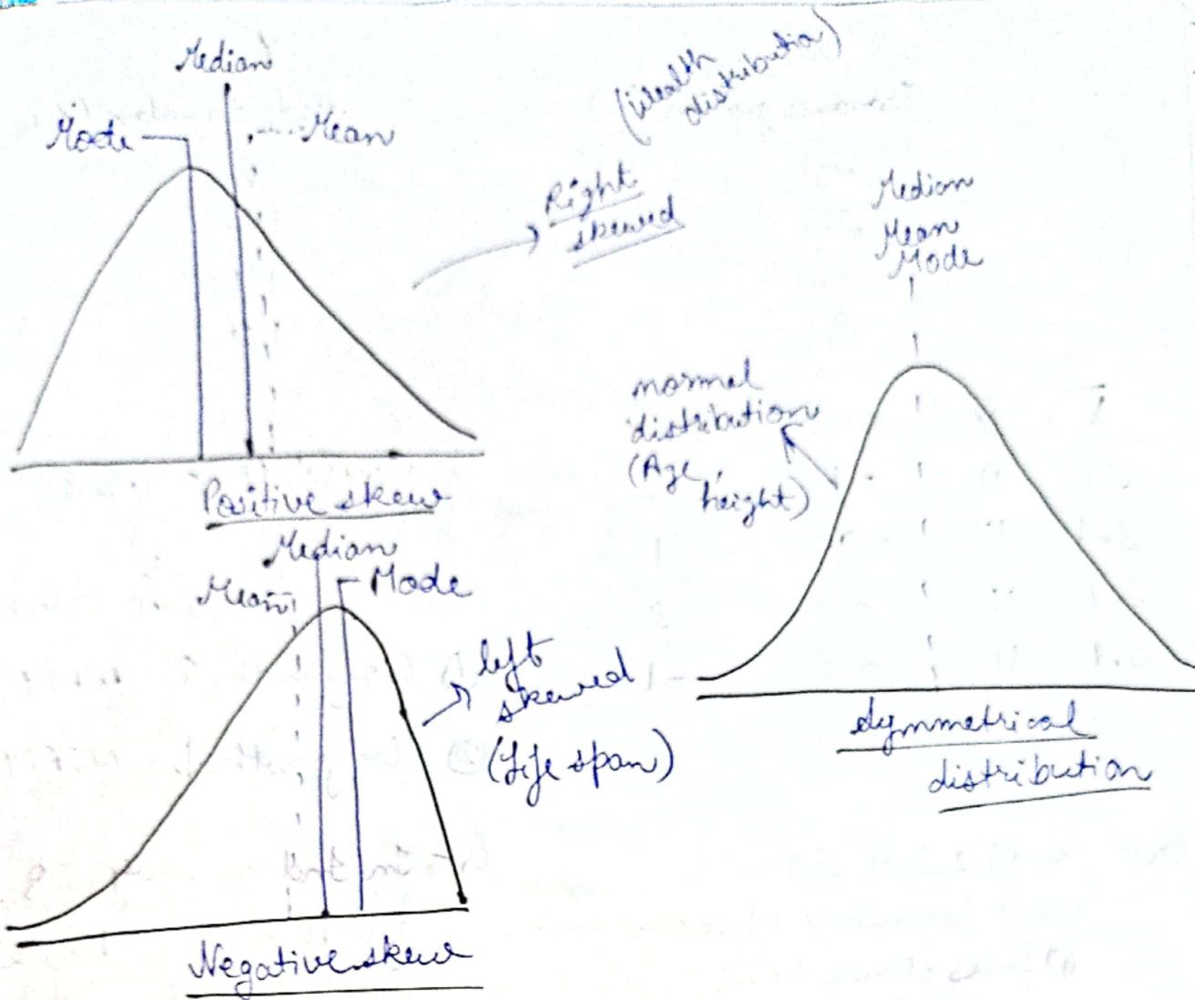
normally distributed



① The distribution of sample means will become more normal as its sample size increases.

② *RULE - sample distri. will be approximately normal if their sample size is

$$n > 30$$



COVARIANCE - Quantify the relationship between X & Y .

$$+ve \rightarrow \left\{ \begin{array}{l} X \uparrow Y \uparrow \\ X \downarrow Y \downarrow \end{array} \right\} \quad \left\{ \begin{array}{l} X \uparrow Y \downarrow \\ X \downarrow Y \uparrow \end{array} \right\} \leftarrow -ve$$

population -

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

sample -

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

$\text{cov}(x, x) = \text{var}(x)$

X
Economic growth (%)

2.1
2.5
4.0
3.6

Y
Nifty 50 index (%)

8
12
14
10

X	Y	$X - \bar{X}$	$Y - \bar{Y}$
3.1	11	-1	3
3.1	11	-0.6	1
3.1	11	0.9	3
3.1	11	0.5	-1

$$\text{Cov}(X, Y) = 1.533$$

positive value

- ① Eco. growth \uparrow NIFTY \uparrow
 ② Eco. growth \downarrow NIFTY \downarrow

