

TOPS TECHNOLOGIES AHMEDABAD

NAME \Rightarrow TANMAY SUTHAR

Course \Rightarrow Data Analysis

12 Dec. Batch

What is Statistics

It is a science of collecting, organizing, analyzing data for better decision making.

What is Data

Facts or pieces of information that can be measured.

Eg:-

Age of class
S 89, 90, 50 3

Age of students
S 24, 22, 21 3

Types of Stats

Descriptive Stats: It consist of organizing & summarizing data

Inferential Stats: Using data, we can make conclusion using some techniques

Eg :

class of 20 students

1st Sem Maths E 86, 70, 90, 55 ... 3

1 what is the avg of class (descriptive)

2 7th Sem (Inferential)

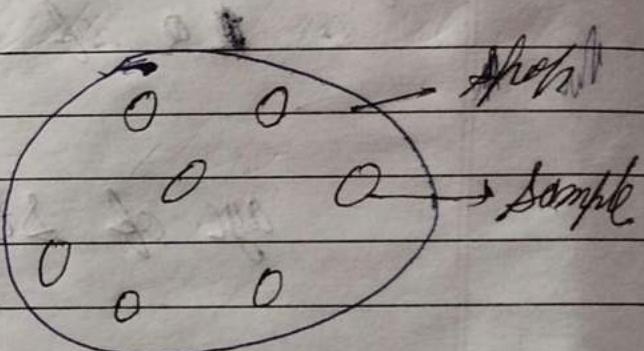
* sample and population

Population = N



whole dataset is

Known as
population



Sample = n

Small subsets of data
taken from population
→ sample

Sampling Techniques

1 Simple Random Sampling

every member of population has an equal chance of getting selected in sample (n)

2 Stratified Sampling

splitting data into non-overlapping groups

Age group — 0 - 20
→ 20 - 40

Gender — M
F

3 Systematic Sampling

4 Convenience Sampling

domain expert

1. Exit poll ?

2.

(*) A method for selecting a sample from a population in a randomized manner.

* Variables

It is a property that can hold / store / take any value

Age : { 8, 10, 15, 20, 25 }

Marks : { 76, 80, 95 }

Types

1 Qualitative : Categorical values

(based on some characteristics
we can derive categorical values)

I Q : 0 - 10 → low →
10 - 50 → Avg
50 - → good

2 Quantitative : Numerical value (measurable numerically)

height : { 162, 159, 155 } . 3

Weight : { 59, 65, 79 } . 3

discrete (int)
whole No.

continuous (float)
decimal No.

- 1 No. of Students 1 Height
165.2, 167.9
- 2 No. of Rank also 2 Weight
165.5 - 160.9 - 3

- 1 Blood Pressure → continuous / discrete
- 2 Marital Status → qualitative
- 3 River length → cont.
- 4 Day length → cont.
- 5 gender → qualitative & category 3

Variable Measurement scales

- 1 ordinal : ordered [rank, order, dates, graduation]
- 2 Nominal : categorical values (colors, classes, degrees)
- 3 Interval : [No zero / absolute point] (ordered as well as value ratios)
- 4 Ratio : zero means nothing

Interval

$20^{\circ}\text{C} : 210^{\circ}\text{C}$

Ratio

$20\text{kg} : 40\text{kg}$

$1 : 2$

$20^{\circ}\text{C} : \underline{210^{\circ}\text{C}}$

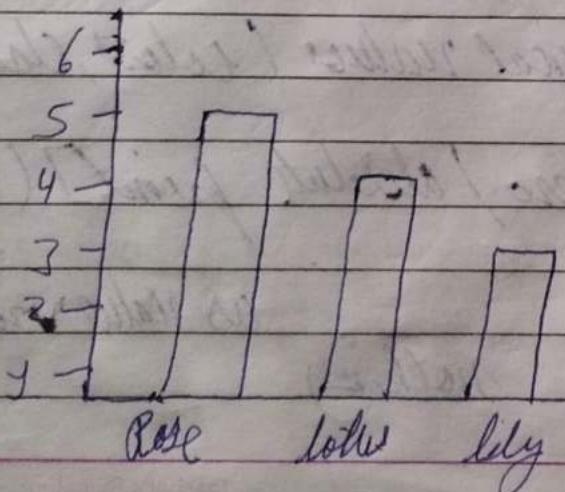
$1^{\circ}\text{C} : 2^{\circ}\text{C}$

Frequency

data : flowers

{ Rose, lily, lotus, Rose, rose, rose,
rose, lotus, lily, lotus, lily, lotus }

flowers	frequency	(Cumulative)
Rose	5	5
lily	3	8
lotus	4	12
	$\overline{12}$	

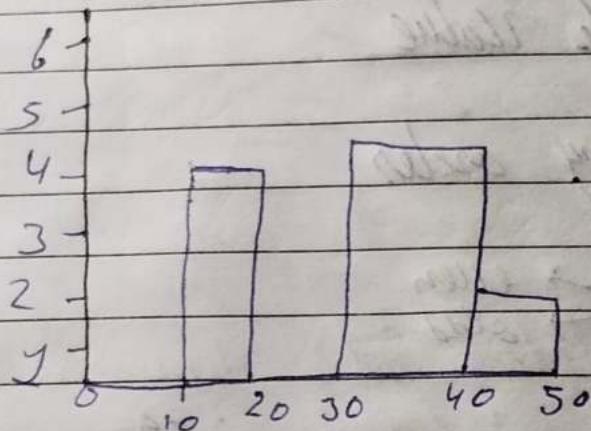


Bar graph
chart

Histogram

Marks = [12, 15, 12, 15, 21, 27, 28, 35, 34, 36
 39, 42, 45]

bins	(10 - 20)	4	→ continuous
	(20 - 30)	3	
	(30 - 40)	4	
	(40 - 50)	2	



Measure of central Tendency

Avg → Mean

$$\text{DOP}$$

$$\bar{x} = \frac{\sum x_i}{N}$$

$$\bar{x} = \frac{\sum x_i}{n}$$

5 2, 3, 5, 3, 2, 1, 3

$$\underline{2+3+5+3+2+1+3}$$

7

\Rightarrow Mean : It refers to the measure used to determine the center of the distribution of the data.

