

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.

↓ Defn

It consists of Organizing, Summarizing, and visualizing data.

- ① Measure of Central tendency.
- ② Measure of dispersion (V, SD, Z score)

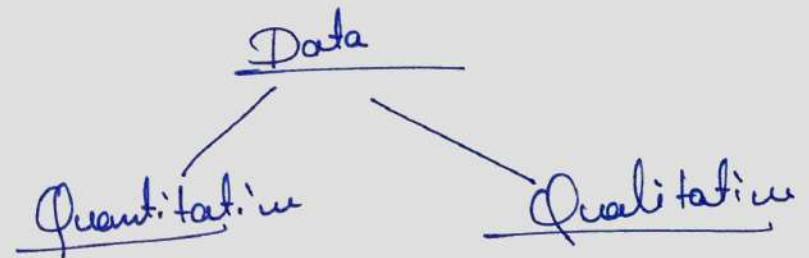
② Diff types of distribution

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

Inferential Stats

Defn: It consists of using data you have measured to form conclusions.

- 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.
- w.r.t age, gender, income etc.

③ Convenience Sampling.

- Participants are selected based on availability and willingness to take part.

④ Cluster Sampling.

- Divide a population into clusters
- eg: districts, schools and randomly selecting the clusters.

⑤ Snowball Sampling.

- Existing study subjects are used to recruit more subjects into the sample.

⑥ Purposive Sampling.

- Selecting samples based on the judgement of the survey taker or researcher.

⑦ Systematic Sampling

- It is a probability sampling method where researchers select members at a regular interval.

Scale of Measurement

- 1) Nominal scale data → Qualitative / categorical data.
Ex: Gender, Color, Labels.
Order or Rank doesn't matter.
- 2) Ordinal scale data → Rank in imp. Order Matters.
Difference cannot be measured unless we take more information.
- 3) Interval scale data → The Rank / order matters.
Difference can be measured.
The Ratio cannot be measured.
No 0 starting point.
Eg: Temperature Variable.
Intelligence Quotient.
- 4) Ratio scale data → The Order will matter.
Difference is measurable (Ratio).
Contains a '0' starting point.
Eg: Student marks.

Ratio Scale data

The Order will matter.
Difference is measurable (Ratio).
Contains a '0' starting point.
Eg: Student marks.

The Rank / order matters.
Difference can be measured.
The Ratio cannot be measured.
No 0 starting point.
Eg: Temperature Variable.
Intelligence Quotient.

Random Variables

Variable (x, y)

$$x + y = 7$$

$$x = 2$$

$$x = 2$$

$$y = 6$$

$$y = 7 - x$$

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

ex: Toss a coin.

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

Rolling a dice.

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

defn: 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.

Covariance & Correlation

It is the measure of the relationship b/w two random variables. It measures how much the variables change together, or the variance b/w them.

→ It measures the direction of a relationship b/w two variables.

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

$\bar{y}, \bar{x} = \text{mean}$

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

Advantages

→ +ve or -ve value

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

Disadvantages

Doesn't have a specific limit value.

② Pearson Correlation Coefficient

$$r_{x,y} = \frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

PCC doesn't work well with non-linear data. So we will use Spearman Correlation.

③ Spearman's rank Correlation Coefficient

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

$R = \text{Rank}$. based on frequency / order / value

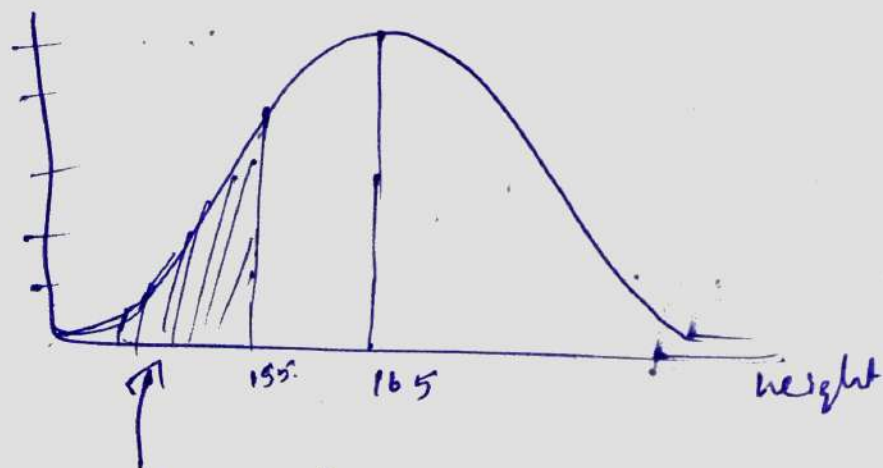
→ If 2 features are highly correlated, it's ok to drop one of the features.

Probability Distribution Function.

Probability density function

Continuous value.

ex: Age, Height.
float



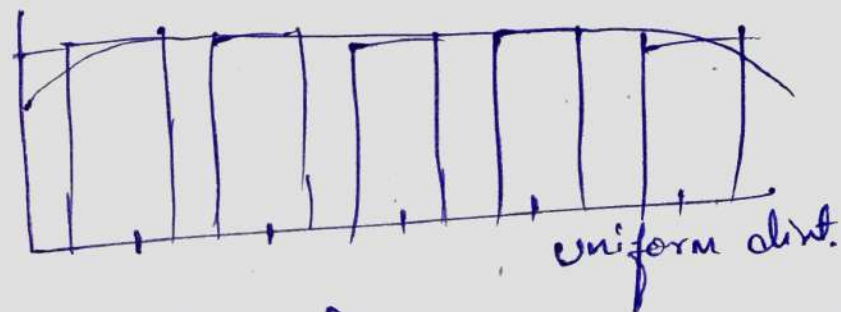
$Pr(H \leq 155)$

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

Probability mass function.

Discrete values.