- Q- In India avg. IQ is 100 with $\sigma = 15$, what is the % of population which you expect to have an IQ - ① lower than 85
 ② higher than 85
 ③ b/w 85 & 100

RW - ① Z score = $\frac{x_i - \mu}{\sigma}$

$$= \frac{85 - 100}{15}$$

$$= -0.25$$

② Z score = $\frac{3.75 - 4}{15}$

$$= -0.25$$

$$= 40\%$$

Normalization - [lower scale \leftrightarrow higher scale] $\rightarrow [0, 1]$

① Min-Max scaler [0-1]

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

We apply this in -

Deep learning & image processing

i.e. convert
pixels (0-255)
to (0-1)

x	\rightarrow	y
1		0
2		0.25
3		0.5
4		0.75
5		1

* Normalization & standardization are a part of Feature scaling

mostly used in Deep learning

mostly used in Machine learning

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 -

Data is choosed according to the categories.

Gender

\rightarrow Male

\rightarrow Female

Education degree

\rightarrow High School

Masters

Phd

③ Systematic sampling -

Select every n^{th} individual out from the population (N)

\hookrightarrow Convenience sampling -

Only those who are interested in the survey will only participate.

* Pearson Correlation Coefficient - (-1 to 1)

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

→ most common way of measuring a linear correlation.

Between 0 & 1 → +ve correlation - when one variable changes, the other variable changes in the 'same direction'.

0 → No correlation - there is 'no relationship' b/w the variables.

Between 0 & -1 → -ve correlation - when one variable changes, the other variable changes in the 'opposite direction'.

* Spearman Rank Correlation -

$$\gamma_s = \frac{\text{cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}$$

→ when the variables are normally distributed or the relationship b/w the variables is not linear.

x	y
10	4
8	6
7	8
6	10

R(x)	R(y)
4	1
3	2
2	3
1	4

→ Ascending Order

PROBABILITY

Probability = $\frac{\text{No. of ways an event can occur}}{\text{No. of possible outcomes}}$

Mutually exclusive - if events can not occur at same time.

(Non-) mutually exclusive - if events can occur at same time.

Eg-1 \rightarrow Mutually exclusive - in basketball

Q - Probability of getting heads or tails on flipping a coin?

$$\text{Ans} - P(H) = \frac{1}{2} = 0.5 \quad \text{and } P(T) = \frac{1}{2} = 0.5$$

Addition rule for mutually exclusive

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

$$\text{i.e. } P(H \text{ or } T) = 0.5 + 0.5 \\ = 1$$

Eg-2 \rightarrow Non-mutually exclusive event -

Q - When picking randomly from a deck of cards, what is the probability of choosing a card that is heart or a Queen?

$$\text{Ans: } P(\text{Heart}) = \frac{13}{52}, \quad P(\text{Queen}) = \frac{4}{52}, \quad P(\text{Q or H}) = \frac{1}{52}$$

Addition rule for non-mutually exclusive

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

$$= \frac{13}{52} + \frac{4}{52} - \frac{1}{52}$$

$$= \frac{16}{52} = 0.307 = 30.7\%$$

Independent events - two events are independent if they do not affect one another.

Dependent events - two events are dependent if they affect one another.

Eg-1 → Independent events

Q - What is the probability of rolling a 6 & then a 3, with a normal 6 sided dice.

Ans - Multiplication rule for independent events

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

$$P(6 \text{ and } 3) = P(6) * P(3) = \frac{1}{36}$$

Eg-2 → Dependent events

Q - What is the probability of drawing a "Queen" and then a "King" from a deck of cards?

Ans - Multiplication rule for dependent events

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

$$P(Q \text{ and } K) = P(Q) * P(K|Q)$$

$$= \frac{4}{52} * \frac{3}{51} = 0.006 = 0.6\%$$

PERMUTATION -

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

n = total # of obj.

r = # of obj. you are picking

COMBINATION -

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

④ Histogram -

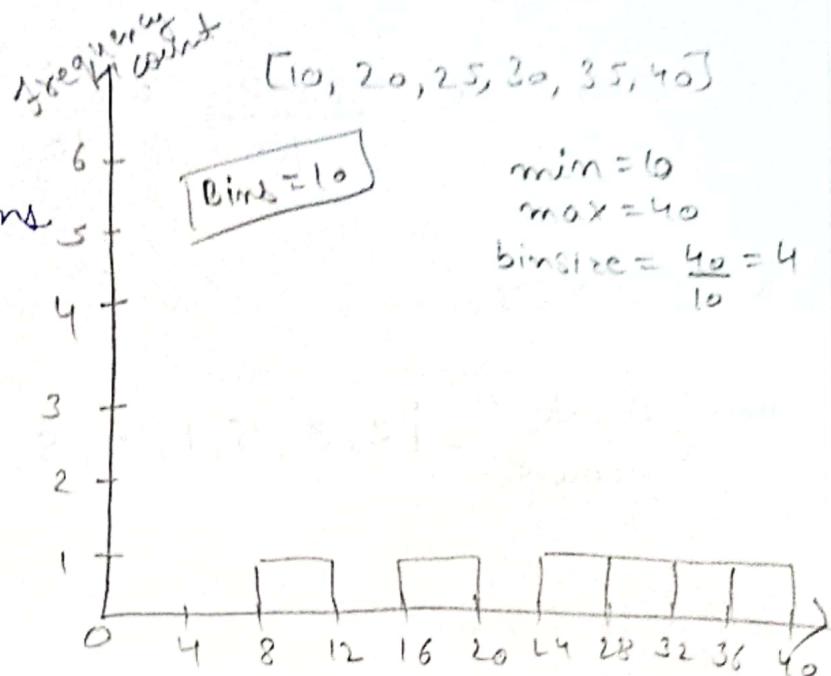
Ages = {0, 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100}

