

Introduction to Statistics

Defn: Stats is the science of collecting, organizing and analysing data.

Data: Facts or pieces of information

Eg: Heights of students in classroom
Salary of people in society.

⇒ Types of Stats

- ① Descriptive Stats
- ② Inferential Stats

→ ① Descriptive Stats

① It consists of organizing, summarizing, and visualizing data.

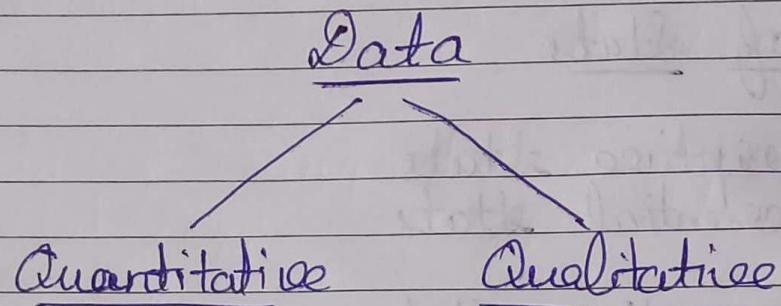
- ① Measure of Central Tendency
- ② Measure of dispersion (V, σ , Z-score)
- ③ Different types of distribution.

Eg: Histogram, pdf, pmf, Gaussian, log normal, Exponential, Binomial, Bernoulli, Poisson.

→ ② Inferential stats

- It consists of using data you have measured to form conclusion.

- ① Z-test
- ② T-test
- ③ Chi square test
- ④ Anova test



⇒ Sampling Techniques

① Random Sampling

- Most common in behavioral research
- Equal chance of being selected in the sample.

② Stratified Sampling

- Dividing the population into subgroups or strata based on certain characteristics or attributes.

- ① w.r.t age, gender, income etc.
- ② Convenience sampling
- ③ Participants are selected based on availability and willingness to take part.
- ④ Cluster sampling
 - ① Divide a population into clusters
 - ② Eg: Districts, schools and randomly selecting the clusters
- ⑤ Snowball sampling
 - ① Existing study subjects are used to recruit more subjects into the sample.
- ⑥ Purposive sampling
 - ① Selecting samples based on the judgement of the survey taker or researcher.
- ⑦ Systematic sampling
 - ① It is a probability sampling method where researchers select ~~number~~ members at a regular interval.

→ Scale of Measurement

- 1) Nominal scale data
- 2) Ordinal scale data
- 3) Interval scale data
- 4) Ratio scale data

1) Nominal scale data:

- ① It is a type of qualitative data that labels variables without using numbers.
- ② Eg: Mode of transportation, genotype, blood type, zip code, gender, place, etc...

2) Ordinal scale data:

- ① It is the 2nd level of measurement that separates the ordering and ranking of data without establishing the degree of variation between them.
- ② Eg: Ranking of school students, Rating of services in restaurants.

3) Interval scale data:

- It is a type of quantitative data that measures values on a scale with equal distances between each value.
- Eg: Temperature in Fahrenheit or Celsius, pH measures, IQ and SAT scores, etc.

4) Ratio scale data:

- It is a type of quantitative data that measures variables on a scale with equal intervals between values and a true zero.
- Eg: Height, money, age, weight, etc..

⇒ Random Variables

- A Random variable is a variable in statistics that assigns numerical values to the outcomes of sample space. The possible values of a random variable depend on the outcomes of a random phenomenon.

Variable (x, y)

$$x + 6 = 8$$

$$x = 2$$

$$x = 2$$

$$y = 6$$

$$8 = y + x$$

- Random variable is a process of mapping the output of a random process or experiment to a number.

- Ex: \rightarrow Toss a coin

$$X = \begin{cases} 0 & \text{if heads} \\ 1 & \text{if tails} \end{cases}$$

- \rightarrow Rolling a dice

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

\Rightarrow Covariance & Correlation

- It is the measure of the relationship between two numbered random variables. It measures how much the variable change together, or the variance between them.

- It measures the direction of a relationship between two variables.

\rightarrow Population Covariance Formula

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

→ Sample Covariance

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

① Advantages:

→ +ve or -ve value

→ Shows the relationship b/w two variables positive / negative.

① Disadvantages:

→ Doesn't have a specific limit values.

② Pearson Correlation Coefficient

$$\rho = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2] \cdot [n\sum y^2 - (\sum y)^2]}}$$

① Pearson Correlation Coefficient doesn't work with non-linear data. So we will use Spearman Correlation.

③ Spearman's Rank Correlation Coefficient

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

R = Rank Based on frequency (order) value

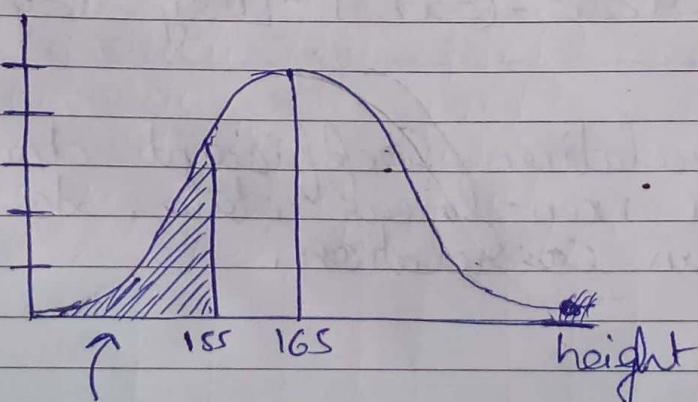
- If 2 features are highly correlated, it's ok to drop one of the features.

→ Probability Distribution Function

① Probability density function

② Continuous function

③ Eg: Age, Height
[float]