$$\{1, 2, 2, 2, 3, 3, 4, 5, 5, 6, 100\} \rightarrow \text{outliers}$$

$$\frac{32}{10} \Rightarrow 3.2 \rightarrow \frac{132}{11} \Rightarrow 12$$

\Rightarrow Median : middle value

\hookrightarrow ascending order

\hookrightarrow data \rightarrow even
 \rightarrow odd

$$\Rightarrow \text{even} \quad \frac{\left(\frac{n}{2}\right)^{\text{th}} + \left(\left(\frac{n}{2}\right)^{\text{th}} + 1\right)^{\text{th}}}{2}$$

dataset : $\{11, 12, 13, 14, 15, 16\}$

$$n = 6 \quad \frac{\left(\frac{6}{2}\right)^{\text{th}} + \left(\frac{6}{2}\right)^{\text{th}} + 1}{2}$$

$$\frac{3^{\text{rd}} + (3+1)^{\text{th}}}{2}$$

$$\frac{3^{\text{rd}} + 4^{\text{th}}}{2} \Rightarrow \frac{13 + 14}{2}$$

27

2

$$\Rightarrow 13.5$$

$$\text{Median} = 13.5$$

odd

$$\frac{(n+1)^{\text{th}}}{2}$$

{ 11, 12, (13), 14, 15 }

$$n = 5$$

$$\frac{(5+1)^{\text{th}}}{2} \div \frac{1^{\text{st}}}{2} = 3 \therefore 13$$

{ 11, 12, 13, 14, 15, 100 }

$$\frac{(\frac{n}{2})^{\text{th}} + ((\frac{n}{2})^{\text{th}} + 1)^{\text{th}}}{2}$$

$$\frac{(-\frac{1}{2})^{24} + ((-\frac{1}{2})^{24} + 1)^{24}}{2}$$

$$\frac{3^{24} + 4^{24}}{2} = \frac{13 + 14}{2} = 13.5$$

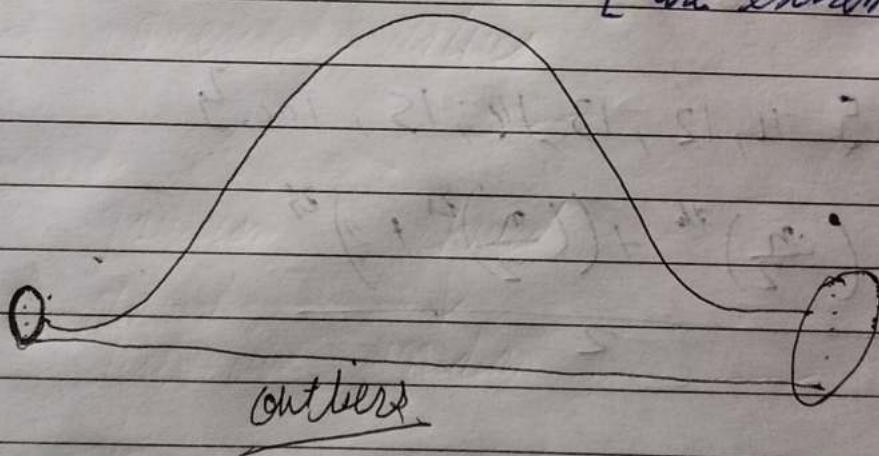
5 21, 23, 25, 29, 32 100 3

$$\frac{25+29}{2} = \frac{54}{2} = 27$$

~~#~~ outlier :-

a data points who doesn't follow pattern or trend of the dataset then it is considered as outlier.

[are extreme points]



Mode

Most frequent value (repeated)

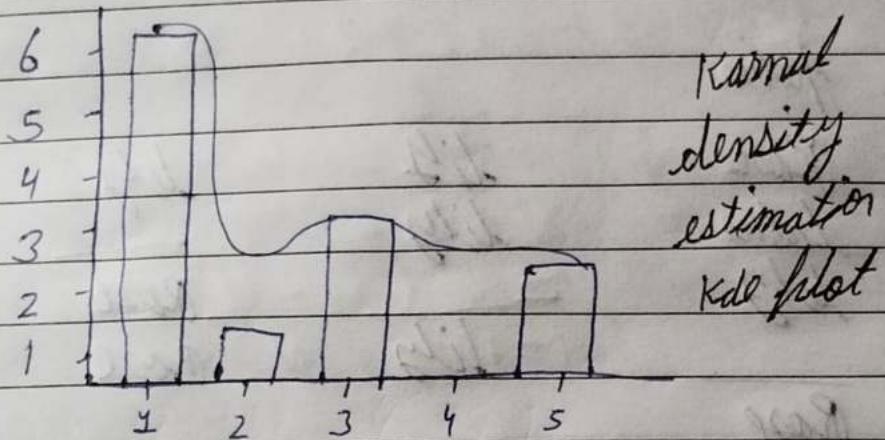
[1, 1, 2, 3, 5, 1, 1, 3, 5, 1]

$$1 = 6$$

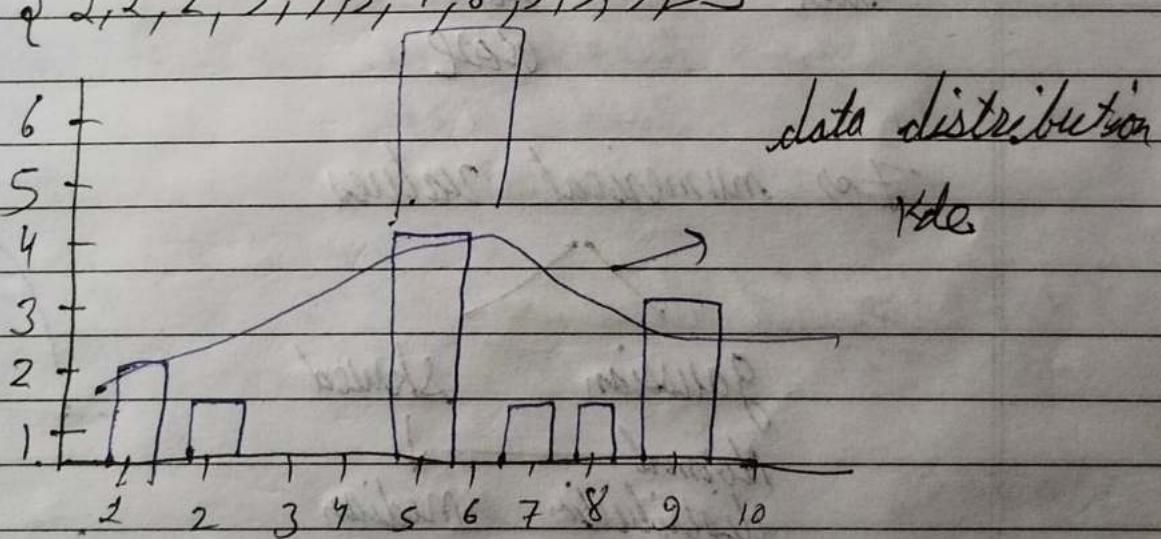
$$3 = 3$$

$$2 = 1$$

$$5 = 2$$



{2, 1, 2, 5, 9, 5, 7, 8, 9, 5, 9, 5}



for categorical missing data

0 - 5 + → } mode

→ ↗ new category "missing"
"unknown"

0?

