

1. Stats

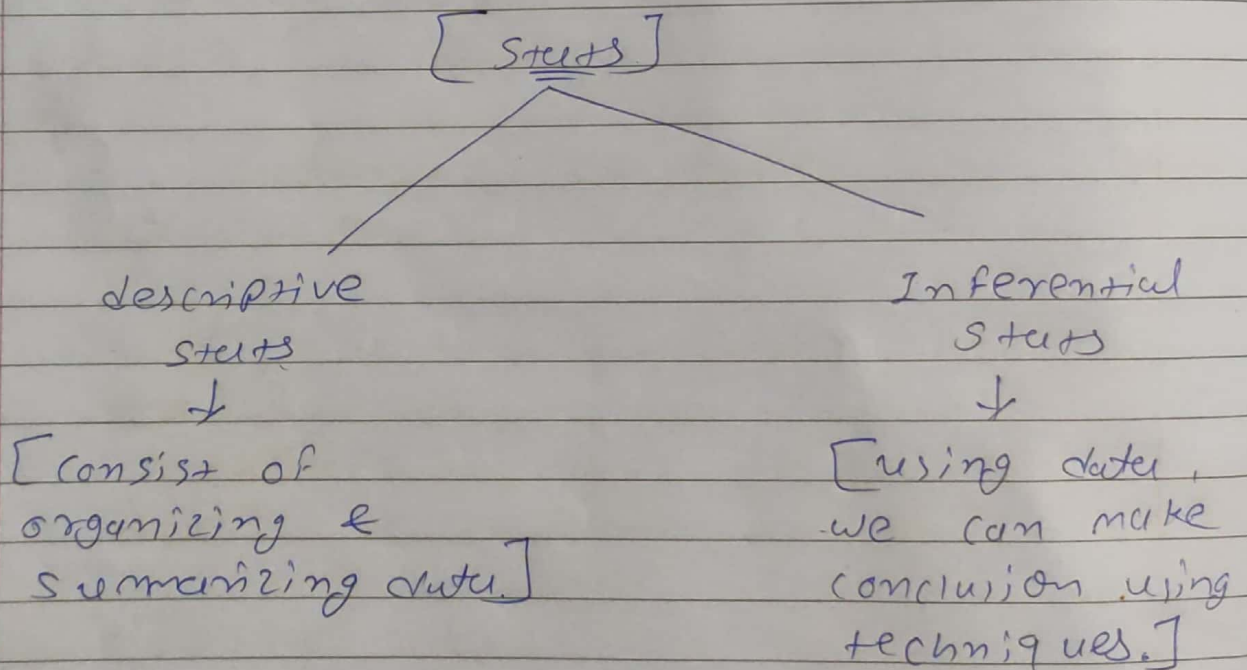
* what is stats?

→ stats is a science of collecting, analyzing & organizing data.

* Data :

→ Facts or pieces of information that can be measured.

* Types of Stats



2.

* -) Sample & population

Population, N / sample, n
↓

-) whole dataset is known as population.

* sampling techniques

1. simple Random sampling

2. stratified sampling

3. systematic sampling

4. convenience sampling

* variables

-) It is a property that can hold or take only ~~only~~ any value.

ex //

Age: { 8, 20, 5, 7, 600, 8078, ... }

3.

marks i. { 76, 80, 95, ... }

* Types of Variables

Variables

1) Qualitative

2) Quantitative

→ categorical values

→ numerical value

Ex

IQ: 0-10 → low

10-50 → Avg

50-70 → good

Ex

Height i. { 16.2, 15.9 }

weight i. { 59, 61 }

4. There are two types of Quantitative.

2) Quantitative

① Discrete (int)

→ whole no.

Continuous (float)

→ decimal no.

ex

→ no. of students

→ no. of peoples.

ex

→ Height

{ 165.2, 167.9, ... }

→ weight

{ 165.5, 50.9, ... }

* Variable measurement scales.