$$P(H \leq 155)$$

$$P(H \leq 155, H \geq 175)$$

② Probability mass function

③ Discrete values

④ Eg: No. of bank acc, Rolling a dice
int

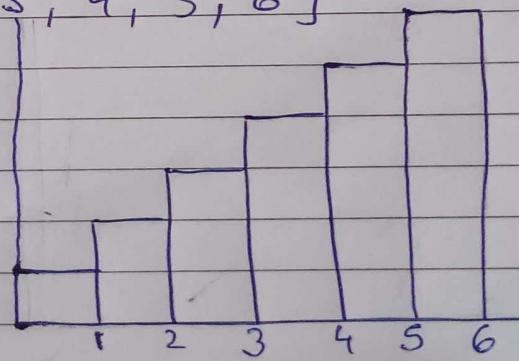


$$Pr(x \leq 4) = Pr(x=1) + Pr(x=2) + Pr(x=3) + Pr(x=4).$$

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

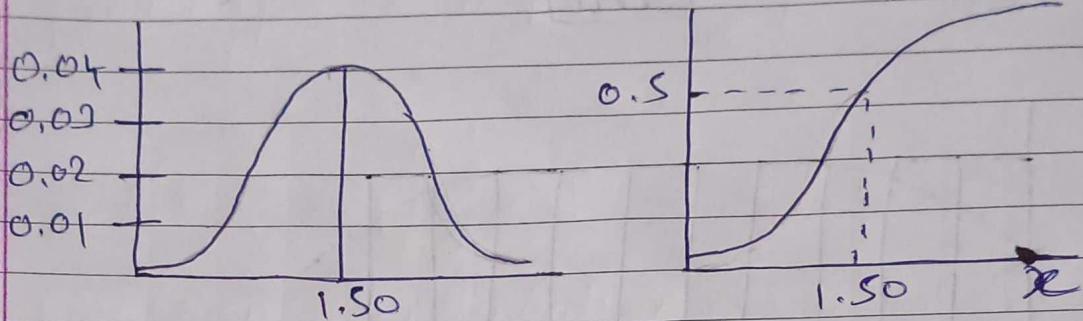
③ Cumulative Distribution Function

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



- Adding the previous to the next value.

pdf & cdf



- Probability Density of pdf \Rightarrow Gradient or derivative of cdf.

pdf is the derivative of cdf
cdf is the integration of pdf

\Rightarrow Types of Probability function

- 1) Normal dist (pdf)
- 2) Bernoulli dist (pmf) [Binary outcome]
- 3) Uniform dist (pmf)
- 4) Poisson dist (pmf)
- 5) Binomial dist (pmf)
- 6) Log Normal dist (pdf)

→ Bernoulli Distribution

- Discrete Random Variable (pmf)
- Outcome are Binary

Eg: → Tossing a coin $\{H, T\}$.

$$P_r(H) = 0.5 = p$$

$$P_r(T) = 0.5 = q$$

→ Whether a person will pass / fail.

$$P_r(\text{pass}) = 0.7 = p$$

$$P_r(\text{fail}) = 1-p = 0.3 = q$$

→ Binomial Distribution

- Discrete Random Variable

- Every Experiment outcome is Binary

- It represents the probability for x success in n trials, given a success probability p for each trial.

- Eg: Tossing a coin 10 times.

notation $\rightarrow B(n, p)$

$$P_r(H) = 0.5 = p$$

$$P_r(T) = 0.5 = q$$

parameter $\rightarrow n \in \{0, 1, 2, 3, \dots\} \Rightarrow$ no. of trials

$p \in \{0, 1\} \rightarrow$ success probability
 $q = 1-p.$

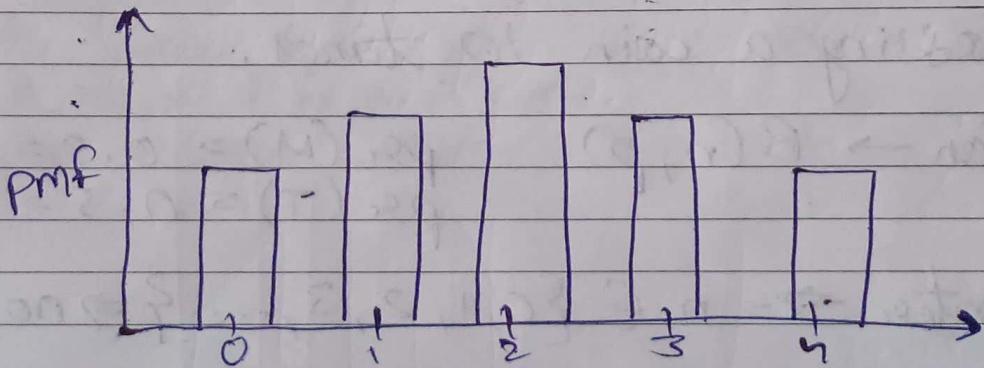
support = $k \in \{0, 1, 2, \dots, n\} \rightarrow$ no. of success

pmf: $P(k, n, p) = {}^n C_k p^k (1-p)^{n-k}$

$${}^n C_k = \frac{n!}{k!(n-k)!}$$

→ Poisson Distribution.

- ① Discrete Random Variable (pmf)
- ② Describes the no. of elements occurring in a fixed time interval.
- ③ Eg: - no. of people visiting hospital every hour
 - no. of people visiting banks, at 11 am.



$\lambda = 3$ = Expected no. of event occurs at every time interval.

Q. What is the probability of no. of people visiting at 3 pm. $\Rightarrow P(X=3)$

$$= \frac{e^{-3} 3^3}{5!} = 0.101 \approx 10\%$$

10% of the people visits the bank at 3 pm.

3) Empirical Rule of Normal dist

$$\begin{array}{ccc} 68 & 95 & 99.7\% \\ 1^{\text{st}} \text{ sd} & 2^{\text{nd}} \text{ sd} & 3^{\text{rd}} \text{ sd} \end{array}$$

⇒ Uniform Distribution

1) Continuous Uniform Distribution

0 A Continuous uniform distribution is a probability distribution that takes values within a specified range. It is defined by the two parameters, a, b where a is the lower limit and, b is the upper limit.

Notation = $U(a, b)$

Parameters = $-\infty < a < b < \infty$

pdf = $\begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$

cdf = $\begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } x \in [a, b] \\ 1 & \text{for } x > b \end{cases}$

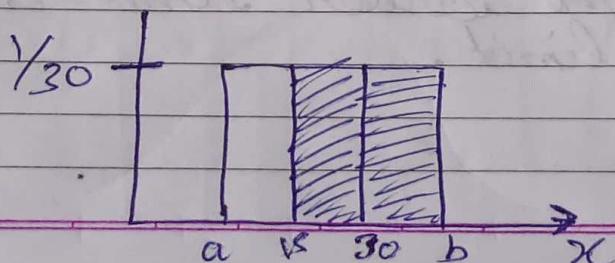
$\frac{x-a}{b-a}$ for $x \in [a, b]$

Mean = $\frac{a+b}{2}$

Median = $\frac{a+b}{2}$

① Example: The no. of candies sold at a shop is uniformly distributed with a max of 40 and a min of 10.