"random"

15 species

Rose

lily

lily

lotus

lily

null

lily

Rose

na

lily

Rose:

Rose

15 - (10)

lily

Rose

Rose

lotus

Rose

For numerical values



Gaussian

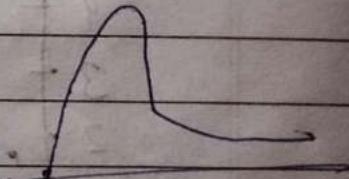
Normal

distribution

skewed

↓

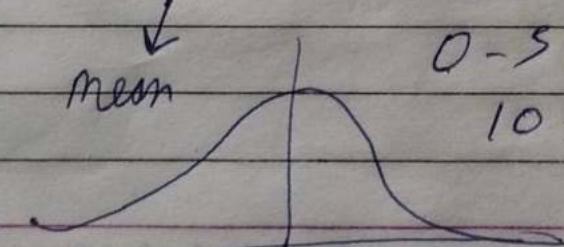
median



0 - 5 +

10 +

mean



mean = median = Mode

measure of dispersion → spread

(15)

$$[5, 1, 1, 1, 2, 3] \rightarrow \frac{5+5}{5} = 2$$

$$[2, 2, 2, 2, 2] = 2$$

→ Variance

It means how far the numbers in a dataset are from the mean (avg.) (how each value differs from a dataset in mean).

High Variance → more spread (far from mean)

low Variance → closer to mean

POP

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

summation

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

s,

Bessel's correction

is the use of $n-1$ instead of n in the formula for the sample variance and sample standard deviation.

Teacher's Signature

x	$x - \bar{x}$	$(x_i - \bar{x})^2$
1	-1.83	3.34
2	-0.83	0.69
2	-0.83	0.69
3	0.17	0.02
4	1.17	1.36
5	2.17	4.70

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

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N} = \frac{10.8}{6}$$

$x_i \rightarrow$ every data point

$\bar{x} \rightarrow$ mean of pop

$N \Rightarrow$ total pop

$\sum =$ summation

$$\sigma^2 = 1.8$$

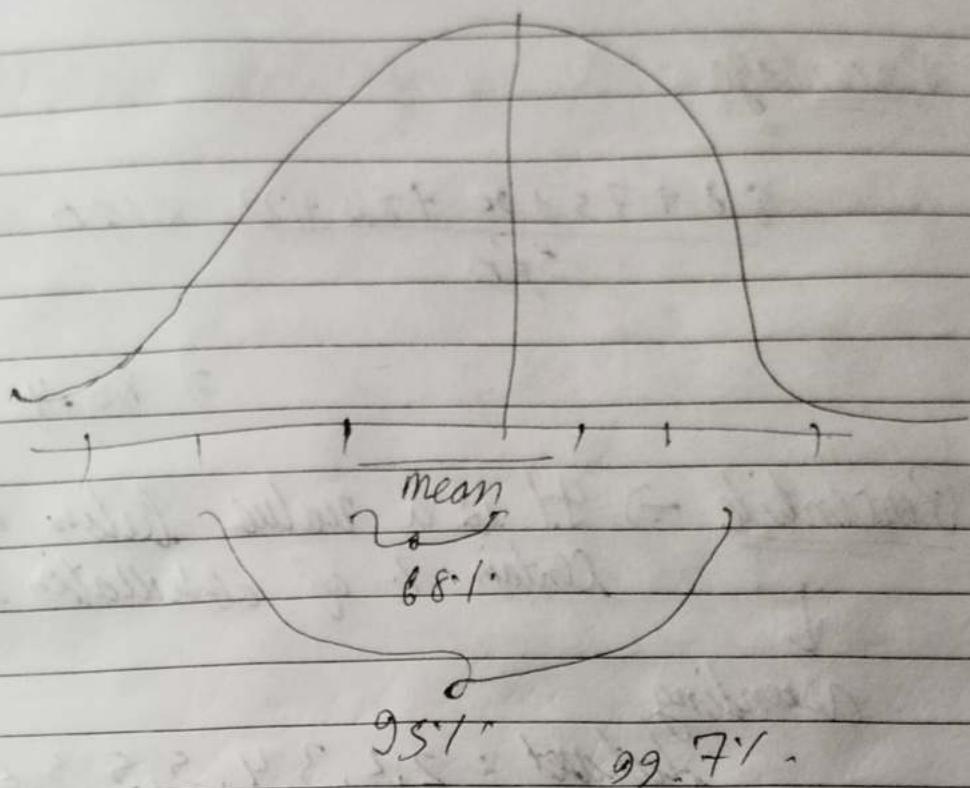
standard deviation

pop

sample

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

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



$$\mu + 1\sigma = 68.1$$

$$\mu + 2\sigma = 95.1$$

$$\mu + 3\sigma = 99.7$$

→ square root of variance

→ It gives measure of spread that is in the same unit as the original data, making it easier to interpret.

Percentage

$$\frac{88 + 75 + 90 + 80 + 99}{500} \times 100$$

$$\Rightarrow 86.4$$

⇒ Percentile ⇒ It is a value below which a certain % of observation lie.