① ordinal :- ordered [rank, graduation]

② Nominal :- categorical values [colors, classes]

③ Interval :- [no zero / abs datapoint]

④ Ratio :- zero means nothing

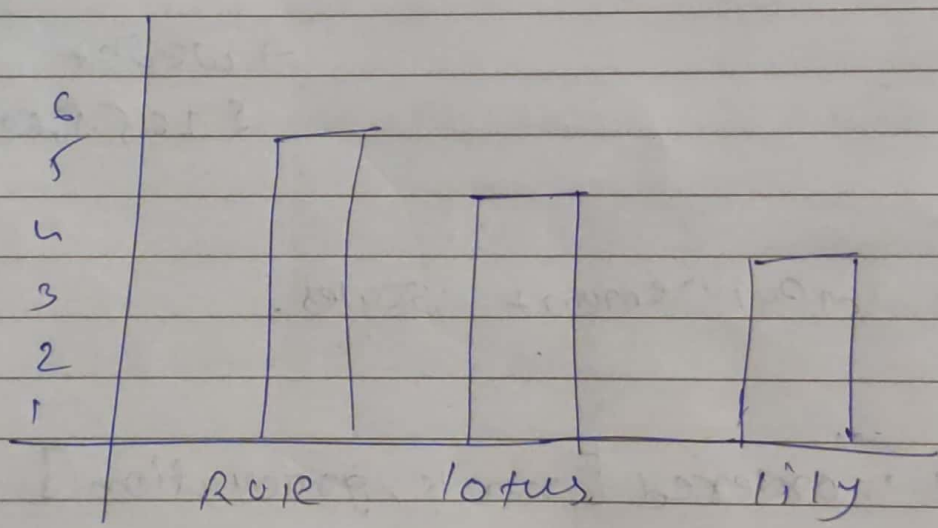
5.

* frequency \Rightarrow

ex

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

Flowers	F	CF
Rose	5	5
lily	3	8
lotus	4	12
	12	



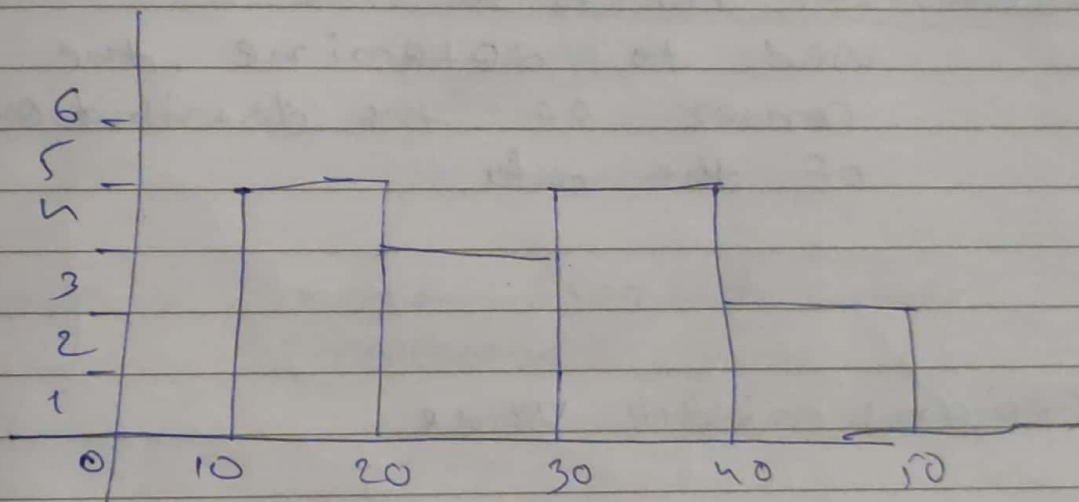
bar graph
chart

Histogram

marks: [12, 15, 17, 18, 21, 22, 27, 31, 34, 36, 39, 42, 45]

bins

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



7.

measure of central Tendency

Avg \rightarrow mean

pop

sum

$$\mu = \frac{\sum x_i}{N}$$

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

\rightarrow mean: It refers to the measure used to determine the center of the distribution of the data.