= Probability of daily sales to fall b/w 15 and 30.



$$\begin{aligned} a &= 10 \\ b &= 40 \\ P(15 \leq x \leq 30) &= \end{aligned}$$

$$\begin{aligned} &= (x_2 - x_1) \times \frac{1}{b-a} \\ &= 30 - 15 \times \frac{1}{40 - 10} \\ &= 15 \times \frac{1}{30} = 0.5 \end{aligned}$$

$$= \Pr(X \geq 30)$$

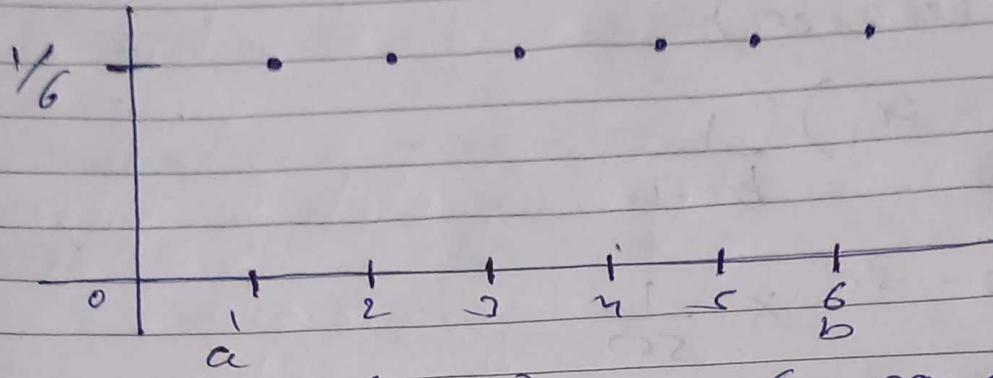
$$(x_2 - x_1) \frac{1}{b-a}$$

$$40 - 30 \times \frac{1}{50}$$

$$\frac{10}{30} \Rightarrow [0.33]$$

2) Discrete Uniform Distribution

- A discrete uniform distribution is a statistical distribution where the probability of outcomes is equally likely and with finite values.
- In a discrete uniform dist, everyone of n value has equal probability $\frac{1}{n}$.
- Eg: Rolling a dice
 - Selecting a card from deck of card
 - Tossing a fair coin.
- Eg: Rolling a dice



Notation: $U(a, b)$ $n=6$ no. of outcomes
 $[n = b-a+1]$

$$P(I) = \frac{1}{n} = \frac{1}{6} \quad \text{Mean} = \frac{a+b}{2}$$

Parameters: a, b $[b \geq a]$ Median = $\frac{a+b}{2}$

$$\text{pmf} = \frac{1}{n}$$

⇒ Standard Normal Distribution

→ Z-score

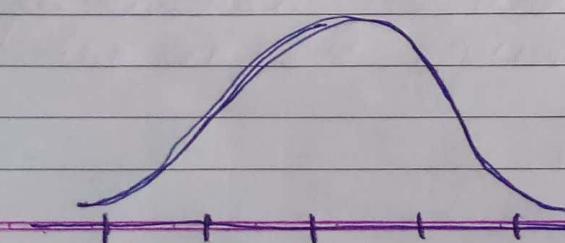
→ Z-stats

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

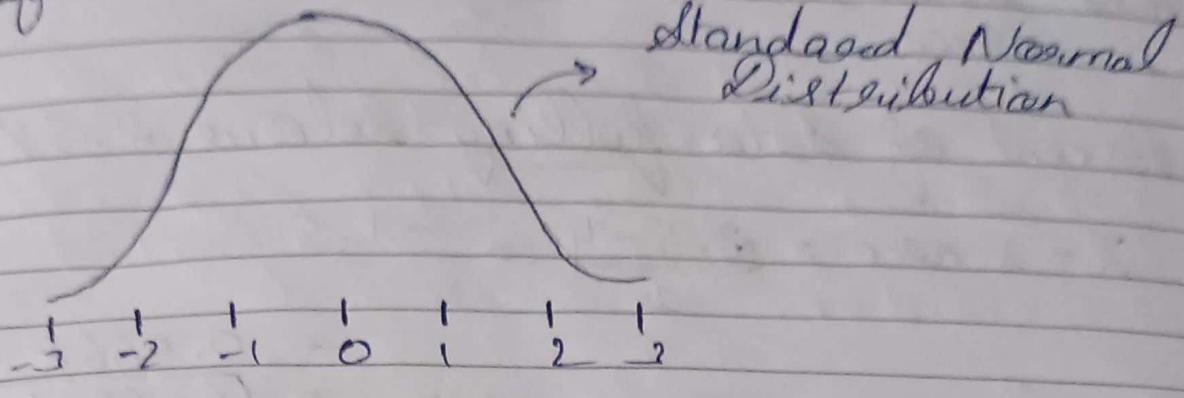
Normally distributed

$$\mu = 3$$

$$\sigma = 1.414$$

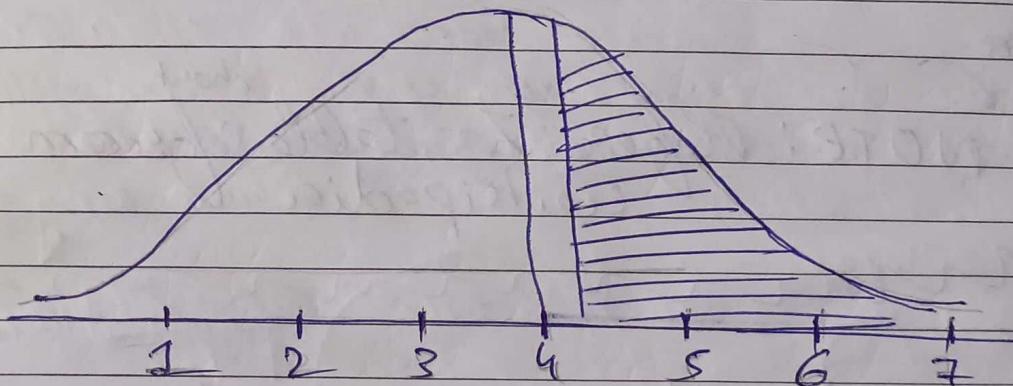


if $\mu = 0$ & $\sigma = 1$



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

Z-table - We use this to find out the area under the curve.