↓
Ascending

dataset = 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

$$n = 20$$

Percentile rank = $\frac{\text{no. of values below } x}{n}$

$$\Rightarrow \frac{16}{20} = \frac{4}{5} \times 100$$

$$80\% \text{ values are below } 10 = 80$$

$$11 = \frac{17}{20} \times 100 = 85$$

$$85\% \text{ values are below } 11$$

what value exists at percentile rank 25?

$$\text{value} = \left(\frac{\text{percentile}}{100} \times n \right) + 1$$

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

$\underbrace{}_3$

$$5 + 1 = 6 \rightarrow \text{index value}$$

$$75 = \left(\frac{75}{100} \times 29 \right) + 1$$

$\underbrace{}_2$

$$\therefore 15 + 1 = 16$$

9

Five Number Summary

1. Minimum

2. 25 percentile \rightarrow first quartile

Median \rightarrow 50 Percentile

75 percentile

Maximum

[1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 15, 27]

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

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

$$IQR \rightarrow Q_3 - Q_1 \rightarrow \text{Inter quartile Range}$$

$$Q_1 = 3$$

$$Q_3 = 8$$

$$IQR = 8 - 3 = 5$$

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

$$= 3 - 1.5(5)$$

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

$$= 8 + 1.5(5)$$

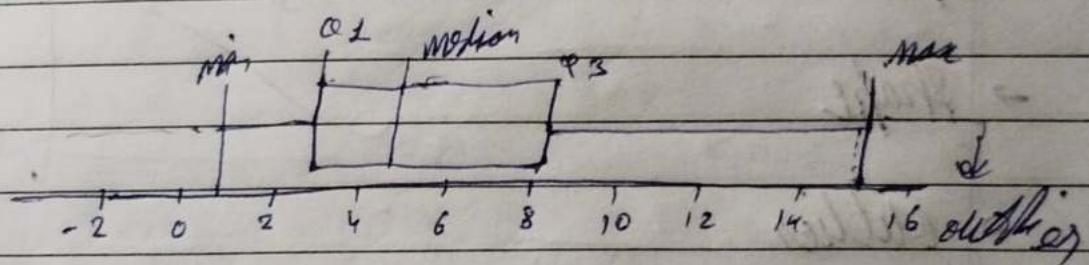
$$1 \text{ min} = 1$$

$$2 \text{ sec} = 3$$

$$3 \text{ motion} = 5$$

$$4 \text{ } \alpha_3 = 8$$

$$5 \text{ max} = 15$$

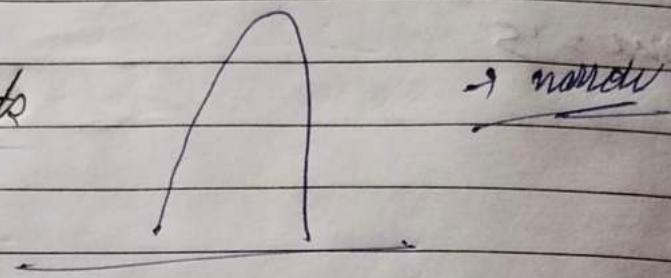


distribution

data distribution

- It refers to a way in which values or data points are spread or arranged.
- It shows how often different values occur in data set & describes the overall pattern of the data.

→ 4 key points



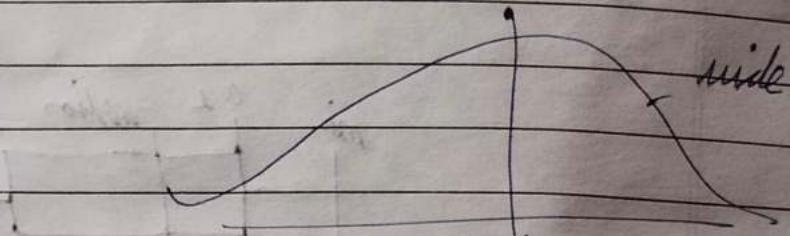
→ narrow

→ center

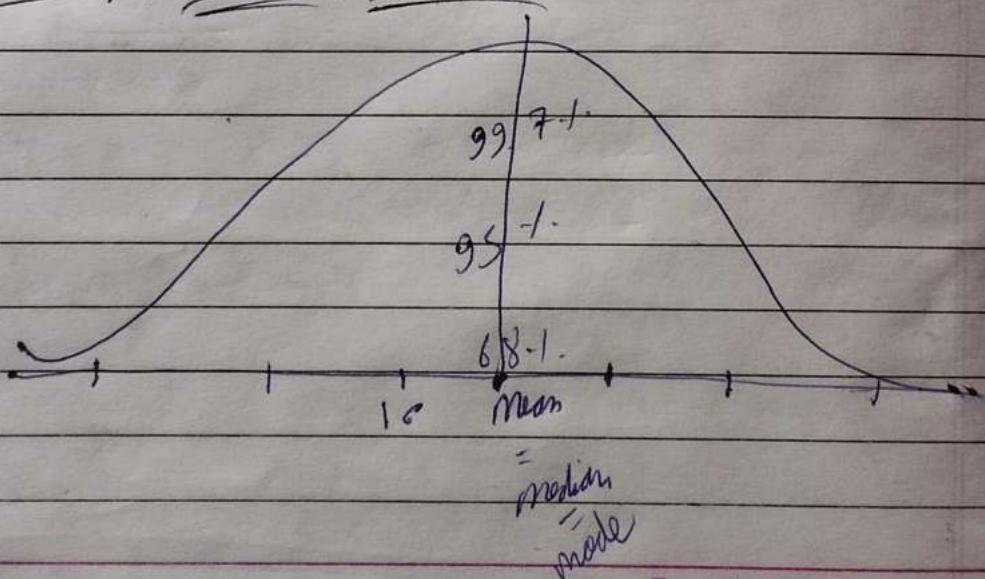
→ spread.