\rightarrow median: middle value

\rightarrow as ascending order

\rightarrow data

even odd

\downarrow

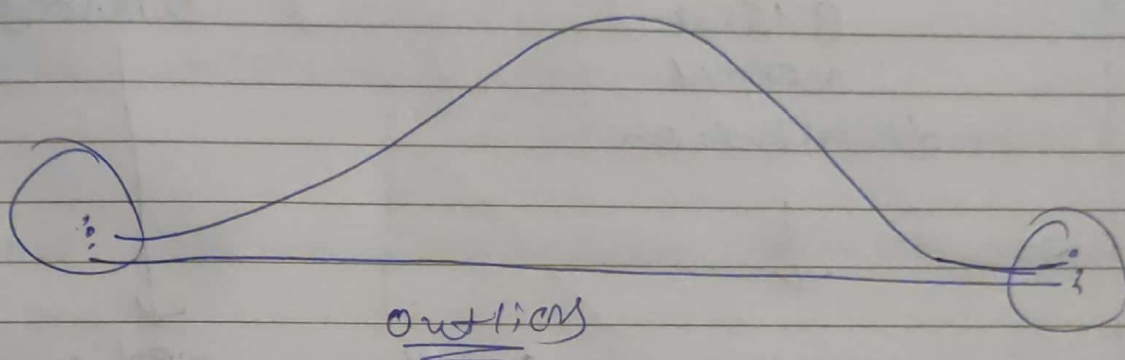
$$\frac{\left(\frac{n}{2}\right)^{th} + \left(\left(\frac{n}{2}\right)^{th} + 1\right)^{th}}{2}$$

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

8.

* outlier: a datapoint who doesn't follow pattern or trend of the data set then it is considered as outlier.

[Case extreme points]



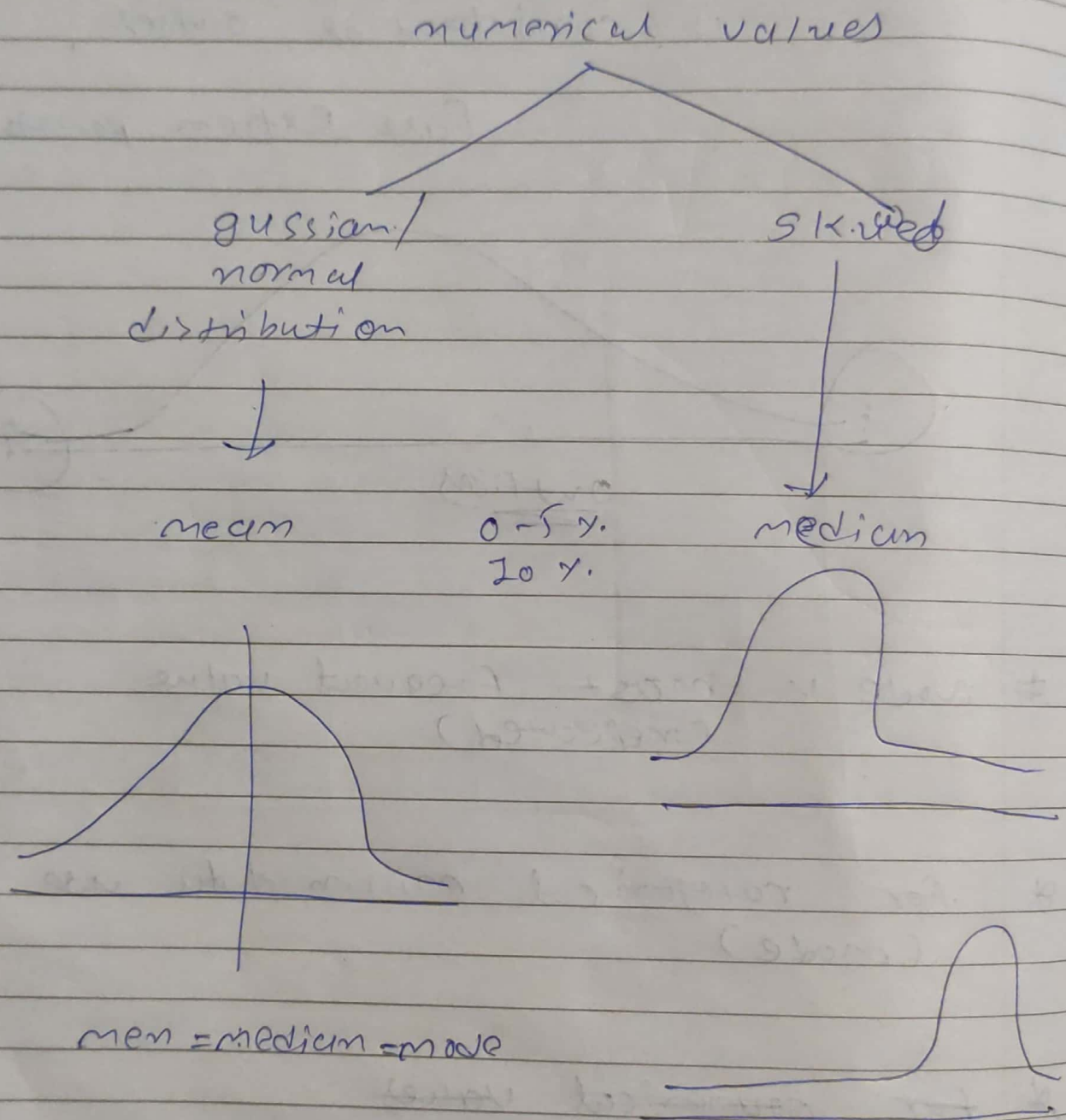
* Mode is most frequent value (repeated)