What % of the data is falling down ~~to~~ 4.5

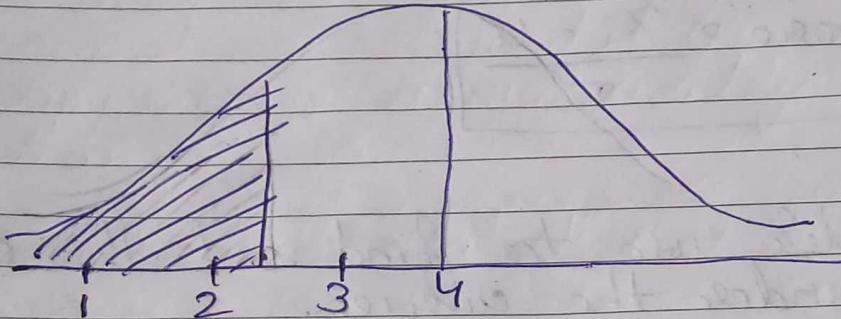
$$Z\text{-score} = \frac{x_i - \mu}{\sigma} = \frac{4.5 - 4}{1} = 0.5$$

$$\Rightarrow 0.6915$$

Area under the curve $\rightarrow 1 - 0.6915 = 0.3085 = 30.85\%$

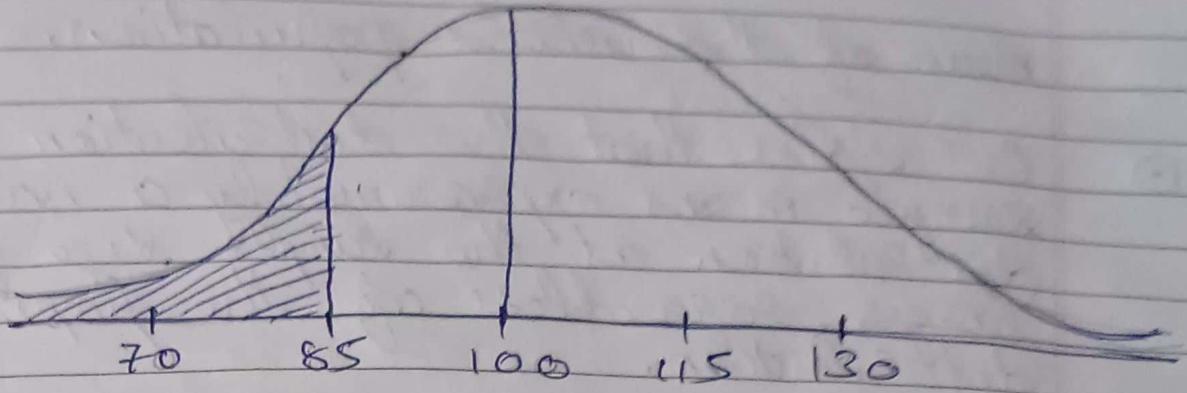
② Percent of data falling below 2.5%.

$$\begin{aligned} \text{Z-score} &= \frac{2.5 - 4}{1} = -1.5 \\ &= 0.0668 \\ &= 6.68\% \end{aligned}$$



~~NOTE~~ NOTE: Refer Z-table ^{chart} from any [wikipedia](https://en.wikipedia.org) site.

In India the avg IQ is 100, with σ of 15. What is the percentage of the population would you expect to have an IQ lower than 85%.



$$Z = \frac{85 - 100}{15} = \frac{-15}{15} = -1 = 0.1587 \Rightarrow 15.87\%$$

Area under the curve $> 85\%$.

$$1 - 0.1587 = 0.8413 = 84.13\%$$

Area between 85 & 100

$$= 0.5 - 0.1587 = 0.3413 = 34.13\%$$

⇒ Central Limit Theorem

- CLT is a statistical theory that states that when a large sample size has a finite variance, the samples will be normally distributed.
- The theorem also states that the mean of the samples will be approximately equal to the

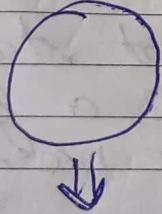
mean of the whole population.

- CLT states that the distribution of sample means approximates a normal distribution as the sample size gets larger, regardless of the population's distribution.
- The theorem holds true regardless of whether the source population is normal or skewed, provided the sample size is sufficiently large than $n=30$.

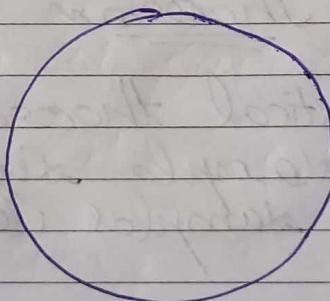
* Inferential Statistics

- Point Estimate: It is an observed numerical value used to estimate an unknown ~~parameter~~ population parameter. (single numerical value)

Sample

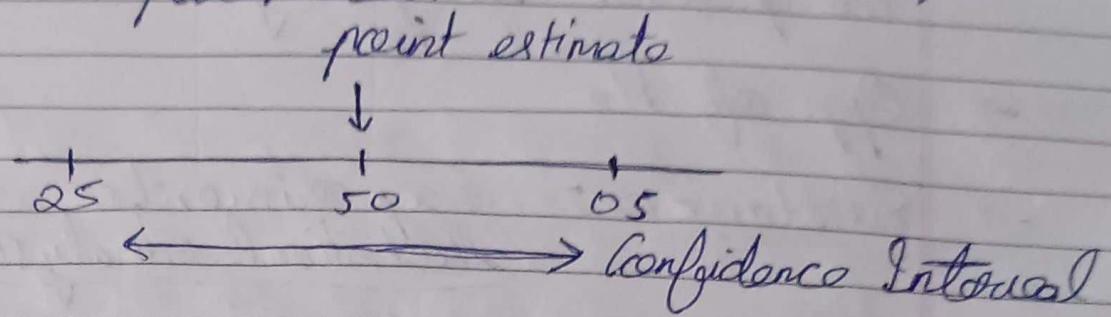


population



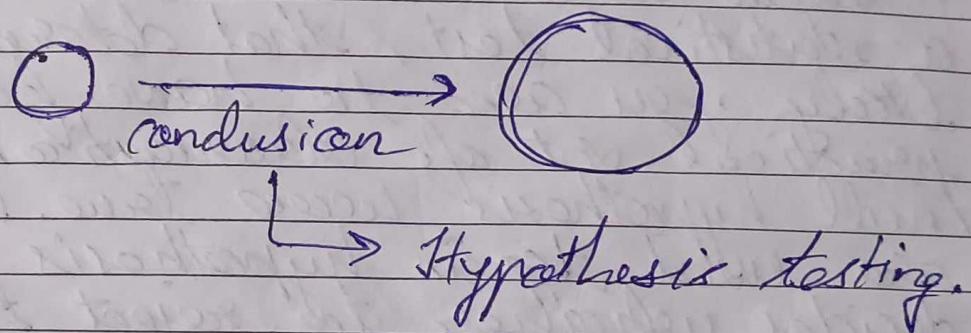
$\bar{x} \leftarrow$ Point estimate $\rightarrow \mu \rightarrow$ population mean

2) Interval Estimate: Range of values used to estimate the unknown population parameter.