① Sort the numbers

② Bins → no. of groups

③ Bin size → size of bins

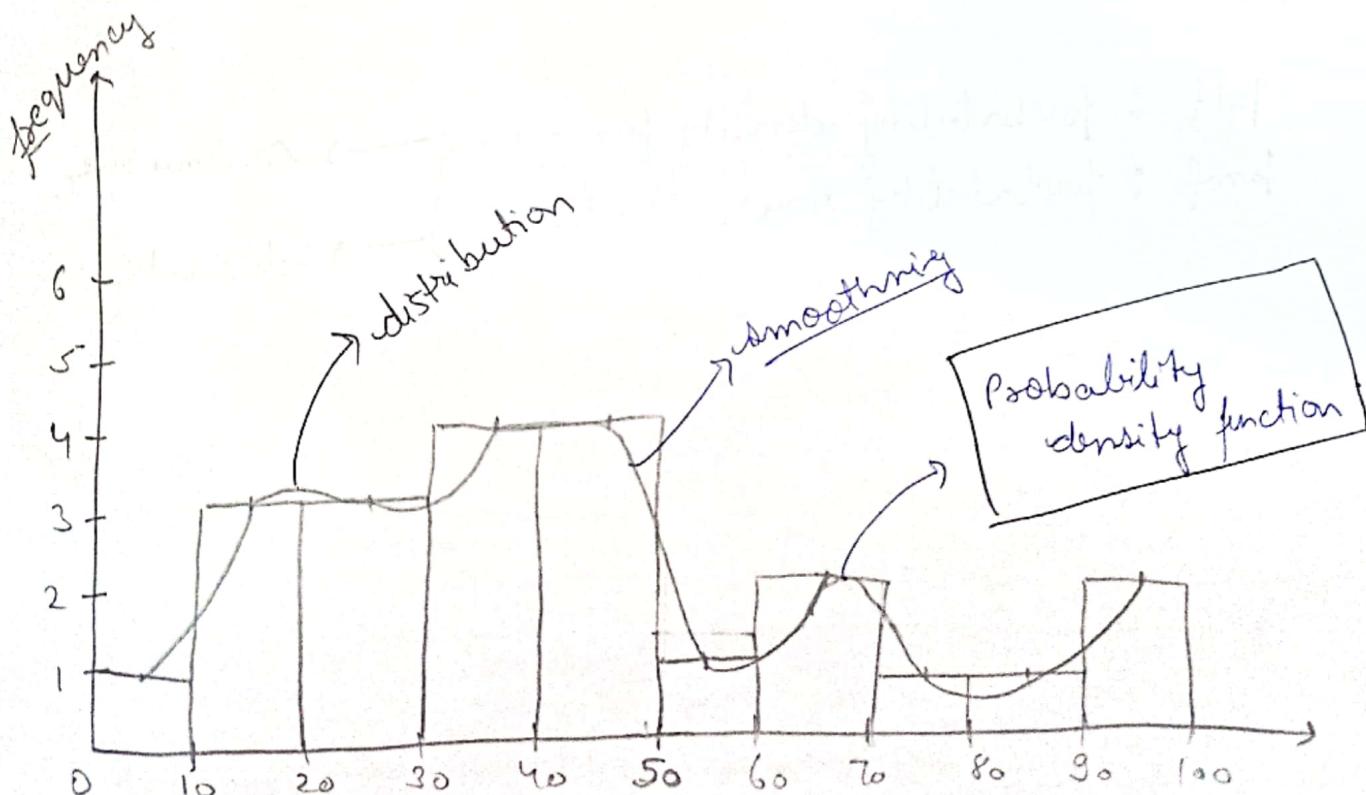
$$\boxed{\text{Bin size} = \frac{\text{Max} - \text{Min}}{\text{Bins}}}$$