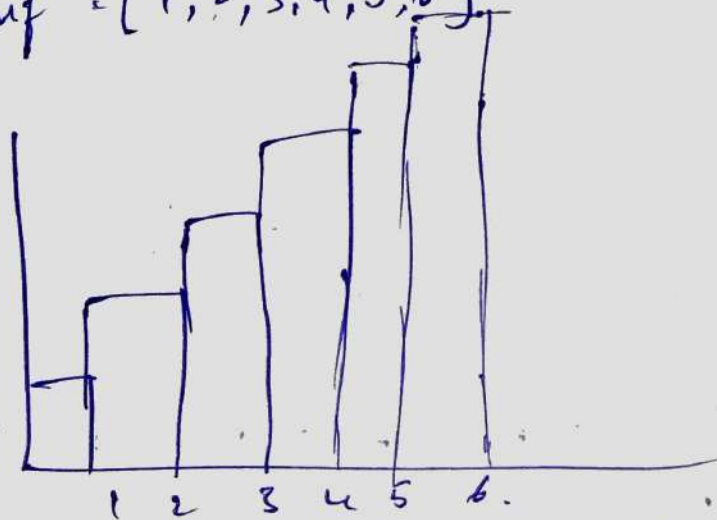
ex: no of bank acc.
int
Rolling a dice.



$$\begin{aligned} Pr(X \leq 4) &= Pr(X=1) + Pr(X=2) + \\ &\quad Pr(X=3) + Pr(X=4) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{4}{6} = \underline{\underline{\frac{2}{3}}} \end{aligned}$$

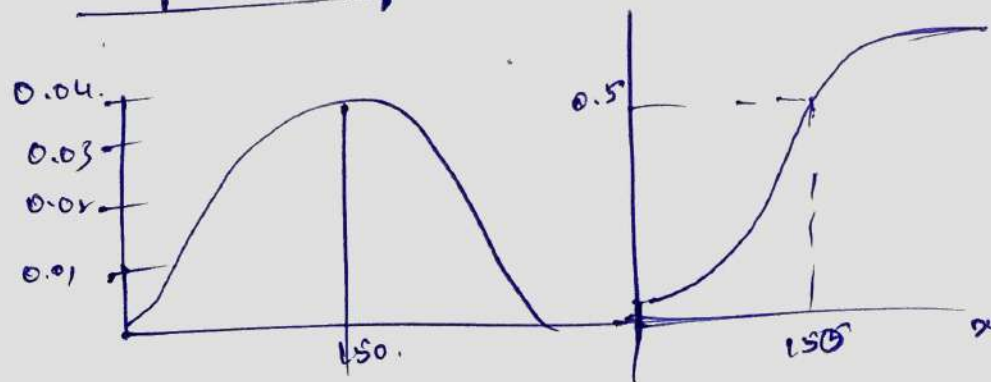
3) Cumulative Distribution Function

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



→ Adding the previous value to the next value

→ pdf & cdf



Probability Density of pdf → Gradient or derivative of cdf.

pdf is the derivative of cdf.

cdf is the integration of pdf.

Types of Probability function

- 1) Normal / Gaussian 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)

① Bernoulli Distribution

→ discrete Random Variable (pmf)

→ Outcomes are Binary

ex:- Tossing a coin {H, T}.

$$Pr(H) = 0.5 = p$$

$$Pr(T) = 1 - 0.5 = q$$

whether a person will pass/fail.

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

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

② Binomial Distribution

→ Discrete Random variable

→ Every Experiment outcome is Binary.

→ This experiment is performed for n trials

Notes:- Refer wikipedia.

eg:- Tossing a coin 10 times.

notation $\rightarrow B(n, p)$

$$Pr(H) = 0.5 = p$$

$$Pr(T) = 0.5 = q$$

parameter = $n \in \{0, 1, 2, 3, \dots\} \Rightarrow$ 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:

$$Pr(k, n, p) = {}^nC_k p^k (1-p)^{n-k}$$

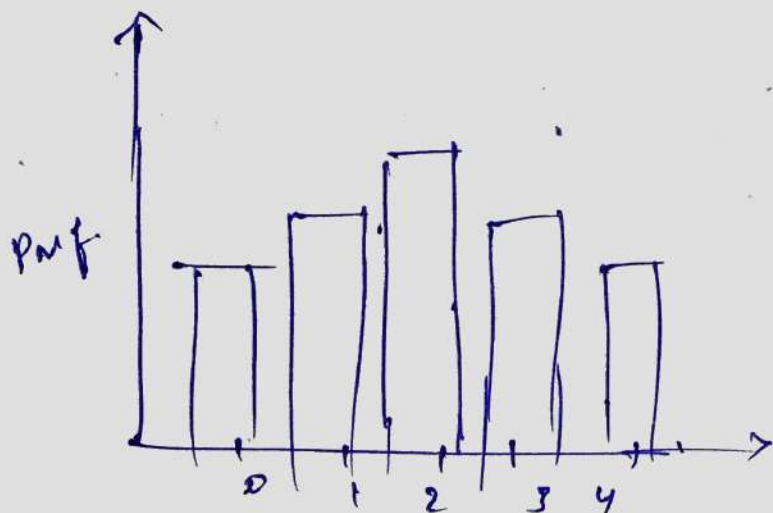
$${}^nC_k = \frac{n!}{k!(n-k)!}$$

Poisson Distribution

- Discrete Random Variable (pmf)
- Describes the no of events occurring in a fixed time interval.

eg: no of people visiting hospital every hour.

no of people visiting banks at 11 am.



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

what is the probability of no of people visiting at 3 pm. $\Rightarrow Pr(X=3)$

$$= \frac{e^{-3} 3^3}{5!} = 0.001 < 10\%$$

10% of the people visit the bank at 3 pm

Empirical Rule of Normal dist

68
1st sd

95
2nd sd

99.7 %
3rd sd

UNIFORM DISTRIBUTION

Continuous Uniform Distribution.

It is a continuous uniform dist in a probability dist that takes values within a specified range. It is defined by two parameters, a, b where a is the lower limit and b is the upper limit.