⇒ Hypothesis & Hypothesis Testing Mechanism

Inferential stats → Conclusion



→ Hypothesis testing mechanism
person claims →

① Null hypothesis (H_0)

→ The person is not guilty.

→ Assumptions you are ~~not~~ beginning with

②

Alternate Hypothesis (H_1)

The person is guilty \rightarrow

\rightarrow Opp of H_0

\rightarrow Evidence - Experiments
Statistical Analysis

Z test, t test, ANOVA test etc.

\rightarrow P-value

① P value is a number calculated from a statistical test, that determine how likely you are to have found a practical set of observations if the null hypothesis were true. P values are used to in hypothesis testing to decides whether to reject the null hypothesis.

② Example: Coin is fair or not

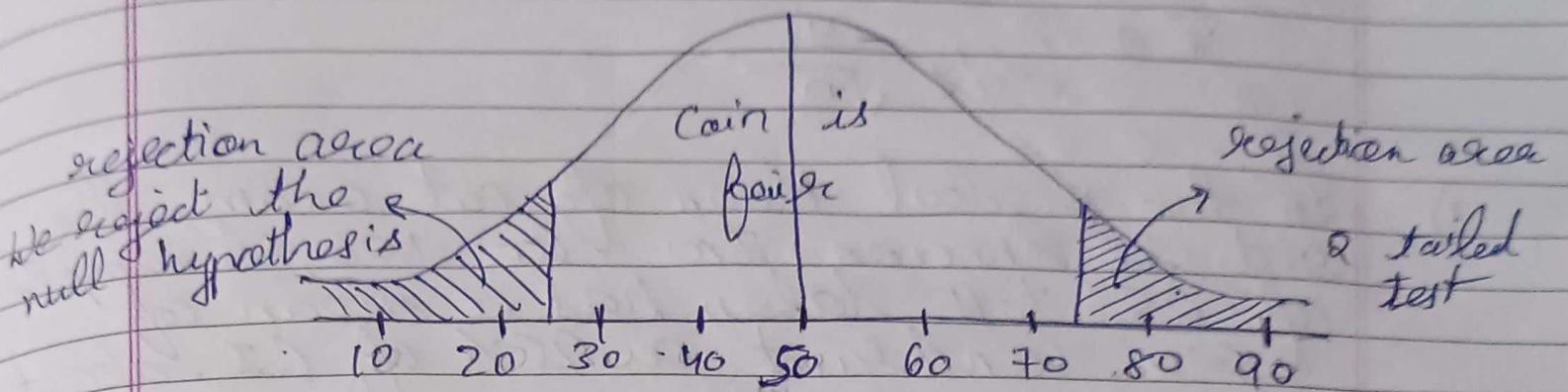
$\{H_1, H_0\}$.

\rightarrow Hypothesis Testing

1) Null hypothesis $H_0 \rightarrow$ Coin is fair

2) Alternate hypothesis $H_1 \rightarrow$ Coin is not fair.

3) Experiment \rightarrow Testing.

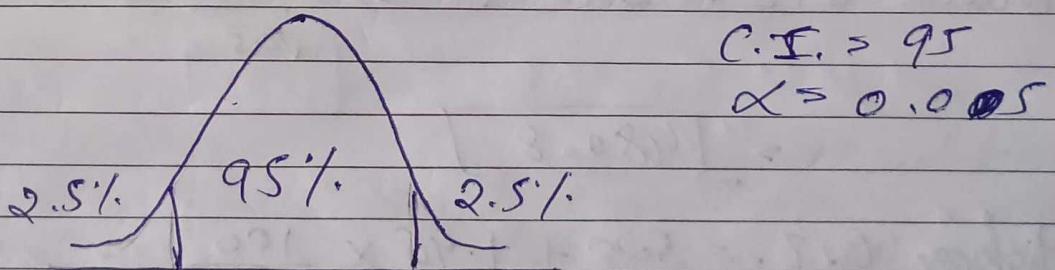


Confidence interval = 95%.

$$\text{Significance Value} = \alpha = 1 - \text{C.I.} \\ = 1 - 0.95 = 0.05$$

if $P \leq \text{Significance Value}$
we reject the null hypothesis
else:
we fail to reject the null hypothesis.

\Rightarrow Confidence Interval & Margin of Error



We construct a confidence to help estimate.
What is the actual value of unknown population mean.

Point estimate \pm margin of error

$$\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

- 1) In the verbal Mean of cat exam, the S.D is known to be 100, A sample of 25 test taken has a mean of 520. Construct a 95% of CI about the mean.

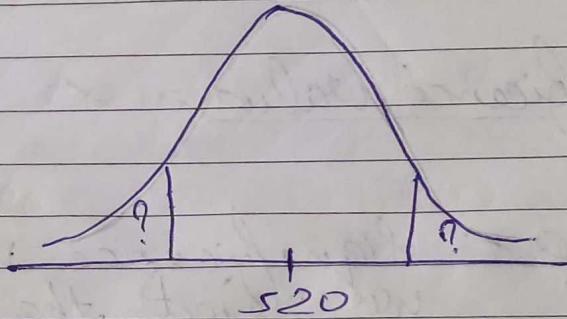
$$\sigma = 100$$

$$n = 25$$

$$\bar{x} = 520$$

$$CI = 95\%$$

$$\alpha = 0.05$$



$$\bar{x} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$Z_{0.025} = -1.96$$

$$\text{Lower C.I.} = 520 - (1.96) \times \frac{100}{\sqrt{25}}$$

$$= \boxed{480.8}$$

$$\text{Higher C.I.} = 520 + 1.96 \times \frac{100}{\sqrt{25}}$$

$$= \boxed{559.2}$$

I am 95% confident that the mean of scores lies b/w 480.8 and 559.2

⇒ Hypothesis testing & Statistical Analysis

- ① Z-test ↗ avg value
- ② t-test ↗
- ③ Chi square → categorical
- ④ ANOVA → Variance

① Z-test

- ① The avg height of all residents in a city is 168 cm with a $\sigma = 3.9$. A doctor believes the mean to be different. He measured the height of 36 individuals and found the avg height to be 169.5 cm.

① State null & alternate hypothesis.

② At a 95% confidence level, is there enough evidence to reject the null hypothesis.