→ shape

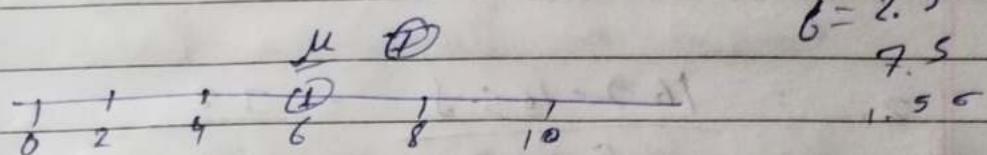
→ outliers



gaussian / Normal distribution



$$\left. \begin{array}{l} \mu + 1\sigma = 68.1 \\ \mu + 2\sigma = 95.1 \\ \mu + 3\sigma = 99.7 \end{array} \right\} \text{empirical rule}$$



Z-score

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

$$\mu = 0 \quad \sigma = 1$$

Standard Normal distribution



$$\mu = 0 \quad \sigma = 1$$

height	weight	$z = \frac{x - \mu}{\sigma}$	
169	60	3.2	10.24
172	65	6.2	38.44
150	45	-15.8	249.64
168	70	2.2	4.84
170	71	4.2	17.64
<u>829</u>			<u>320.8</u>

Teacher's Signature

$$\text{M}_n = \frac{829}{5} = 165.8$$

$$\text{C}_4 = \frac{320.8}{5} = 64.16$$

$$\text{C}_7 = 8.009$$

$$Z_1 = \frac{169 - 165.8}{8.009} = \frac{3.2}{8.009} = 0.399 \Rightarrow 0.4$$

$$Z_2 = \frac{6.2}{8.009} = 0.7$$

$$Z_3 = \frac{-15.8}{8.009} = -1.9$$

$$Z_4 = \frac{2.2}{8.009} = 0.27$$

$$Z_5 = \frac{4.2}{8.009} = 0.52$$

$$\Rightarrow -0.001 \\ \Rightarrow -0.01$$

Normalization

$$x = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \rightarrow \begin{matrix} \text{in max} \\ \text{scales} \end{matrix}$$

skewed distribution

positively skewed distribution

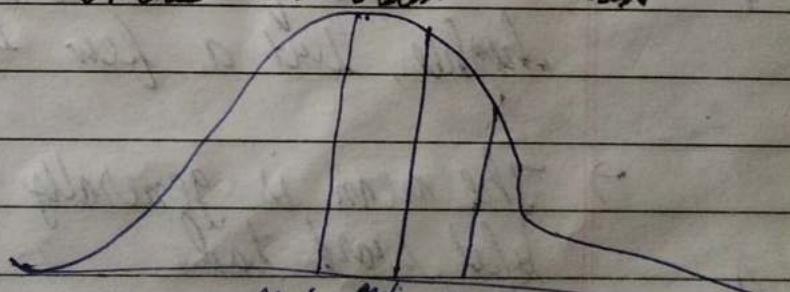
What?

In positively skewed distribution, most values are concentrated on the lower end with a long tail extending to the right. A few high values pull the average to the right of the median.

Skewness \rightarrow A distortion or asymmetry that deviates from the symmetrical bell curve

\rightarrow also called right-skewed or right-tailed distribution

Mode < median < mean



Mean > median > mode

Teacher's Signature

When?

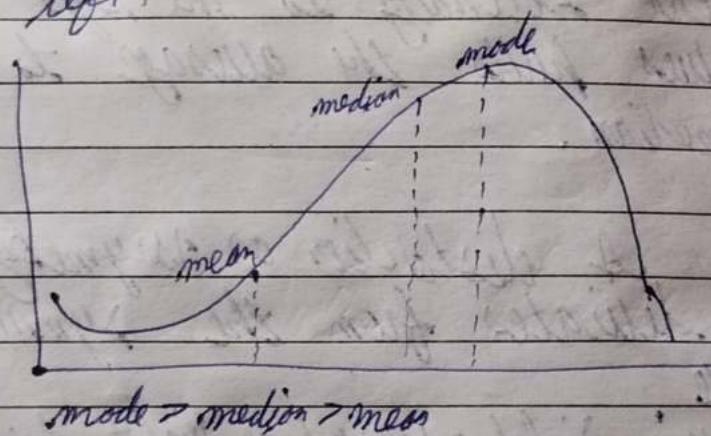
useful for data with rare but significant high values such as income levels where few individuals earn much more than the rest.

- The mean is greater than the median due to outlier on the higher end.

* Negatively skewed distribution

What?

most values are clustered at higher end, with a few values creating a long tail to the left.



When?

used for datasets where values are typically higher, but a few lower values exist (retirement age)

- The mean is generally less than median due to leftward tail.

→ Also called left skewed or left tailed.

why?

It helps detect cases where data points are generally higher but occasionally much lower.

Exponential distribution (continuous)
 $\lambda = \text{constant rate}$

It describes the time b/w events in a process where event occurs independently at a constant rate λ .

$$f(x) = \lambda e^{-\lambda x} \quad x \geq 0$$

* Bernoulli distribution (discrete)

It models a single experiment with two possible outcomes.

Success ($x=1$) , failure ($x=0$)

Binomial distribution (discrete)

Binomial dist. extends Bernoulli to n independent trials.

$$P(X=k) \Rightarrow \binom{n}{k} p^k \cdot (1-p)^{n-k}$$

(so)

X = random variable (no. of success out of n trials)

n = total no. of trial (100)

K = No. of success ($0 \leq K \leq n$)

$p = P(\text{Success in 1 trial}) = 0.5$

$$\binom{n}{K} = \frac{n!}{(n-K)!K!}$$

$n = 5$ $K = 3$ $p = 0.5$

$$P(3) = \frac{5!}{2!3!} \times 0.5^{(3)} \times 0.5^{(2)}$$

\Rightarrow Uniform distribution (continuous)

$[a, b]$

$$f(x) = \frac{1}{b-a} \quad a < x < b$$

The probability of any value within the range $[a, b]$ is same.

* uniform distribution (discrete)

all outcomes are equally likely

$$P(x) = \frac{1}{n} \rightarrow \text{total no.}$$

* confidence Interval

$$\bar{x} = 50 \rightarrow \text{point estimate} \rightarrow \mu$$

$$\begin{array}{c} 40 - 60 \rightarrow 20 \\ \text{---} \quad \text{---} \\ 5 \end{array} \quad \begin{array}{l} 2 \text{ batch} \\ - 50 + 5 \\ 45 - 55 = 50 \end{array}$$