Notation = $U(a, b)$

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

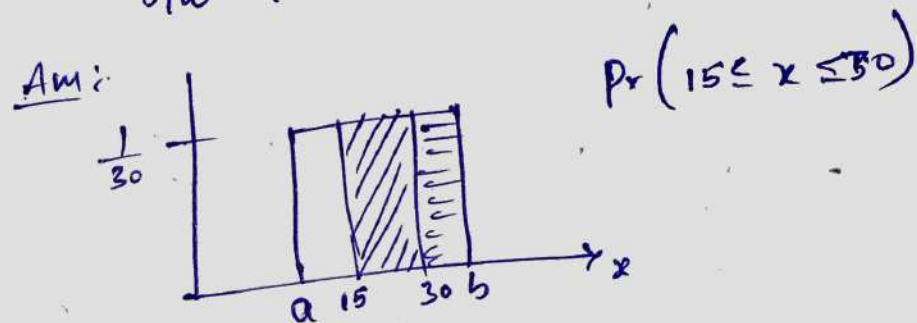
$$pdf = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

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

$$\text{mean} = \frac{a+b}{2} \quad \text{median} = \frac{a+b}{2}$$

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

\rightarrow probability of daily sales to fall b/w 15 and 30. $a=10$
 $b=40$



$$\rightarrow Pr(x > 30)$$

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

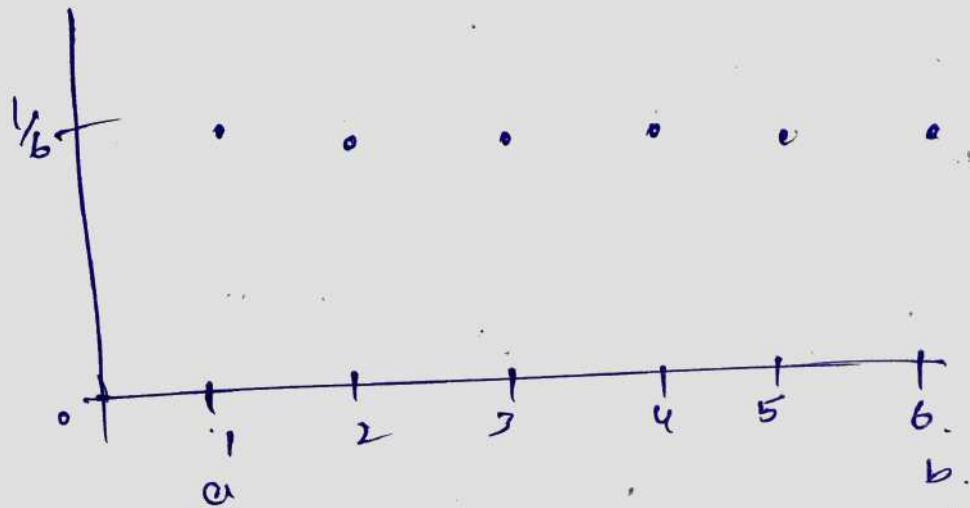
$$40 - 30 \times \frac{1}{30}$$

$$\frac{10}{30} = \boxed{0.33\%}$$

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

② Discrete Uniform Distribution (pmf)

Ex → Rolling a Dice



Notation: $U(a, b)$

$$\begin{array}{l} n = 6 \quad \text{no of outcomes} \\ \boxed{n = b - a + 1} \end{array}$$

$$Pr(i) = \frac{1}{n} = \frac{1}{6}$$

parameter a, b $\boxed{b > a}$

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

$$\text{mean} = \frac{a+b}{2}$$

$$\text{median} = \frac{a+b}{2}$$

Definition: A discrete uniform distribution is a statistical distribution when the probability of outcomes is equally likely and with finite values. In a discrete uniform dist, Every one of a n values has equal probability $\frac{1}{n}$.

Ex:- Rolling a die
Selecting a card from deck of cards
Flipping a fair coin.

① STANDARD NORMAL DISTRIBUTION

Z-SCORE

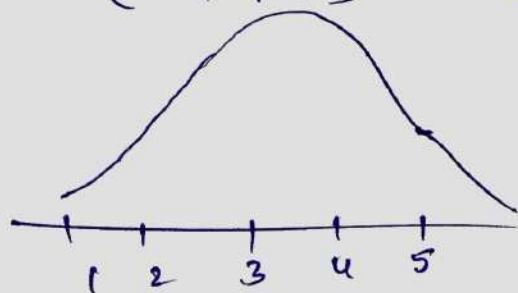
Z-Score

⇒ Normally Distributed.

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

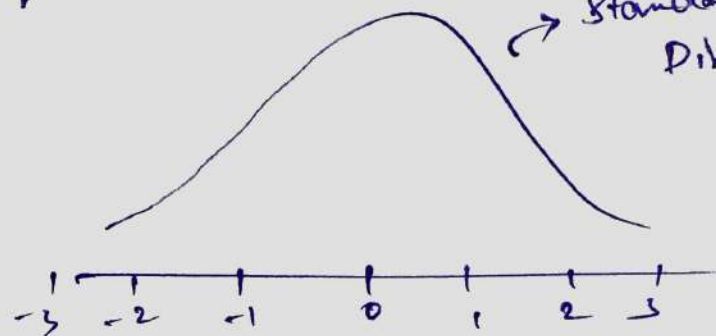
$$\mu = 3$$

$$\sigma = 1.414$$



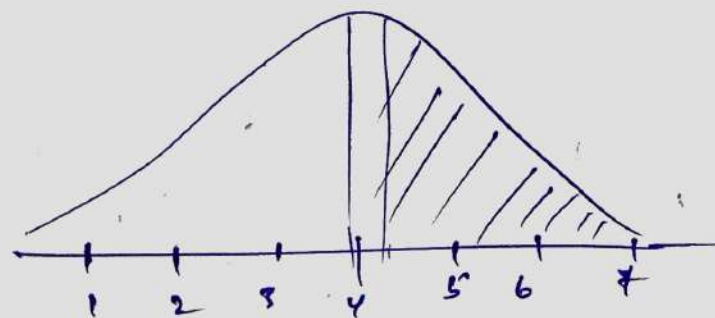
if $\mu = 0$ & $\sigma = 1$

⇒ Standard Normal Distribution.