Ans.

$$\mu = 168 \text{ cm}$$

$$\sigma = 3.9$$

$$n = 36$$

$$\bar{x} = 169.5 \text{ cm}$$

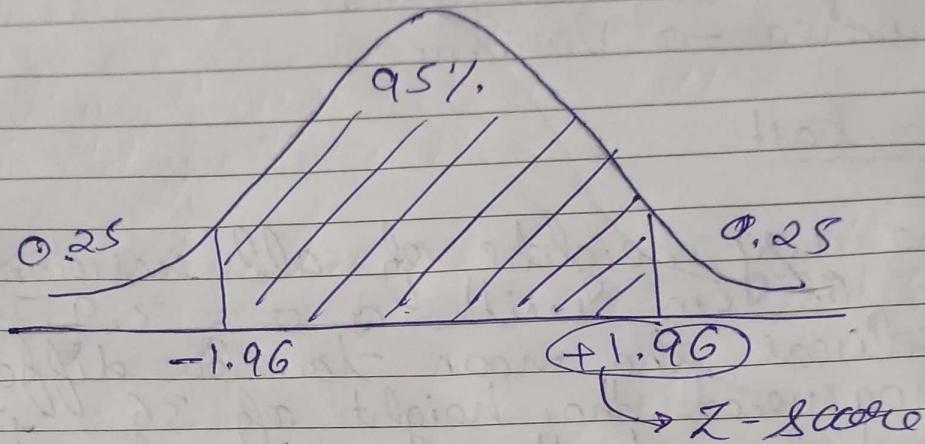
Whenever population sd is given then we should definitely use Z test.

- a) Null Hypothesis $H_0: \mu = 168 \text{ cm}$
b) Alternative Hypothesis $H_p: \mu \neq 168 \text{ cm}$.

$$CI = 0.95$$

$$\alpha = 0.05$$

Decision Boundary



Statistical Analysis

$$Z_{\text{test}} = \frac{\bar{x} - \mu}{\sqrt{\sigma^2 / n}} \rightarrow \text{Standard error}$$

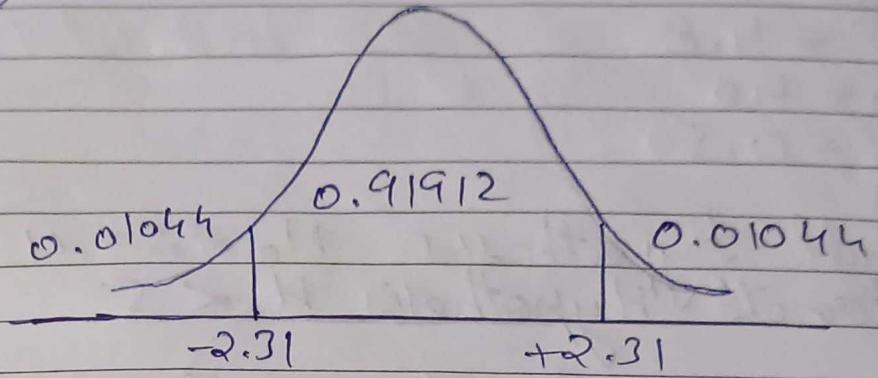
$$\Rightarrow \frac{169.5 - 168}{3.9 / \sqrt{36}} = 0.31$$

Conclusion \Rightarrow If Z-test value is less than -1.96 or greater than 1.96 , we reject the ~~null~~ null hypothesis.

$2.31 > 1.96 \Rightarrow$ Reject the null hypothesis.

The doctor is absolutely right.

⇒ P-value



$$P\text{ value} = 0.1044 + 0.1044 = 0.02088$$

if P value < significance value

$0.02088 < 0.05 \Rightarrow$ Reject the null hypothesis.

- ② A factory manufactures bulbs with an avg warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will manufacture in less than 5 years. He tests a sample of 40 bulbs and find the avg time to be 4.8 years

- State null & alternate hypothesis.
- At a 2% significance level, is there

Page
Date

enough evidence to support the idea that the warranty should be extended?

Ans.

$$\mu = 5$$

$$\bar{x} = 4.8$$

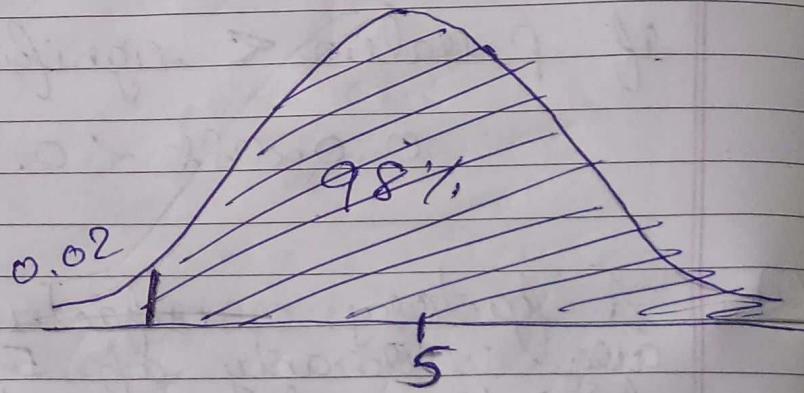
$$n = 40$$

$$\sigma = 0.50$$

1) Null Hypothesis $H_0: \mu = 5$

2) Alternate Hypothesis $H_1: \mu < 5$

3) Decision Boundary $\approx C.I. = 0.98$
 $\alpha = 1 - 0.98 = 0.02$

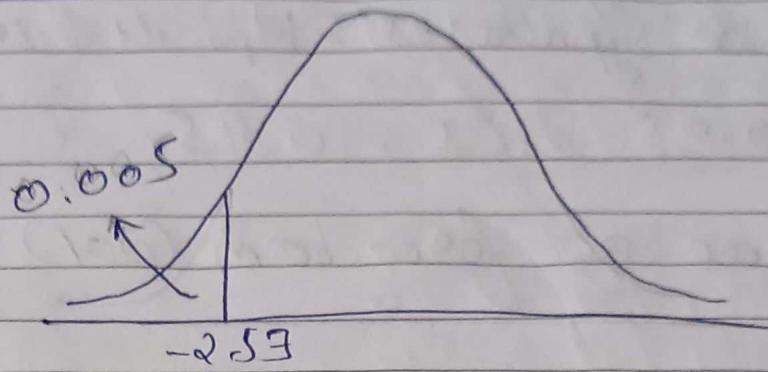


a) Z-test

$$= \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{4.8 - 5}{0.50 / \sqrt{40}} = -2.53$$

Z test $< -2.05 \Rightarrow$ Reject the null hypothesis.

5) P-value:



$p \text{ value} < \text{significance value}$
We reject the H_0 .

⇒ T-test

① In the population, the avg IQ is 100. A team of researchers want to test a new medication to see if it has either a +ve or -ve effect on intelligence or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence? C.I. = 95%.