For age dataset -

$$\text{bins} = 20 \quad \text{bin size} = \frac{100}{20} = 10$$

* CONTINUOUS



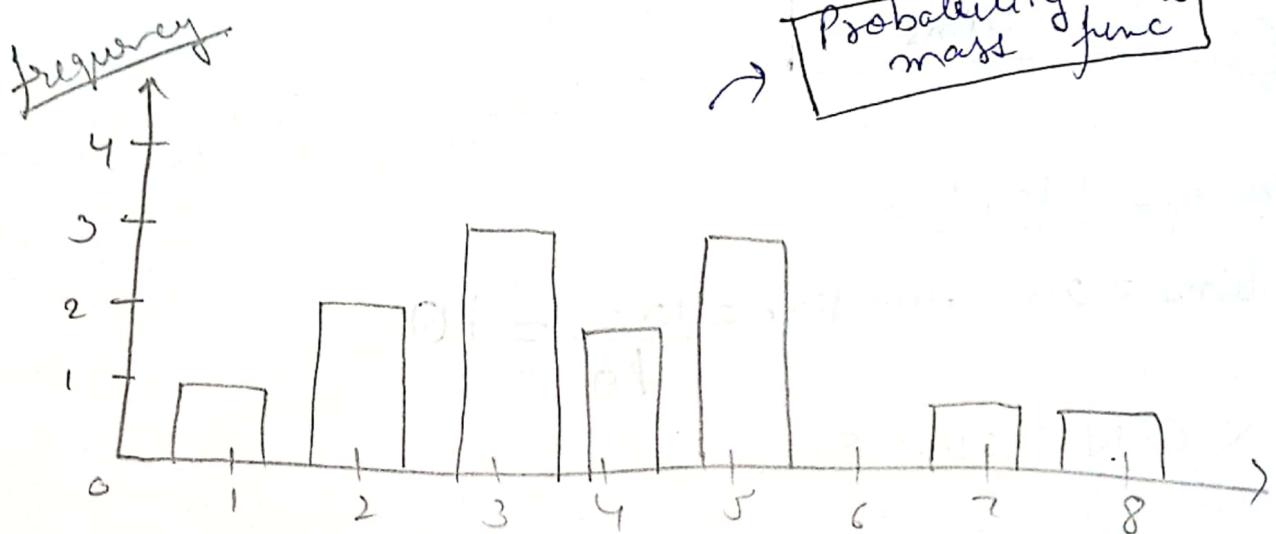
Weight = [30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 77
80, 90, 95]

bins = 10

$$\therefore \text{bin size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5 \rightarrow \text{continuous value}$$

* DISCRETE -

No. of Bank accounts = [2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5]



* pdf : probability density function } \rightarrow continuous
 pmf : probability mass function } \rightarrow discrete

Inferential statistics -

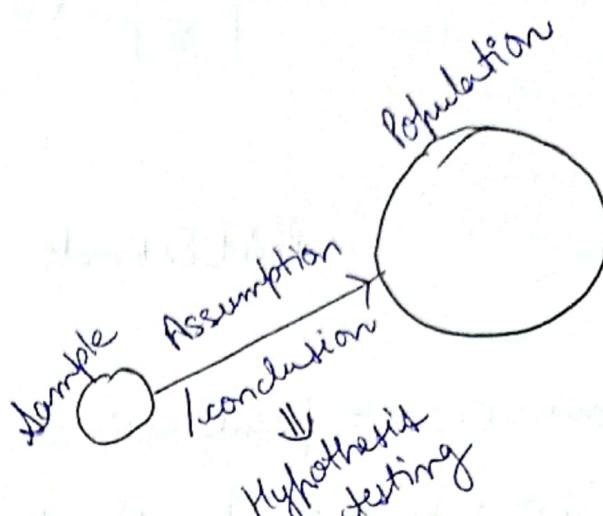
- Hypothesis testing
- p-value
- Confidence interval
- Significance value

* Hypothesis testing

Steps of hypothesis testing -

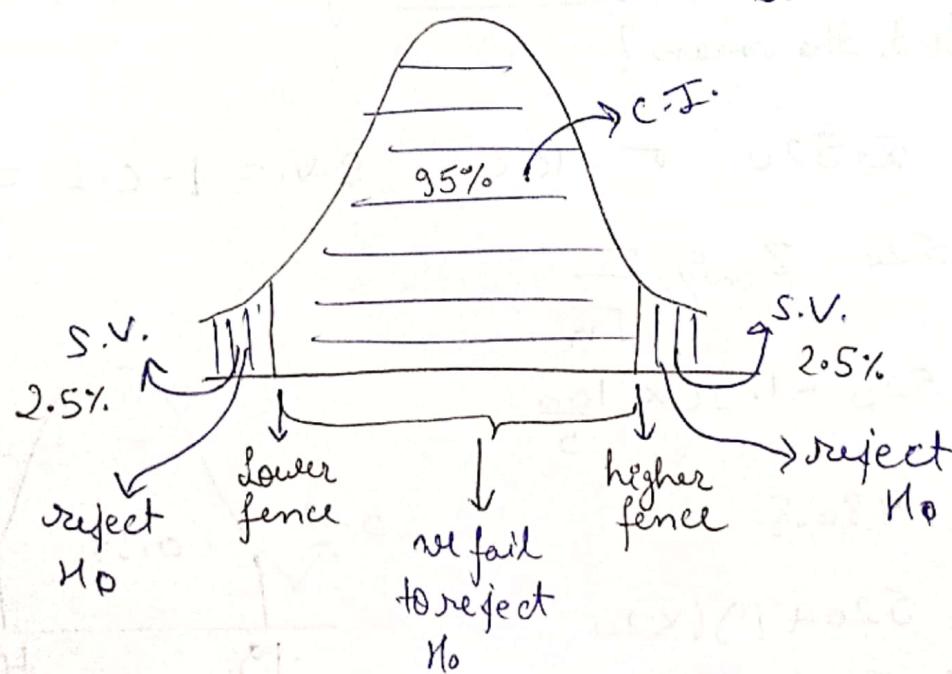
- ① Null hypothesis (H_0)
- ② Alternate hypothesis (H_1)
- ③ Perform experiments