$$\Rightarrow Z\text{-Score} = \frac{X_i - \mu}{\sigma}$$

Z-table - we use to find out the area under the curve



what % of the data is falling below 4.5

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

$$\Rightarrow 0.6915$$

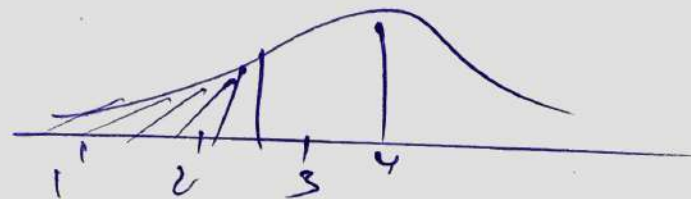
$$\text{Area under the curve} = 1 - 0.6915 = 0.3085 = 30.85\%$$

② percent of data falling below 2.5 ?

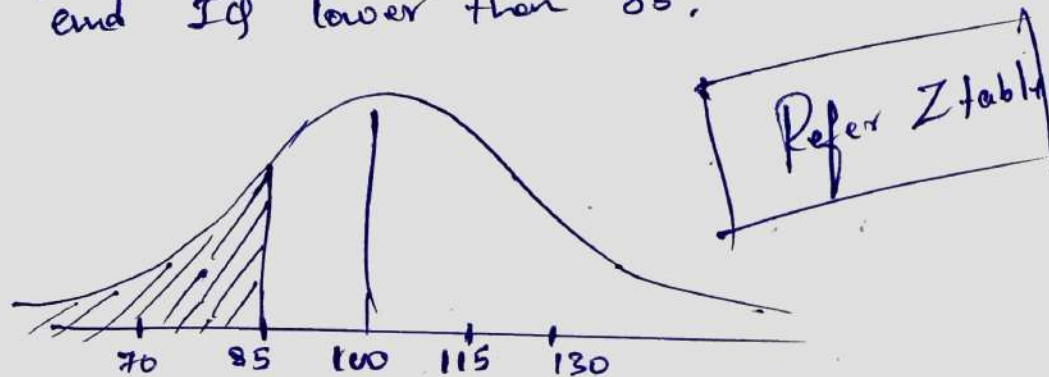
$$Z\text{-score} = \frac{2.5 - 4}{1} = -1.5$$

$$= 0.0668$$

$$= \boxed{6.68\%}$$



In india. the avg IQ is 100, with σ of 15. what is the percentage of the population would you expect to have and IQ lower than 85.



$$Z_{\text{score}} = - \frac{85 - 100}{15} = \frac{-15}{15} = -1 = 0.1587$$

$\boxed{15.87\%}$

$$= \cancel{0.4602}$$

$\boxed{46.02\%}$

Area under the curve > 85 .

$$1 - 0.1587$$

$$= 0.8413 = \boxed{84.13\%}$$

Area between 85 & 100

$$0.5 - 0.1587$$

$$= 0.3413 = \boxed{34.13\%}$$

z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
-3.4	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0003	.0002
-3.3	.0005	.0005	.0005	.0004	.0004	.0004	.0004	.0004	.0004	.0003
-3.2	.0007	.0007	.0006	.0006	.0006	.0006	.0006	.0005	.0005	.0005
-3.1	.0010	.0009	.0009	.0009	.0008	.0008	.0008	.0008	.0007	.0007
-3.0	.0013	.0013	.0013	.0012	.0012	.0011	.0011	.0011	.0010	.0010
-2.9	.0019	.0018	.0018	.0017	.0016	.0016	.0015	.0015	.0014	.0014
-2.8	.0026	.0025	.0024	.0023	.0023	.0022	.0021	.0021	.0020	.0019
-2.7	.0035	.0034	.0033	.0032	.0031	.0030	.0029	.0028	.0027	.0026
-2.6	.0047	.0045	.0044	.0043	.0041	.0040	.0039	.0038	.0037	.0036
-2.5	.0062	.0060	.0059	.0057	.0055	.0054	.0052	.0051	.0049	.0048
-2.4	.0082	.0080	.0078	.0075	.0073	.0071	.0069	.0068	.0066	.0064
-2.3	.0107	.0104	.0102	.0099	.0096	.0094	.0091	.0089	.0087	.0084
-2.2	.0139	.0136	.0132	.0129	.0125	.0122	.0119	.0116	.0113	.0110
-2.1	.0179	.0174	.0170	.0166	.0162	.0158	.0154	.0150	.0146	.0143
-2.0	.0228	.0222	.0217	.0212	.0207	.0202	.0197	.0192	.0188	.0183
-1.9	.0287	.0281	.0274	.0268	.0262	.0256	.0250	.0244	.0239	.0233
-1.8	.0359	.0351	.0344	.0336	.0329	.0322	.0314	.0307	.0301	.0294
-1.7	.0446	.0436	.0427	.0418	.0409	.0401	.0392	.0384	.0375	.0367
-1.6	.0548	.0537	.0526	.0516	.0505	.0495	.0485	.0475	.0465	.0455
-1.5	.0668	.0655	.0643	.0630	.0618	.0606	.0594	.0582	.0571	.0559
-1.4	.0808	.0793	.0778	.0764	.0749	.0735	.0721	.0708	.0694	.0681
-1.3	.0968	.0951	.0934	.0918	.0901	.0885	.0869	.0853	.0838	.0823
-1.2	.1151	.1131	.1112	.1093	.1075	.1056	.1038	.1020	.1003	.0985
-1.1	.1357	.1335	.1314	.1292	.1271	.1251	.1230	.1210	.1190	.1170
-1.0	.1587	.1562	.1539	.1515	.1492	.1469	.1446	.1423	.1401	.1379
-0.9	.1841	.1814	.1788	.1762	.1736	.1711	.1685	.1660	.1635	.1611
-0.8	.2119	.2090	.2061	.2033	.2005	.1977	.1949	.1922	.1894	.1867
-0.7	.2420	.2389	.2358	.2327	.2296	.2266	.2236	.2206	.2177	.2148
-0.6	.2743	.2709	.2676	.2643	.2611	.2578	.2546	.2514	.2483	.2451
-0.5	.3085	.3050	.3015	.2981	.2946	.2912	.2877	.2843	.2810	.2776
-0.4	.3446	.3409	.3372	.3336	.3300	.3264	.3228	.3192	.3156	.3121
-0.3	.3821	.3783	.3745	.3707	.3669	.3632	.3594	.3557	.3520	.3483
-0.2	.4207	.4168	.4129	.4090	.4052	.4013	.3974	.3936	.3897	.3859
-0.1	.4602	.4562	.4522	.4483	.4443	.4404	.4364	.4325	.4286	.4247
-0.0	.5000	.4960	.4920	.4880	.4840	.4801	.4761	.4721	.4681	.4641