Ans. If population s.d. is not given, use t-test

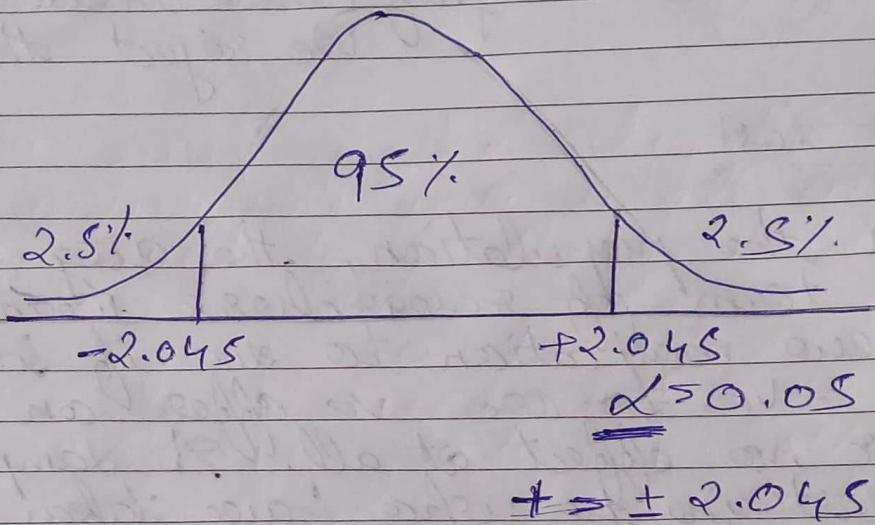
$$\mu = 100, n = 30, \bar{x} = 140, S = 20, C.I. = 0.95 \\ t = 0.05$$

① Null Hypothesis: $H_0: \mu = 100$
 Alternate Hypothesis: $H_1: \mu \neq 100$ (2 tailed)

② $\alpha = 0.05$ CI = 0.95

③ Degree of freedom (dof) = $n - 1 = 30 - 1 = 29$

④ Decision Rule



Conclusion: The t -test is less than -2.045 and greater than 2.045 .
 Reject the null hypothesis.

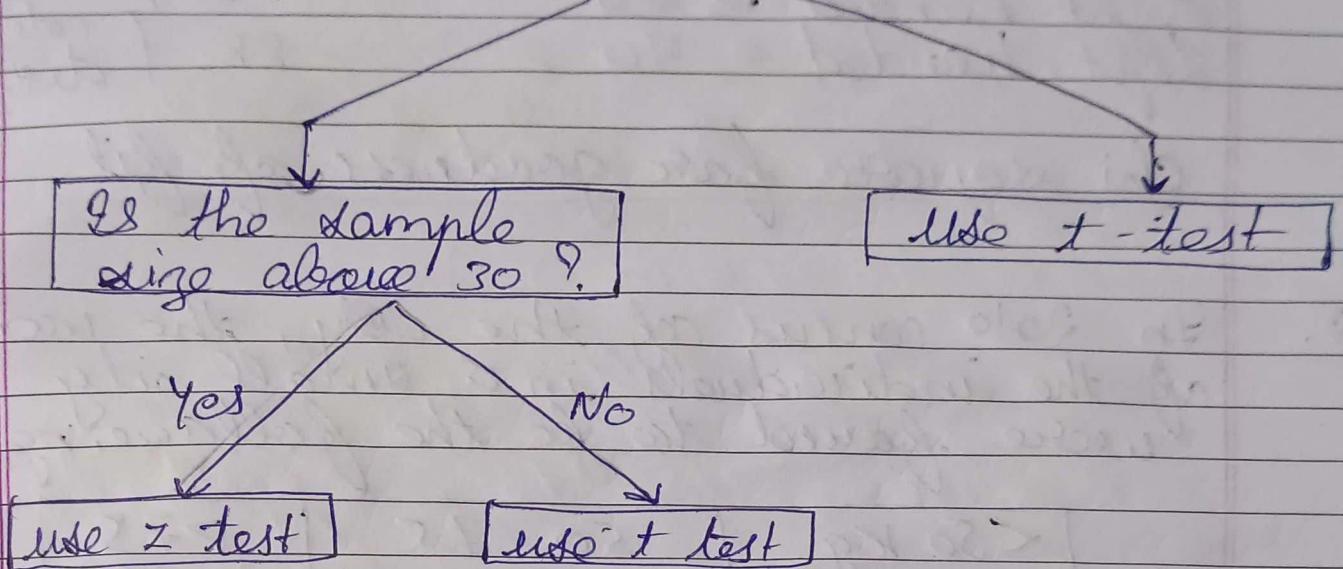
⑤ t -test Statistics

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \geq \frac{140 - 100}{20/\sqrt{30}} \geq 10.96$$

$10.96 > 2.0452$. So, we reject the null hypothesis.

⇒ When to use T-tests & z-tests

Do you know the population std σ ?



⇒ Chi-Squared Test

- ① The chi squared test for goodness of fit claim about population proportion (categorical test).
- ② It is a non parametric test that is performed on categorical data [nominal, ordinal].

① In a student class of 100 students, 30 are right handed. Does this class fit the theory? 12% of people are right handed.

<u>observed</u>		$\frac{6}{12}$	theory
Right handed	30		
Left handed	70	12	categorical

chi square for goodness of fit

Q. In 2010 census of the city, the weight of the individuals in a small city were found to be the following.

$\leq 50 \text{ kg}$	$50 - 75$	> 75
20%	30%	50%

In 2010, weight of $n = 500$ individuals were sampled, below are the results.

≤ 50	$50 - 75$	> 75
140	160	200

using $\alpha = 0.05$, would you conclude the population difference of weights has changed in the last 10 years.

Ans. Expected value

$\leq 50 \text{ kg}$	$50 - 75$	> 75
100	46150	250
100		

$H_0 \rightarrow$ The data meets the expectation.

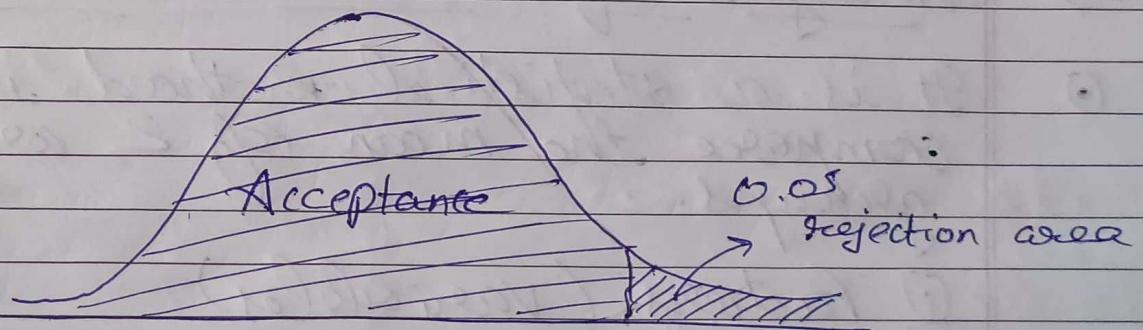
$H_1 \rightarrow$ The data does not meet the expectation.

$$\alpha = 0.05 \quad CI = 95\%$$

degrees of freedom

$$df = k - 1 = 3 - 1 = 2$$

Decision Boundary \rightarrow Chi-squared test



distribution is always right skewed

critical value = 5.991

If χ^2 is > 5.991 , we reject the H_0 .