↳ we fail to reject the null hypothesis [within C.I.]
↳ we reject the null hypothesis [outside C.I.]

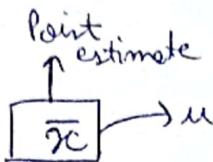


C.I.
[confidence interval]

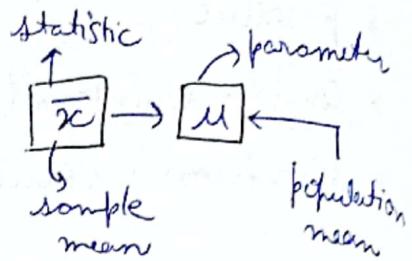
S.V.
[Significance value]
 $S.V. = 1 - C.I.$



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



$$\begin{cases} \bar{x} \geq \mu \\ \bar{x} \leq \mu \end{cases}$$



$$\text{Point Estimate} \pm \text{Margin of Error} = \text{Parameter} \rightarrow \text{Population mean}$$

Lower C.I. = point estimate - Margin of error

Higher C.I. = point estimate + Margin of error

$$\text{Margin of error} = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \rightarrow \begin{array}{l} \text{population SD} \\ \text{standard error} \end{array}$$

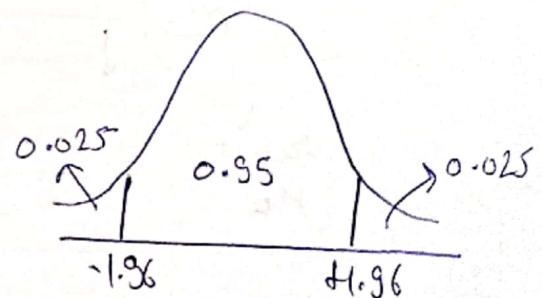
α = significance value

Q - On the quant test of CAT Exam, a sample of 25 test takers has a mean of 520 with a population S.D. of 100. Construct a 95% C.I. about the mean?

$$n = 25 \quad \bar{x} = 520 \quad \sigma = 100 \quad S.V. = 1 - C.I. = 0.05$$

$$\begin{aligned} \text{Lower C.I.} &= 520 - Z_{0.05/2} \frac{\sigma}{\sqrt{n}} \\ &= 520 - 1.96 \times \frac{100}{5} \\ &= 480.8 \end{aligned}$$

$$\begin{aligned} \text{Higher C.I.} &= 520 + 1.96 \times 20 \\ &= 559.2 \end{aligned}$$



* Population S.D - Z test + $n \geq 30$
 Sample S.D - T test + $n < 30$

Q- 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?

$$\bar{x} = 520 \quad s = 80 \quad n = 25 \quad C.I. = 95\% \quad S.V. = 0.05$$

$$\begin{aligned} \text{Lower C.I.} &= 520 - t_{0.05/2} \left(\frac{80}{\sqrt{25}} \right) && \text{Degree of freedom} \\ &= 520 - 2.064 \times 16 && = n-1 = 25-1 \\ &= 486.916 && = 24 \end{aligned}$$

$$\text{Higher C.I.} = 553.024$$

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

- a) State Null & Alternate hypothesis
- b) At 95% C.I., is there enough evidences to support machine is working properly or not.

Ans. Step 1 -

$$\text{Null hypothesis } \mu = 80$$

$$\text{Alternate hypothesis } \mu \neq 80$$

Step 2 -

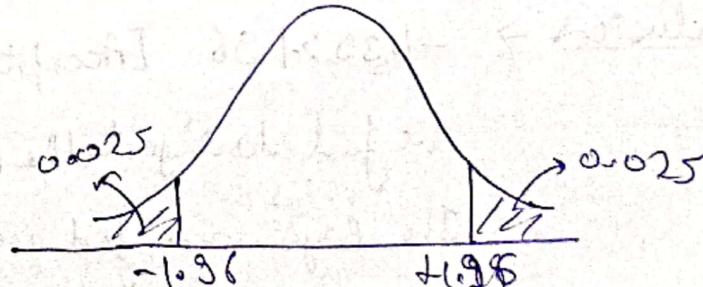
$$C.I. = 0.95 \quad S.V. (\alpha) = 1 - 0.95 = 0.05$$

Step 3 -

$$n = 40 \quad s = 2.5$$

Z-test

Decision boundary



Calculate test statistics -

$$Z = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \rightarrow \text{standard error}$$
$$= \frac{78 - 80}{\frac{2.5}{\sqrt{40}}} = -5.05$$

Conclusions -

Decision rule - If $Z = -5.05$ is less than -1.96 or greater $+1.96$, Reject the Null Hypothesis with 95% C.I.