⑬ CENTRAL LIMIT THEOREM.

After taking multiple samples from the population

Sample dist of mean $\bar{X} \approx N(\mu, \sigma/\sqrt{n})$

n can be any value.

Mean will be approx to μ

Standard dev will be $\boxed{\sigma/\sqrt{n}}$ Standard error

* QQ plots \rightarrow to check whether the dist is normal or not

Defn:- CLT is a Statistical theory that states that when a large sample size has a finite variance, the samples will be normally distributed.

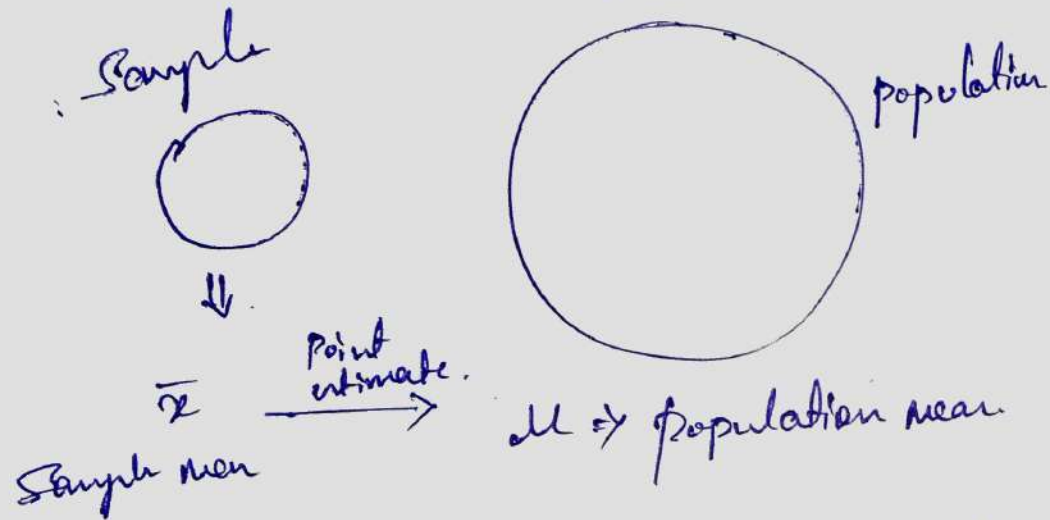
\rightarrow The theorem also states that the mean of the samples will be approximately equal to the mean of the whole population.

\rightarrow CLT assumes that all samples are identical in size, and regardless of the population's actual distribution

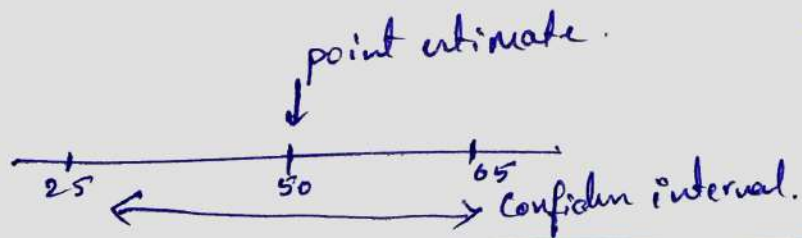
\rightarrow The theorem holds true regardless of whether the source population is normal or skewed, provided the sample size is sufficiently larger than $n=30$

INFERENCE STATISTICS

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

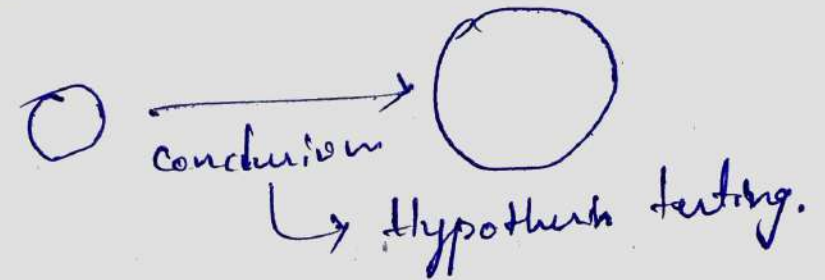


2) Interval Estimate → Range of values used to estimate the unknown population parameter.



Hypothesis & Hypothesis Testing Mechanism

Inferential Stats → Conclusion



Hypothesis testing mechanism

Person Crime →
① Null hypothesis (H_0)
⇒ The person is not guilty.
→ Assumptions you are beginning with.

② Alternate hypothesis (H_1) The person is guilty.
→ Opp of H_0 .