$$40 - 50 - 60 + 10 - 10$$

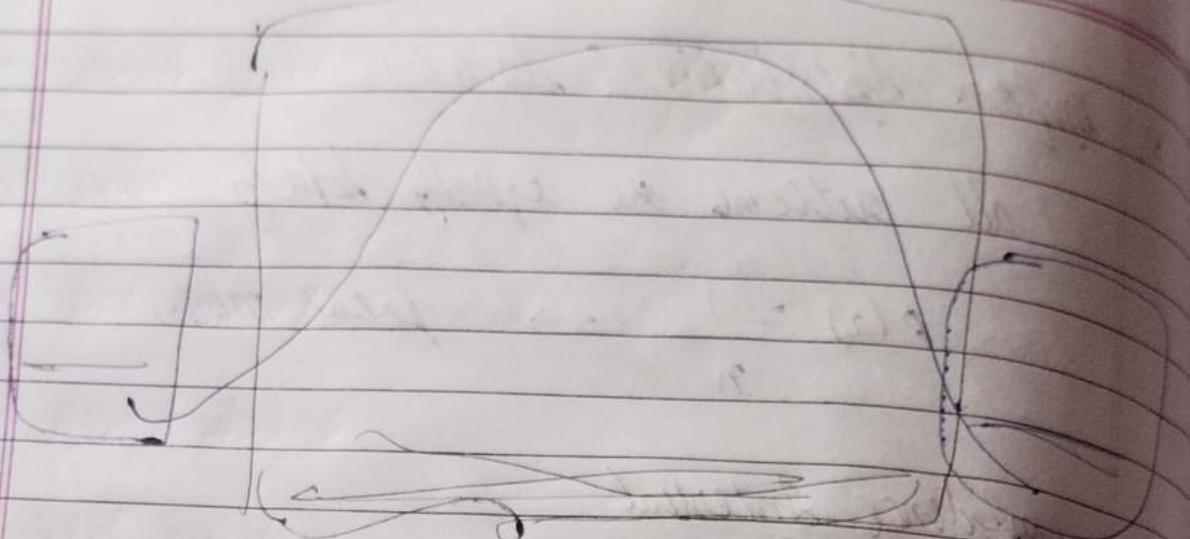
$$\begin{array}{c} +5 \\ - \\ \hline \end{array} \quad \begin{array}{c} \pm 10 \\ \text{---} \end{array} \quad \text{Margin of error}$$

50 \rightarrow point estimation

$$\begin{array}{c} 45 - 55 \\ \text{---} \\ 40 - 60 \end{array} \quad \text{confidence interval}$$

It is a range of values within which we expect a particular population parameter to fall.

Confidence Interval = point estimate \pm margin of error



$H_0 \text{ accept}$
 $10s \rightarrow X$

$H_0 \text{ reject}$

* Hypothesis testing

$$10 \times D \rightarrow [100 - \text{test com}]$$

10 1000

A Statistical hypothesis test is a method of statistical inference used to decide whether the data at hand is sufficient to support a particular hypothesis.

Hypothesis testing allows us to make probabilistic statements about population parameters.

Null hypothesis : H_0

The null hypothesis assumes that there is no significant relationship or effect b/w two variables [in simpler terms \rightarrow it says nothing new is happening].

It serves as a starting point for HT & represents 'stat' vs 'the assumption of no effect until proven otherwise'.

The purpose of HT is to gather evidence to reject or fact null hypothesis in favour of alternate hypothesis, which claims there is significant effect or relationship.

Alternate hypothesis H_a or H_1

It is a statement, that contradicts the NH & claims there is significant effect or relationship.

Rejection Region Method

H_0 & H_a

(1.05)

(5%)

0.5

$\alpha \rightarrow$ value

\hookrightarrow Significance level $\rightarrow 0.5\%$.

③ assumptions

$n_D \rightarrow n_1, n_2$

6. Teacher's Signature

4 decide test

2 - test

t - test

5 value

6 Test conduct

7 Reject / Accept

8 state results

Q1 50 → units per day

45 - 53

 $\sigma = 5$
Training

30 amb/ → 53 units per day.

1 $H_0: \mu = 50$ $H_a: \mu > 50$ 2 $\alpha = 0.05$ (5%)

one tailed test

3 data normal, σ , random
 $n = 30$

4 z-test

5 z-score = $\frac{x - \mu}{\sigma}$ prob.

$$\frac{x_i - \bar{X}}{\sigma/\sqrt{n}} = \frac{53 - 50}{\sigma/\sqrt{30}}$$

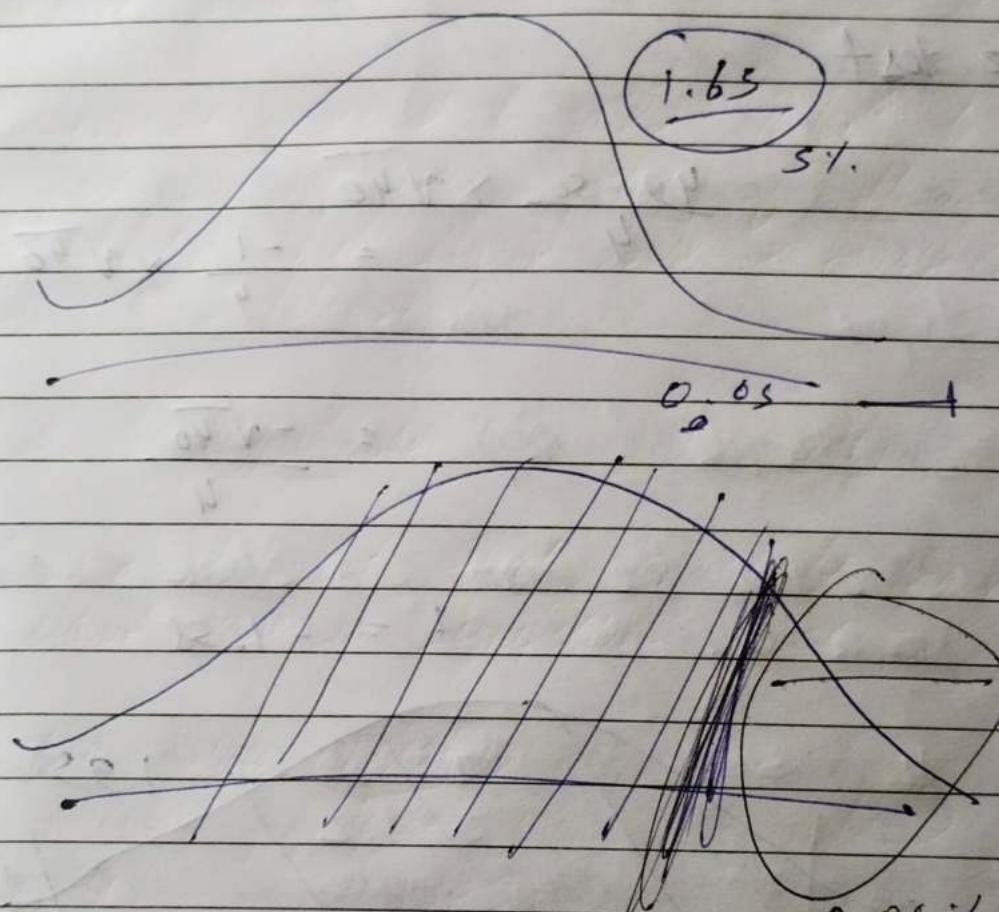
$\alpha = 51.$

$$= \frac{3}{5} \times \sqrt{30}$$

$$= \frac{9}{28} \times \frac{30^6}{5}$$

$$z = 3.28$$

(6)