Reject the Null hypothesis \rightarrow there is some fault in the machine.

- Q - A complain was registered, the boys in a Govt. school are underfed. Avg. weight of the boys of age 10 is 32 kgs with $S.D. = 9$ kgs. A sample of 25 boys were selected from the govt. school and the avg. weight was found to be 29.5 kgs? with $C.I. = 95\%$, check it is true or false.

A.S. $n = 25 \quad \mu = 32 \quad \sigma = 9 \quad \bar{x} = 29.5$

Step 1 -

$$H_0 : \mu = 32$$

$$H_1 : \mu \neq 32$$

Step 2 - C.I. = 0.95 $\alpha = 1 - 0.95 = 0.05$

Step 3 -

$$Z \text{ score} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{29.5 - 32}{\frac{9}{\sqrt{25}}} = -1.39$$

Conclusion $\rightarrow -1.39 > -1.96$ [Accept the Null Hypothesis 95% C.I.]
We fail to reject the Null Hypothesis
The Boys are fed well.

A factory manufactures cars with a warranty of 5 years or more on the engine and transmission. An engineer believes that the engine or transmission will malfunction in less than 5 years. He tests a sample of 40 cars and finds the average time to be 4.8 years with a standard deviation of 0.50.

- ① State the null & alternate hypothesis.
- ② At a 2% significance level, is there enough evidence to support the idea that warranty should be revised?

$n=40 \quad \bar{x} = 4.8 \text{ years} \quad s=0.50$

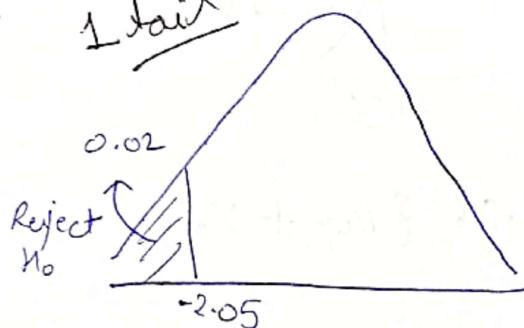
Step 1 -

$$H_0: \mu \geq 5 \quad \{\text{Null hypothesis}\}$$

$$H_1: \mu < 5 \quad \{\text{Alternate hypothesis}\}$$

Step 2 - $\alpha = 0.02 \quad C.I. = 1 - 0.02 = 0.98 = 98\%$

Step 3 - 1 tail



Step 4 - $Z\text{score} = \frac{\bar{x} - \mu}{s/\sqrt{n}}$

$$Z\text{score} = \frac{4.8 - 5}{\frac{0.50}{\sqrt{40}}} = \frac{-0.2}{0.08} = -2.52$$

$$Z\text{score} = -2.52$$

Step 5 - Conclusion: $-2.52 < -2.05$ — Reject H_0
 \rightarrow warranty needs to be revised.

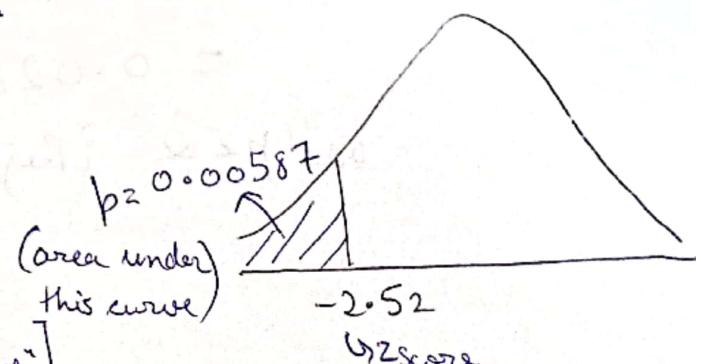
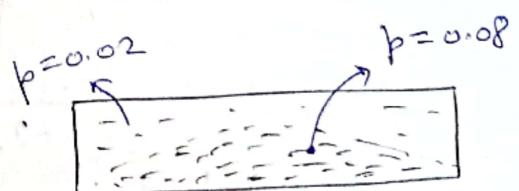
P-value

$p\text{-value} \neq \text{significance value}$
 ↓
 should always be less than or equal to s.v.
 ↓ derived by C.I.

as $p\text{-value} < \alpha (0.02)$

so we reject H_0

[the area under the curve we get through Z-score is "p-value"]



a- The avg. weight of all residents in a town xyz is 168 pounds. A nutritionist believes that the true mean to be different. She measured the weight of 36 individuals and found the mean to be 169.5 pounds with a standard deviation of 3.9.

a) Null & alternate hypothesis

b) 95%. Is there enough evidence to discard the H_0 ?