→ Evidence . Experimental Statistical Analysis . 1)

Z-test, t-test, Anova test etc

(14)

P value

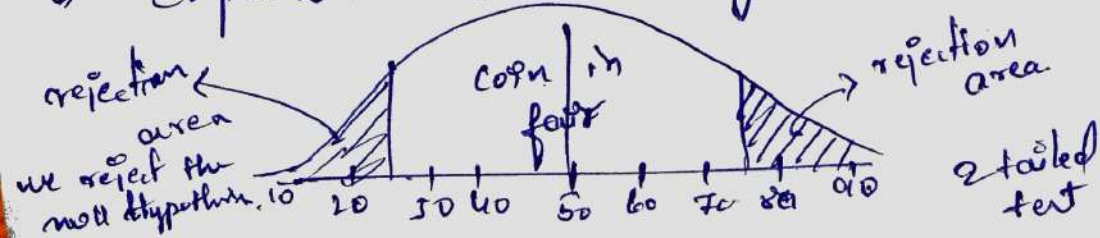
P value is a number calculated from a statistical test, that determine how likely you are to have found a practical set of observations if the null hypothesis were true. P values are used in hypothesis testing to decide whether to reject the null hypothesis.

Example: Coin is fair or not
(H, T).

Hypothesis Testing

Null hypothesis $H_0 \rightarrow$ Coin is fair
Alternate hypothesis $H_1 \rightarrow$ Coin is not fair.

Experiment \rightarrow Tossing



Confidence interval = 95%.

$$\text{Significance value} = \alpha = 1 - CI \\ = 1 - 0.95 = 0.05$$

if $p < \text{Significance value}$.

we reject the null hypothesis

else:

we fail to reject the null hypothesis

Confidence Interval & Margin of Error



$$CI = 95 \\ \alpha = 0.05$$

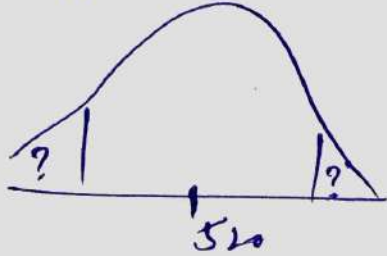
We construct a Confidence to help estimate what is the actual value of unknown population mean

Point estimate \pm margin of error.

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

→ In the verbal exam of cat exam, the s.d is known to be 100, A sample of 25 test takers has a mean of 520. Construct a 95% of CI. about the mean.

$$\begin{aligned}\sigma &= 100 \\ n &= 25 \\ \bar{X} &= 520 \\ CI &= 95\% \\ \alpha &= 0.05.\end{aligned}$$



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

$$\begin{aligned}\text{Lower C.I.} &= 520 - (1.96) \times \frac{100}{\sqrt{25}} \\ &= \underline{\underline{480.8}}\end{aligned}$$

$$Z_{0.025} = -1.96$$

$$\text{Higher CI} = 520 + 1.96 \times \frac{100}{\sqrt{25}} = \underline{\underline{559.2}}$$

I am 95% Confident that the mean cat score 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

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 null hypothesis

Ans) \Rightarrow

$$\begin{aligned}\mu &= 168 \text{ cm.} \\ \sigma &= 3.9 \\ n &= 36 \\ \bar{x} &= 169.5 \text{ cm.}\end{aligned}$$

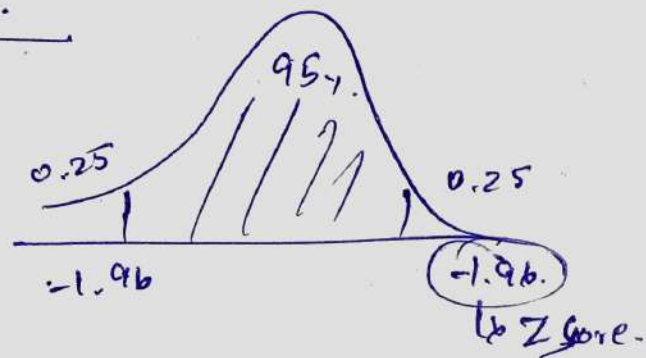
whenever population sd is given, then we should definitely use z-test.

- 15) Null hypothesis - $H_0: \mu = 168 \text{ cm}$
 Alternate hypothesis - $H_1: \mu \neq 168 \text{ cm}$

$$CI = 0.95$$

$$\alpha = 0.05$$

Decision boundary.



Statistical Analysis

$$Z_{\text{test}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \Rightarrow \text{Standard error}$$

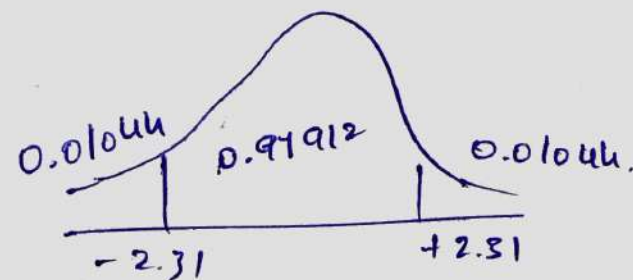
$$= \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = \boxed{2.31}$$

Conclusion \Rightarrow If Z_{test} value is less than -1.96 or greater than 1.96 , we reject the

null hypothesis.

$2.31 > 1.96 \Rightarrow$ Reject the null hypothesis
 The doctor is absolutely right.