7 reflection of H_0

8 $\mu > 50$

2 $\text{ally} - 50g$
 $\sigma = 4g$

$40g \rightarrow$
 $\pi M = 4g$

1 $H_0: \mu = 50g$

2 $\alpha = 0.05$

3 $n \geq 30$

z-test

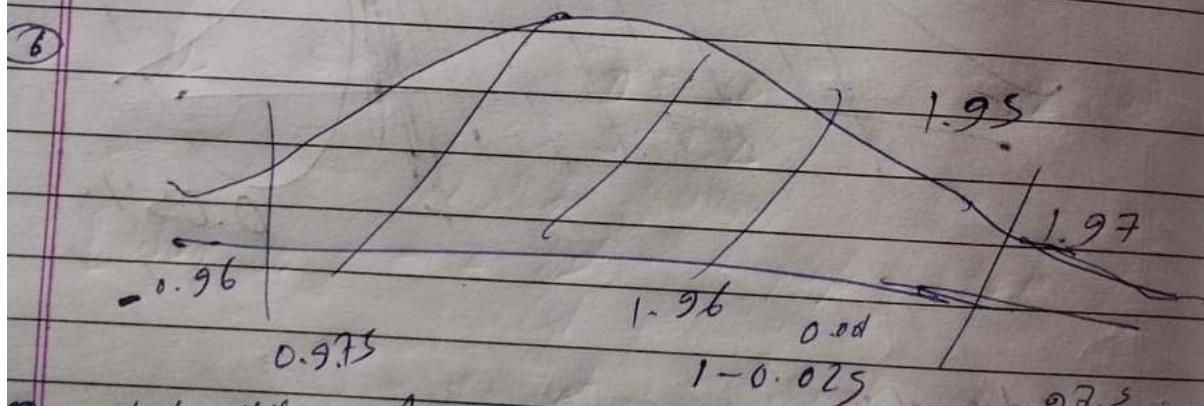
4 z-test

5

$$\frac{49 - 50}{4} \times \sqrt{40} \\ = -\frac{1}{4} \times \sqrt{40}$$

$$= \frac{-\sqrt{40}}{4}$$

1 = -1.58



⑨ Null Hypo Akuhi
 $H = 50g$

2 errors

Type 1

Type 2

H_0 true

H_0 false

reject H_0

Type - I

correct

accept H_0

correct

Type 2

Type - I

false +ve

4

H_0 reject $\rightarrow H_0$ true (correct)

rejecting H_0 when H_0 is actually correct

Type - 2 \rightarrow false - ve

accept H_0 when H_0 is actually incorrect

$P > 0.05 \xrightarrow{2} \text{Null accept.}$

$P < 0.05 \xrightarrow{} \text{Null reject.}$

$P < 0.01 \xrightarrow{} \text{Strong evidence}$

$0.01 \leq P < 0.05 \xrightarrow{} \text{Moderate evi.}$

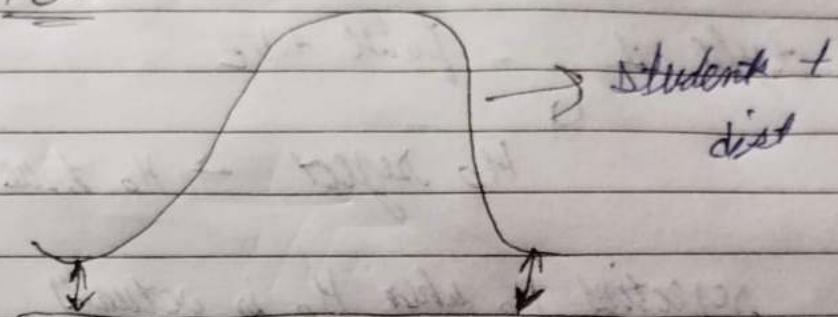
$0.05 \leq P < 0.1 \xrightarrow{} \text{Weak evidence.}$

$P \geq 0.1 \xrightarrow{} \text{No evi.}$

$P < \alpha$ → H_0 reject
 $P > \alpha$ → H_0 accept

P value → It is a measure of the strength of the evidence against the null hypothesis.

T - Test



3 types

One-sample t-test

Compares the mean of a single sample to a known μ .

$$\bar{M} = 50 \quad n = 30 \quad \bar{x} = 49.72 \\ s = 1.2$$

[when we $s \rightarrow z$ -test (prob → standard dev)]

$s + z$ -test] (sample std.)

$$\textcircled{1} \quad \mu = 50.9$$

$$n = 25$$

$$\text{std} = 1.2 \text{ kg}$$

$$\bar{x} = 49.7 \text{ kg}$$

$$H_0: \mu = 50$$

$$H_1: \mu \neq 50$$

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{49.7 - 50}{1.2/\sqrt{25}} = \frac{0.3}{1.2} \times 5$$

$$= -1.5$$

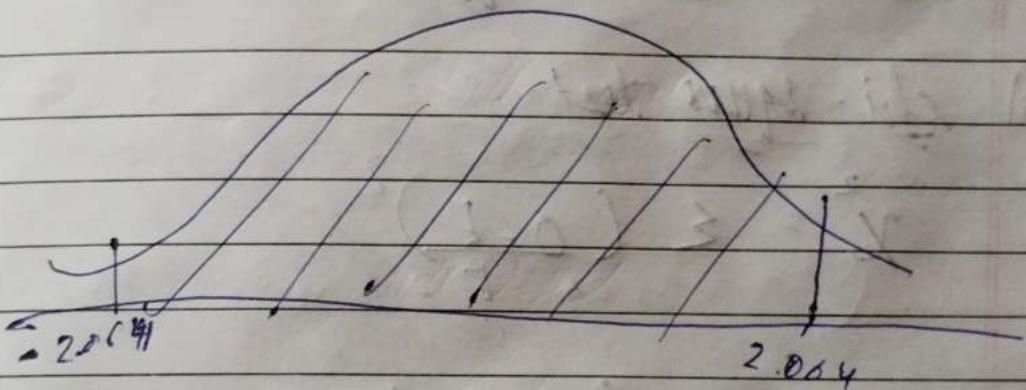
1.2

$$= -1.25$$

df = degree of freedom

$$= n - 1 \quad df = 24$$

$$t_{\text{critical}} = 2.064$$



H_0 accept

(sig)