A.S.

$$\bar{x} = 169.5 \quad s = 3.9 \quad n = 36 \quad \mu = 168$$

$$C.I. = 0.95$$

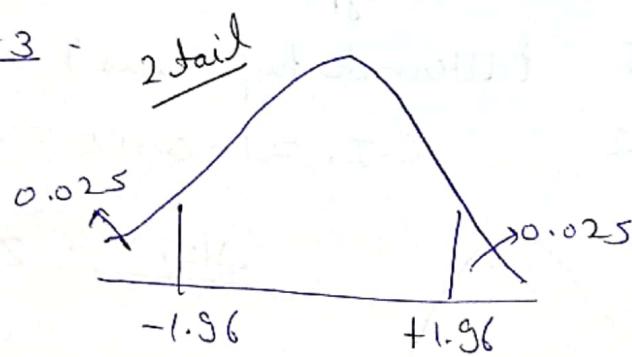
Step-1 - $H_0: \mu = 168$

$H_1: \mu \neq 168$

Step-2 - $C.I. = 0.95$

$$\alpha = 0.05$$

Step-3 -



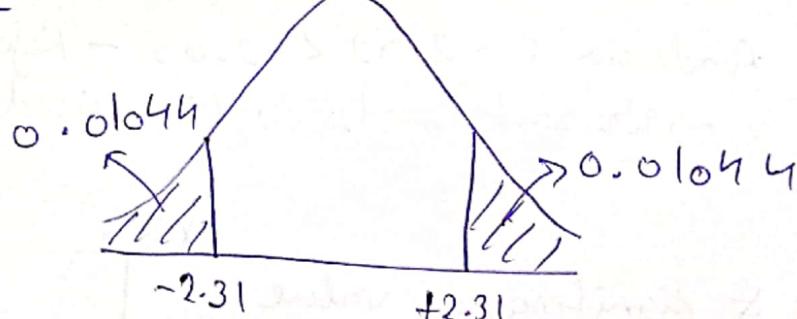
Step-4

$$Z\text{ score} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} \\ \approx 2.31$$

Step-5 -

$$2.31 > 1.96 \quad \{ \text{Reject } H_0 \}$$

P-value -



$$P\text{ value} = 0.01044 + 0.01044 \\ = 0.02088$$

as $P < \alpha \quad \{ \text{Reject } H_0 \}$

Q- A company manufactures bikes batteries with an avg. life span of 2 years or more years. An engineer believes this value to be less. Using 10 samples, he measures the avg. life span to be 1.8 years with a S.D. of 0.15.

a) State the Null and alternate hypothesis

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

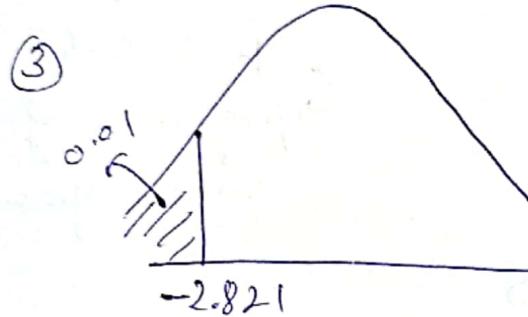
$$\text{Ans} \quad ① \quad H_0 : \mu \geq 2$$

$$H_1 : \mu < 2$$

$$② \quad C.I = 0.99 \quad n = 10 \\ S.V. = 0.01$$

as $n < 30$

$$\text{degree of freedom} = n - 1 \\ = 10 - 1 = 9$$



$$④ \quad t = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{1.8 - 2}{0.15/\sqrt{10}}$$

$$⑤ \quad -4.216 < -2.82 \\ \{ \text{Reject } H_0 \}$$

\Rightarrow The avg. life span of battery is less than 2 years

* Z test with proportions -

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

p_0 - population proportion

\hat{p} - sample proportion

$$q_0 = 1 - p_0$$

$$\hat{p} = \frac{x}{n}$$

\hookrightarrow proportion

$p\text{value} > \alpha \quad \{ \text{Fail to reject } H_0 \}$
else $\{ \text{Accept } H_0 \}$

Q - A tech company believes that the percentage of residents in town XYZ that owns a cell phone is 70%. A marketing manager believes that this value to be diff. He conducts a survey of 200 individuals and found that 130 responded Yes. Owning a cell phone?

① State H_0 & H_1 .

② At a 95% CI, is there enough evidence to reject H_0 ?