P-value



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

if $p\text{-value} < \text{Significance value}$

$0.02088 < 0.05 \Rightarrow$ Reject the null hypothesis

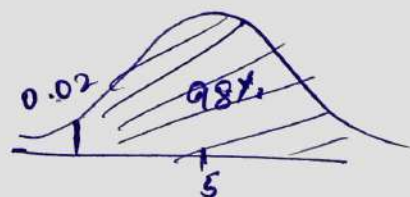
$$Z_{\text{score}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

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.
 a) State null & Alternate hypothesis.
 b) At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised.

$$\begin{aligned} \mu &= 5 \\ \bar{x} &= 4.8 \\ n &= 40 \\ \sigma &= 0.50 \end{aligned}$$

1) Null hypothesis $H_0: \mu = 5$
 2) Alternate hypothesis $H_1: \mu < 5$

3) Decision boundary $CI = 0.98$
 $\alpha = 1 - 0.98 = 0.02$

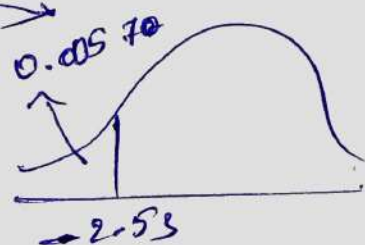


4) Z-test

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

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

5) P-value



P-value $<$ significant 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?
 $CI = 95\%$

\Rightarrow If population sd is not given, use t-test.

$$\mu = 100, n = 30, \bar{x} = 140, s = 20, CI = 0.95$$

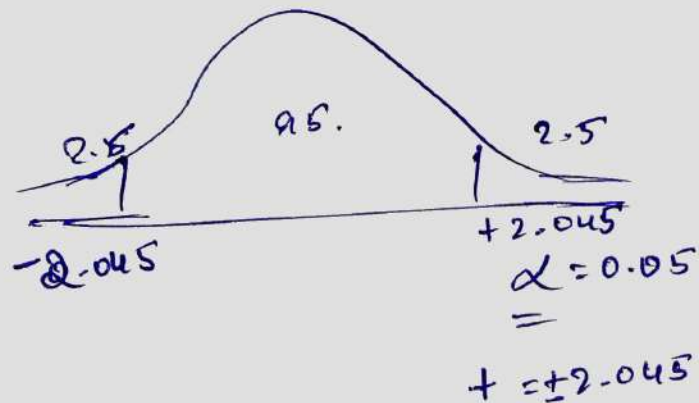
$$\alpha = 0.05$$

1) Null hypothesis $H_0: \mu = 100$
 $H_1: \mu \neq 100$ (2-tailed)

2) $\alpha = 0.05$ $CI = 0.95$

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

④ Decision Rule



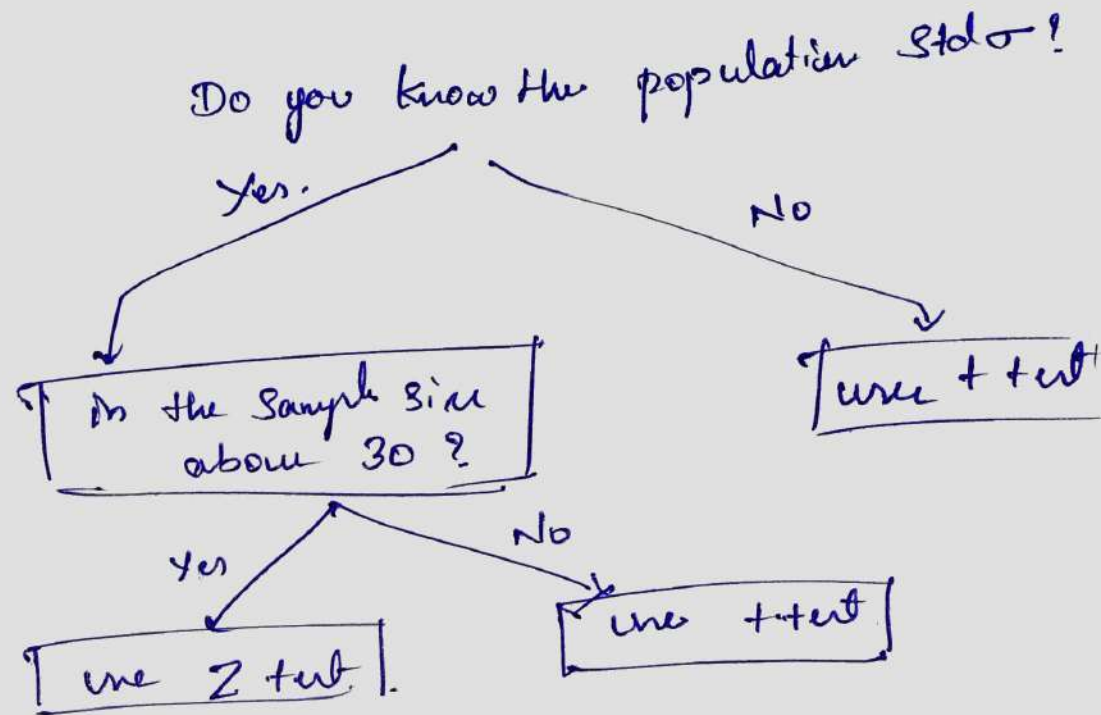
Concl: the t -test is less than -2.045 and greater than 2.045 . Reject the null hypothesis.

⑤ T-test Statistic.

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

$10.96 > 2.0452$ so we reject the null hypothesis.

When to use T-test & Z test



CHI SQUARED TEST

The chi squared test for Goodness of fit claim about population proportion (categorical test).

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.

Right handed = 30
Left handed = 70

Observed.	E
12	} theory categorical distribution
88	

Chi square for goodness of fit

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

<50 kg	50-75	>75
20%	30%	50%

In 2020, 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 values

<50 kg	50-75	>75
100	150	250

$H_0 \rightarrow$ The data meets the expectation.
 $H_1 \rightarrow$ The data does not meet the expectation

7

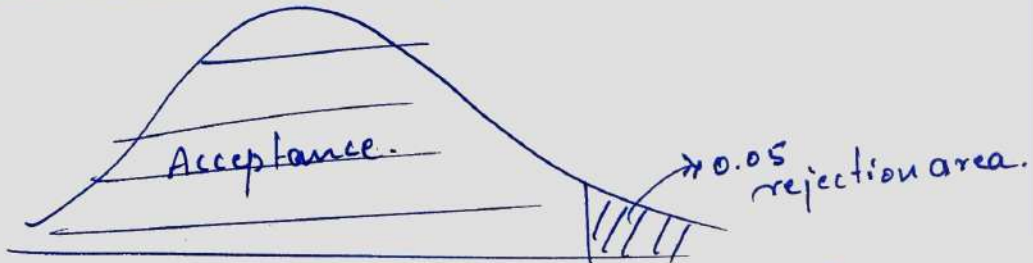
$$\alpha = 0.05$$

$$CI = 95\%$$

degree of freedom.