Calculate Chi Square Test statistics

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

$$= \frac{(140 - 100)^2}{100} + \frac{(160 - 150)^2}{150} + \frac{(200 - 250)^2}{250}$$

$$= \frac{1600}{100} + \frac{100}{150} + \frac{2500}{250}$$

$$= 16 + 0.66 + 10$$

$$\therefore \chi^2 = 26.66$$

$26.66 > 5.991$, Reject H_0 .

⇒ Analysis of Variance (ANOVA)

- It is a statistical method used to compare the mean of 2 or more groups.

- 1 Factors (variables)
- 2 Variables

Eg Medicine → factors

levels → 5 mg 10 mg 20 mg [damage]

Mode of payment — factors

Cash phone Imps NEFT [cheques]

→ Types of ANOVA

① One Way ANOVA: One factor with at least 2 levels, then levels are independent.

eg: Doctor want to test a new medication to decrease headache. They split the participants into 3 condition want to degrees [10, 20, 30].

→ Doctors ask the patients to rate the headache (1-10).

Medication → Factor

<u>10 mg</u>	<u>20 mg</u>	<u>30 mg</u> → levels
5	7	2
9	8	7
—	—	—
—	—	—

②

Repeated Measures Anova

One factor with at least 2 levels, levels are dependent.

Running \rightarrow Factor

Day 1	Day 2	Day 3
8 km	5 km	6 km

③

Factorial Anova

Two or more factors (each of which with at least 2 levels), levels can be either independent or dependent.

Running \rightarrow Factor

	Day 1	Day 2	Day 3	
Gender	8	5	6	dependent
Male	7	4	3	
Female	6	5	4	
	3	2	1	

\Rightarrow

Hypothesis Testing in Anova

\rightarrow

Null Hypothesis: $H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_K$

→ Alternate Hypothesis: H_1 : At least one of the mean is not equal.

→ P-test Statistics

$P = \frac{\text{Variation B/w samples}}{\text{Variation within samples}}$

① Doctors want to test a new medication which reduces headache. They split the participants into 3 condition [15, 30, 45]. Later on the doctor ask the patient to rate the headache between [1 - 10]. Are there any differences between the 3 condition using $\alpha = 0.05$?

<u>15 mg</u>	<u>30 mg</u>	<u>45 mg</u>
9	7	4
8	6	3
7	6	2
8	7	3
8	8	4
9	7	3
8	6	2

① Define null hypothesis

$$H_0: \mu_{15} = \mu_{30} = \mu_{45}$$

② Alternate Hypothesis:

At least one mean is not equal.

③ State significance value

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

④ Calculate degree of freedom

$$N = 21 \quad a = 3 \quad n = 7$$

$$df \text{ between} = a - 1 = 2$$

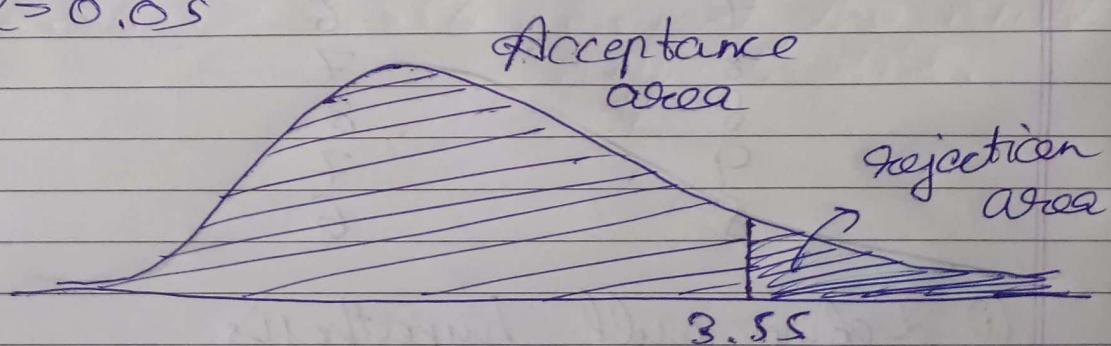
$$df \text{ within} = N - a = 18$$

$(2, 18)$ = F test table

$$df \text{ total} = N - 1 = 21 - 1 = 20$$

Decision boundary

$$\alpha = 0.05$$



if f test is > 3.55 , we reject the H_0 .

5. Calculate F test Statistics

	SS	df	MS	F
between	98.67	2	49.34	<u>86.56</u>
within	10.29	18	0.54	
total	108.95	20		

$$\textcircled{1} \quad \text{SS}_{b/w} = \frac{\sum (Ea_i)^2}{n} - \frac{T^2}{N}$$

$$15 \text{ mg} = 9+8+7+8+8+9+8 = 52 \quad \left. \quad \right\} T^2$$

$$30 \text{ mg} = 7+6+6+7+8+7+6 = 47 \quad \left. \quad \right\} T^2$$

$$45 \text{ mg} = 4+3+2+3+4+3+2 = 21 \quad \left. \quad \right\} T^2$$

$$= \frac{57^2 + 47^2 + 21^2}{7} - \frac{57 + 47 + 21}{21}$$

$$= \boxed{98.67}$$

$$\textcircled{2} \quad \text{SS}_{\text{within}} = \sum y^2 - \frac{\sum (Ea_i)^2}{n}$$

$$= \sum y^2 - \left[\frac{57^2 + 47^2 + 21^2}{7} \right]$$

$\sum y^2$ = sum of squares of each example

$$= 853 - \left[\frac{57^2 + 47^2 + 21^2}{7} \right]$$

$$= 10.29$$

$$\textcircled{3} \quad SS_{\text{total}} = \sum y^2 - \frac{T^2}{N}$$

$$853 - \frac{125^2}{21} = 108.95$$

\textcircled{4} MS

Mean Sq Between
Mean Sq Within

$$F = \frac{49.34}{0.54} = 86.56$$

$86.56 > 3.556$, Reject the null hypothesis.