Ans. Step-1

Null hypothesis $\Rightarrow P_0 = 0.70$

$$n = 200 \quad x = 130$$

Alternate hypothesis $\Rightarrow P_0 \neq 0.70$

$$\hat{p} = \frac{130}{200}$$

$$q_0 = 1 - P_0 = 0.30$$

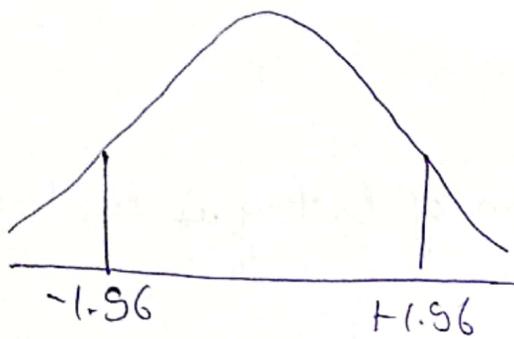
$$\hat{p} = 0.65$$

Step-2

$$C.I. = 0.95 \quad \alpha = 0.05$$

\hookrightarrow proportion of saying "Yes"

Step-3



Step-4

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

Step-5

$$\text{Conclusion: } -1.54 > -1.96$$

{Fail to reject the null hypothesis}

P-Value



$$p\text{-value} = 0.12356$$

Q - A car company believes that the % 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.

H₀:

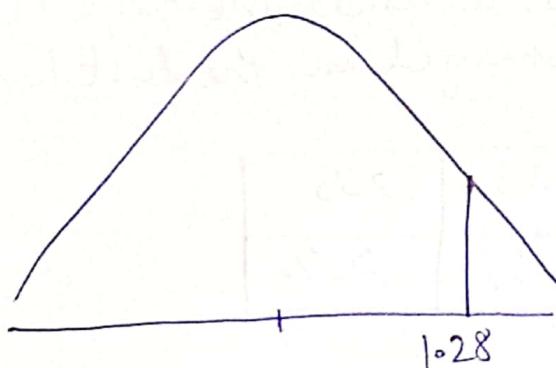
$$H_0 : p_0 \leq 0.60$$

$$H_1 : p_0 > 0.60$$

$$\hat{p} = \frac{170}{250} = 0.68$$

$$q_0 = 0.40$$

$$Z\text{score} = \frac{0.68 - 0.60}{\sqrt{\frac{0.6 \times 0.4}{250}}} = \frac{0.08}{0.0309} = 2.588$$



Reject the Null hypothesis

* Chi-square test -

→ Chi-square test claims about population proportion.
It is a non-parametric test that is performed on categorical data.

↳ ordinal / nominal data / data

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

Q - In the 2000 U.S. census the age of individuals in a small town found to be following -

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

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

<18	18-35	>35
121	288	91

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

Ans.

	<18	18-35	>35
Expected \rightarrow	20%	30%	50%

$n=500$

	<18	18-35	>35
Obs. \rightarrow	121	288	91
Exp. \rightarrow	100	150	250

Step-1 \rightarrow Null hypothesis - H_0 : The data meets the exp. distribution
 H_1 : The data does not meet the expected distribution

Step-2 : $\alpha = 0.05$ C.I. = 95%

Step-3 : degree of freedom {categories} $df = c-1 = 3-1 = 2$

Step-4 : decision boundary = 5.991
{chi-square-table}

Step-5 : Test statistics -

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

$$\chi^2 = 232.494$$

Step-6 : Conclusion -

$$\chi^2 > 5.99 \quad \{\text{Reject } H_0\}$$

* Anova test - (F-test)

→ The t-test works well when dealing with two groups, but sometime we want to compare more than two groups at the same time, so here we use F-test.

→ Extension of t-test Lz-test.

$$F = \frac{\text{measure of effect (MSEffect)}}{\text{measure of error (MSError)}}$$

$F \geq 1$, variation due to effect \geq due to error

$F = 1$, variation due to effect = due to error

$F \leq 1$, variation due to effect \leq due to error

* Bernoulli distribution -

- Used when we want to model the outcome of a single trial of an event.

Eg - Tossing of coin (1 time) - $S = \{H, T\}$

p - probability of success

q - probability of failure

For getting head -

$$P(H) = p = 0.5$$

For not head -

$$P(\text{not } H) = q = 0.5$$

$$\begin{aligned} p &= 1-q \\ q &= 1-p \end{aligned}$$

* Binomial Distribution -

- Used when we want to model the outcome of multiple trials of an event.

Eg - Tossing a coin (twice)

$$S = \{HH, HT, TH, TT\}$$

For $X = \text{no. of heads}$

$$P(X=0) = 1/4 = 0.25$$

$$P(X=1) = P(HT) = 1/4 + 1/4 = 0.5$$

$$P(X=2) = P(HH) = 1/4 = 0.25$$

Similarly if tossed 3-times

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

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

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

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

$$P(X) = \frac{n!}{(n-x)!x!} p^x q^{n-x}$$

n = no. of trials

x = no. of success

p = prob. of success

$q = 1 - p$ = prob. of failure

* POWER LAW / PARETO PRINCIPLE

- b/w two quantities
- 80-20 rule
- i.e. 80% of consequences come from 20% of causes