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

Decision Boundary \rightarrow chi squared test



distribution is always right skewed

$$\text{Critical value} = 5.991$$

If χ^2 is > 5.991 , we reject the H_0

Calculate chi square test statistic

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

Table: Chi-Square Probabilities

The areas given across the top are the areas to the right of the critical value. To look up an area on the left, subtract it from one, and then look it up (ie: 0.05 on the left is 0.95 on the right)

df	0.995	0.99	0.975	0.95	0.90	0.10	0.05	0.025	0.01	0.005
1	---	---	0.001	0.004	0.016	2.706	3.841	5.024	6.635	7.879
2	0.010	0.020	0.051	0.103	0.211	4.605	5.991	7.378	9.210	10.597
3	0.072	0.115	0.216	0.352	0.584	6.251	7.815	9.348	11.345	12.838
4	0.207	0.297	0.484	0.711	1.064	7.779	9.488	11.143	13.277	14.860
5	0.412	0.554	0.831	1.145	1.610	9.236	11.070	12.833	15.086	16.750
6	0.676	0.872	1.237	1.635	2.204	10.645	12.592	14.449	16.812	18.548
7	0.989	1.239	1.690	2.167	2.833	12.017	14.067	16.013	18.475	20.278
8	1.344	1.646	2.180	2.733	3.490	13.362	15.507	17.535	20.090	21.955
9	1.735	2.088	2.700	3.325	4.168	14.684	16.919	19.023	21.666	23.589
10	2.156	2.558	3.247	3.940	4.865	15.987	18.307	20.483	23.209	25.188
11	2.603	3.053	3.816	4.575	5.578	17.275	19.675	21.920	24.725	26.757
12	3.074	3.571	4.404	5.226	6.304	18.549	21.026	23.337	26.217	28.300
13	3.565	4.107	5.009	5.892	7.042	19.812	22.362	24.736	27.688	29.819
14	4.075	4.660	5.629	6.571	7.790	21.064	23.685	26.119	29.141	31.319
15	4.601	5.229	6.262	7.261	8.547	22.307	24.996	27.488	30.578	32.801
16	5.142	5.812	6.908	7.962	9.312	23.542	26.296	28.845	32.000	34.267
17	5.697	6.408	7.564	8.672	10.085	24.769	27.587	30.191	33.409	35.718
18	6.265	7.015	8.231	9.390	10.865	25.989	28.869	31.526	34.805	37.156
19	6.844	7.633	8.907	10.117	11.651	27.204	30.144	32.852	36.191	38.582
20	7.434	8.260	9.591	10.851	12.443	28.412	31.410	34.170	37.566	39.997

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

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

$$= 16 + 0.66 + 10$$

$$\chi^2 = 26.66$$

$26.66 > 5.991$, Reject H_0 .

① Analysis of Variance (ANOVA)

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

① Factors (variables)

② Levels

Eg Medicine → factor

levels → 5mg 10mg 20mg [dosage]

Mode of payment — factor

Upay Phone Impa NEFT [levels]

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 conditions w.r.t to dosages [10, 20, 30].

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

Medication		→ Factor
10 mg	20 mg	30 mg → levels
5	7	2
9	8	7
—	—	—
—	—	—

② Repeated Measures Anova

one factor with atleast 2 levels, levels are dependent.

Running → Factor.

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

③ Factorial Anova.

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

Runrig → Factor

	<u>Day 1</u>	<u>Day 2</u>	<u>Day 3</u>
<u>Gender</u>	8	5	6
<u>Male</u>	7	4	3
<u>Female</u>	6	5	4
	3	2	1

↓ dependent

Hypothesis testing in Anova.

Null hypothesis $H_0 : \mu_1 = \mu_2 = \mu_3 = \dots = \mu_k$

Alternate hypothesis H_1 : Atleast one of the mean is not equal.

F test Statistics

$$F = \frac{\text{Variation b/w samples}}{\text{Variation within samples.}}$$

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 conditions 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

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

④ Calculate degree of freedom

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

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

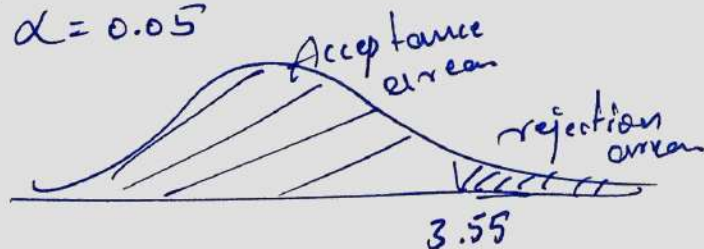
$$df \text{ within} = N - a = 18$$

$$(2, 18) = F_{\text{table}}$$

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

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		

$$① SS_{b/w} = \frac{\sum (\sum a_i)^2}{n} - \frac{T^2}{N}$$

$$15 \text{ mg} = 9 + 8 + 7 + 8 + 8 + 9 * 8 = 57$$

$$30 \text{ mg} = 7 + 6 + 6 + 7 + 8 + 7 + 6 = 47$$

$$45 \text{ mg} = 4 + 3 + 2 + 3 + 4 + 3 + 2 = 21$$

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

$$= \boxed{98.67}$$

② SS within $= \sum y^2 - \frac{\sum (\sum a_i)^2}{n}$

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

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

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

$$= \boxed{10.29}$$

③ SS total $= \sum y^2 - \frac{T^2}{N}$

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

④ MS.

$$\frac{\text{Mean Sq. between}}{\text{Mean Sq. within.}}$$

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

86.56 > 3.556, Reject the null hypothesis