* for categorical missing data use (mode)

* for numerical values

9.

* for numerical values.



* Variance

→ It measures how ~~the~~ far the numbers in a dataset are from the mean (avg)

→ High variance → more spread (far from the mean)

→ low variance → closer ~~for~~ to mean

pop

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

summation

sample

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

x_i = every data pt

μ = mean of pop

N = total pop

Σ = summation

11.

DOMS

Page No.

Date

/

/

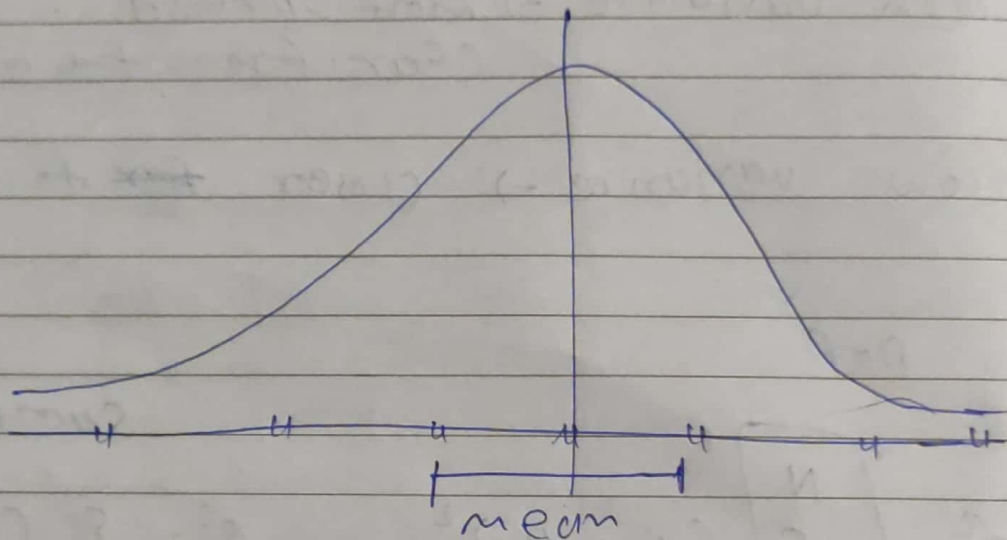
* standard deviation

pop

sample

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

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



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

q. 5 y.