2 Independent Two Sample t-test

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \rightarrow \text{standard error}$$

3 Paired t-test

$$t = \frac{\bar{d}}{s_d / \sqrt{n}}$$

$$\bar{d} = 2 \rightarrow \text{mean difference}$$

$$s_d = \text{std} \rightarrow \text{diff}$$

\Rightarrow chi-square test

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$O \rightarrow$ observed frequency

$E \rightarrow$ expected frequency
grand total

$$d = (r-1) \times (c-1)$$

r = no. of rows

c = no. of columns

	satisfied	not satisfied	Total
high school	50	70	120
college	90	60	150
PG	<u>20</u>	<u>10</u>	<u>30</u>
	<u>160</u>	<u>140</u>	<u>300</u>

EF	S	NT	Total
----	---	----	-------

$$\text{HS} \quad \frac{120 \times 160}{300} = 64 \quad \frac{120 \times 140}{260} = 56$$

$$\text{C} \quad \frac{150 \times 160}{360} = 80 \quad \frac{180 \times 140}{300} = 70$$

$$\text{PG} \quad \frac{160 \times 30}{300} = 16 \quad \frac{140 \times 30}{300} = 14$$

S	64	56	120
C	80	70	150
PG	16	14	30

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

$$S.S = \frac{(50-64)^2}{64} = \frac{196}{64} = 3.06$$

$$S.N.S = \frac{(70-56)^2}{56} = \frac{196}{56} = 3.5$$

$$C.S = \frac{(90-80)^2}{80} = \frac{100}{80} = 1.25$$

$$C.N.S = \frac{(60-70)^2}{70} = \frac{100}{70} = 1.42$$

$$P.W.S = \frac{(20-16)^2}{16} = \frac{16}{16} = 1$$

$$P.G.N.S = \frac{(10-14)^2}{14} = \frac{16}{14} = 1.14$$

$$\chi^2 = 3.06 + 3.5 + 1.25 + 1.42 + 1 + 1.14 \\ = 11.37$$

$$df = (3-1) \times (2-1)$$

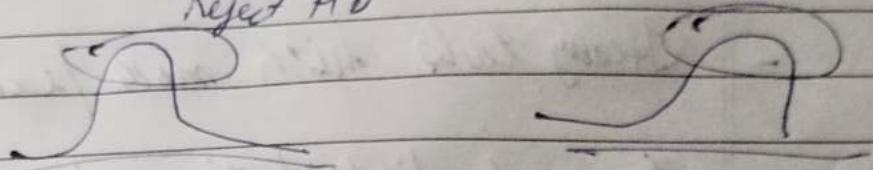
$$= 2$$

$$\alpha = 0.05$$

$$\chi^2_{\text{critical}} = 5.991$$

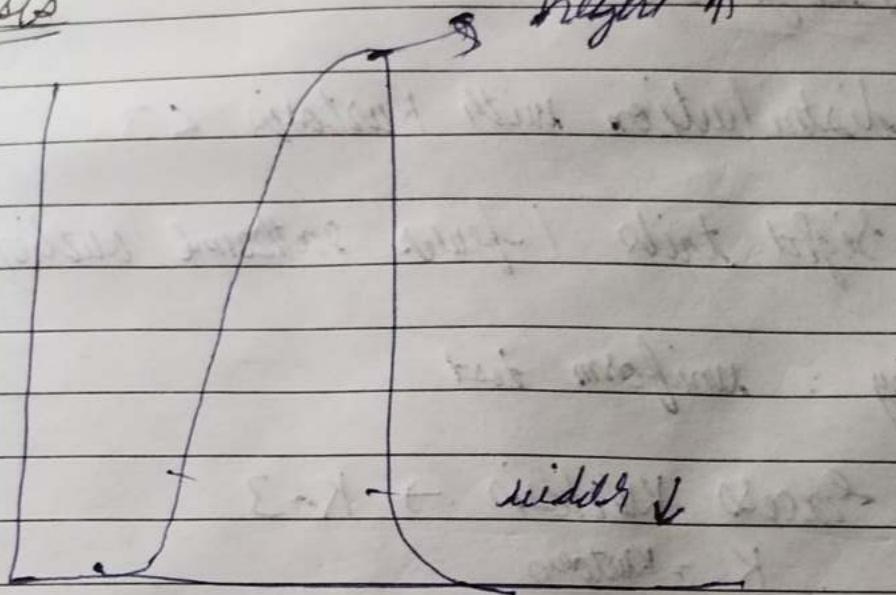
$$[\chi^2 = 11.37]$$

Reject H₀



Kurtosis

height ↑



- Kurtosis measures "tailedness" of a distribution or how extreme the outliers are

3. Mesokurtic

- Tails are similar to N.P

- distribution with kurtosis ≈ 3

e.g.: Standard Normal distribution

2. Leptokurtic

- A distribution with kurtosis > 3
- Heavy tails with more extreme outliers
Ex - t-distribution with very small σ_f

3. Platykurtic

- distribution with kurtosis < 3
- Light tails (fewer extreme outliers)
- Ex : uniform dist

Excess kurtosis $\rightarrow K - 3$
 $K = \text{kurtosis}$

$EK > 0 \rightarrow \text{Leptokurtic}$

$EK < 0 \rightarrow \text{Platykurtic}$

$$K = \frac{n \cdot \sum (x_i - \bar{x})^4}{\left(\sum (x_i - \bar{x})^2 \right)^2} \cdot \frac{3}{n}$$

n = no. of observations

x_i = each data point

\bar{x} = Mean

- high kurtosis (leptokurtic)
- More extreme outliers
- higher likelihood of rare, extreme values.
- by financial returns during market crisis / crashes

Lower kurtosis (Platykurtic)

- fewer extreme outliers
- data is evenly spread.
- kurtosis near 3 (mesokurtic)

Similar to normal distribution