99.7%

$$\begin{aligned}\mu \pm 1\sigma &= 68\% \\ \mu \pm 2\sigma &= 95\% \\ \mu \pm 3\sigma &= 99.7\%.\end{aligned}$$

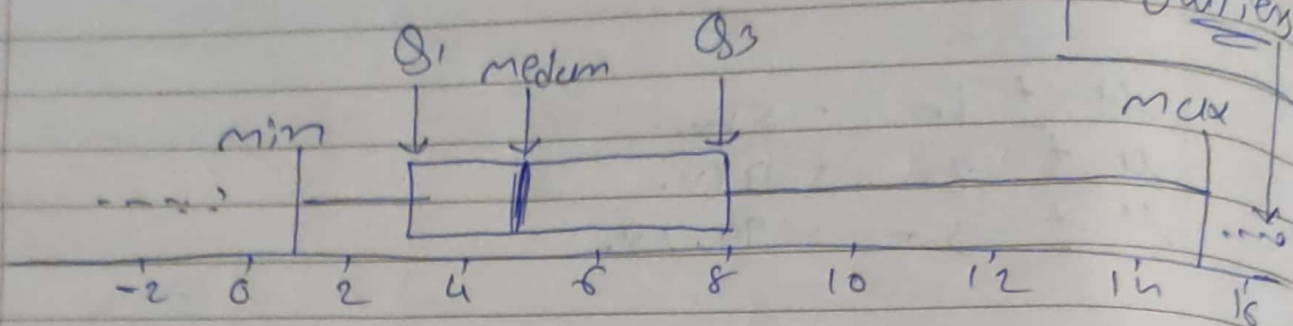
→ square root of variance

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

* Five number summary

- ① minimum Q_0
- ② 25 percentile → First Quartile $[Q_1]$
- ③ median → 50 percentile Q_2
- ④ 75 percentile → Third Quartile $[Q_3]$
- ⑤ maximum Q_4

13.



Box plot

* Z-score

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

$$\mu = 0$$

$$\sigma = 1$$

* Normalization

$$X_n = \frac{X - x_{\min}}{x_{\max} - x_{\min}}$$

14.

* Exponential distribution

- It describes the time b/w events in a process where events occur independently & at a constant rate λ .

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

$$x \geq 0$$

* Bernoulli distribution (discrete)

- It models a single experiment with two possible outcomes.

Success ($x=1$) ; Failure ($x=0$)

* Uniform distribution (continuous)

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

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

15.

* uniform distribution (discrete)

all outcomes are equally likely

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

Hypothesis testing

→ 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.

Rejection Region method

- ① H_0 & H_1
- ② α - value
- ③ assumptions
- ④ derive test value $\left\{ \begin{array}{l} \rightarrow Z\text{-test} \\ \rightarrow T\text{-test} \end{array} \right.$
- ⑤
- ⑥ Test conduct
- ⑦ Reject / Accept
- ⑧ state results

16.

* 2 types of errors

	Type 1 ✓ H_0 true	Type - 2 ✗ H_0 False
reject H_0	Type-1	correct
accept H_0	correct	Type-2

⇒ Type-2 False +ve

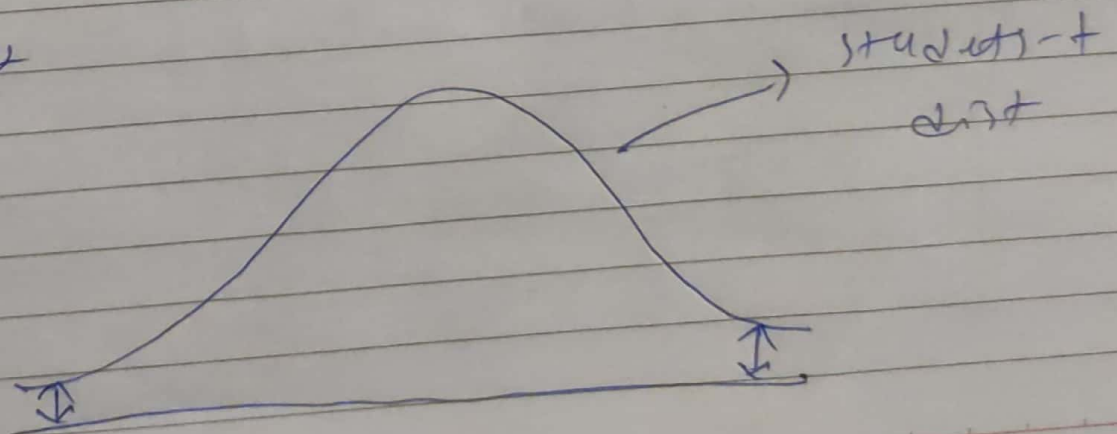
↳ H_0 reject $\rightarrow H_0$ true (correct)

∴ reject H_0 when H_0 is actually correct

⇒ Type-2 False -ve

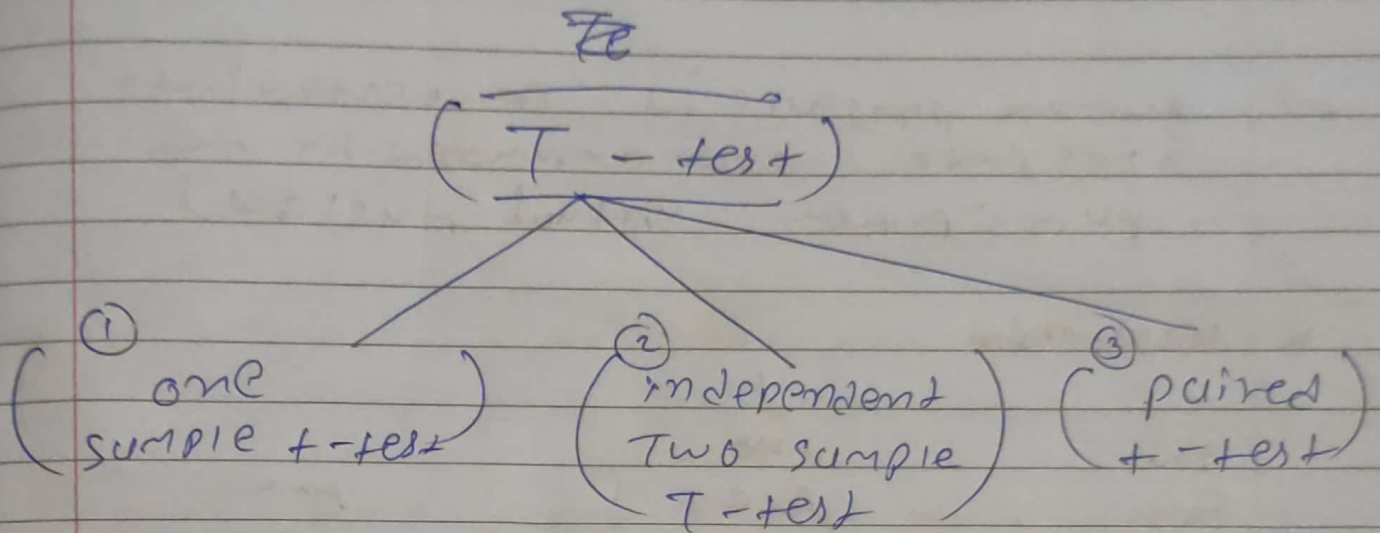
∴ accept H_0 when H_0 is actually incorrect

t-test



17.

→ 3 types



① one sample t-test

→ compares the mean of a single sample to a known μ

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

② Independent Two sample t-test

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \leftarrow \begin{matrix} \text{Standard Error} \\ \text{Difference} \end{matrix}$$

③ Paired t-test

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

$\bar{d} \Rightarrow 2 - 7$ mean diff

$sd \Rightarrow std \rightarrow diff$

~~18~~ 19.

Chi-square test

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

O \Rightarrow observed frequency

E \Rightarrow Expected frequency

$$E = \frac{\text{Row Total} \times \text{column total}}{\text{Grand Total}}$$

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

r = no. of rows

c = no. of cols