

Lecture - 07

- P value and Significance Value.
- Distribution.
- Central limit theorem.
- Bernoulli's distribution.
- Binomial distribution.
- Pareto's distribution. [Power law distribution]
- Log Normal distribution.
- F Test (ANOVA)

P value and Significance Value

↳ Derive the P Value.

(Ques. No.)

1. The average weight of all residents in Bangalore City is 168 pound with a standard deviation 3.9. We take a sample of 36 individuals and the mean is 169.5 pounds with 95% of Confidence Interval.

Given,

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

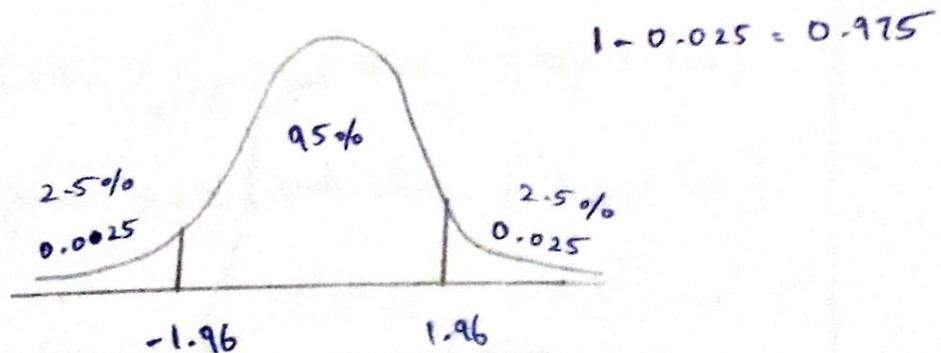
$$\alpha = 0.05$$

1. Null hypothesis, (H_0) $\Rightarrow \mu = 168$

2. Alternate hypothesis, (H_1) $\Rightarrow \mu \neq 168$

3. Significance Value = 0.05

4. Decision boundary, [obviously Two-tail]



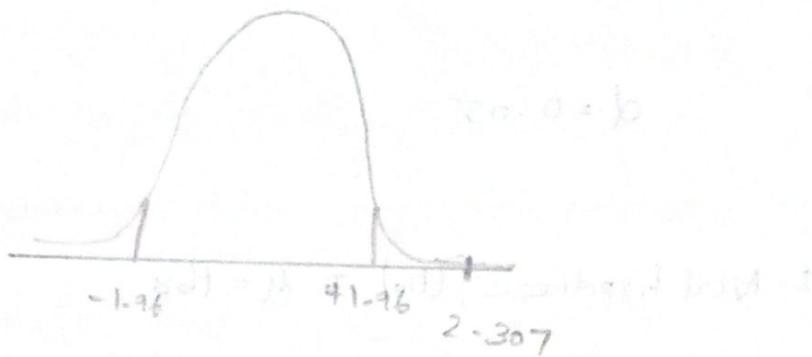
5. Test statistics,

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = \frac{1.5}{3.9} \times 6 = 2.307$$

$$Z = 2.307$$

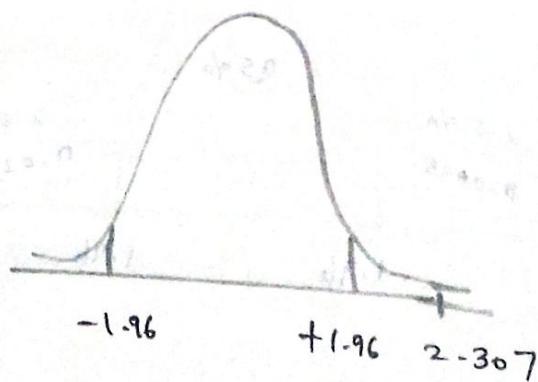
6. State Own decision,

$$Z = 2.307 > 1.96$$



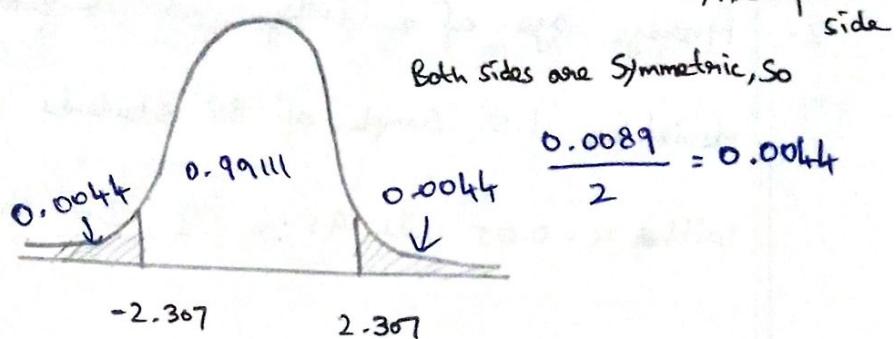
So, obviously reject the null hypothesis.

7. State Own decision with the help of p-value,



Some extend the Curve,

Two Tail,



Add,

↳ See Z score value in Z table

$$0.0044 + 0.99111 + 0.0044 \quad 0.99111$$

$$= 1 \text{ [whole distribution]} \quad (\text{Area Under the Curve})$$

So, the P-value is calculated with the help of two-tail,

$$P \text{ value} = 0.0044 + 0.0044$$

$$\boxed{P_{\text{val}} = 0.0088} \quad [\text{Probability of data fall within the distribution}]$$

We decide,

$P \text{ value} < \text{Significant Value}$ [Reject H_0]

$P \text{ value} > \text{Significant Value}$ [Accept H_0]

$P \text{ value} < 0.05$

$$\boxed{0.0088 < 0.05}$$

[Reject the null hypothesis]

So, The average weight of all residents in Bangalore City is 168 is statistically accepted.

2. Average age of a College is 24 years with a Standard deviation 1.5. Sample of 35 students Mean is 25 years. with $\alpha = 0.05$ with 95% CI, Does the age vary?

Given,

$$\mu = 24 \quad \sigma = 1.5 \quad \bar{x} = 25 \quad n = 35 \quad \alpha = 0.05$$

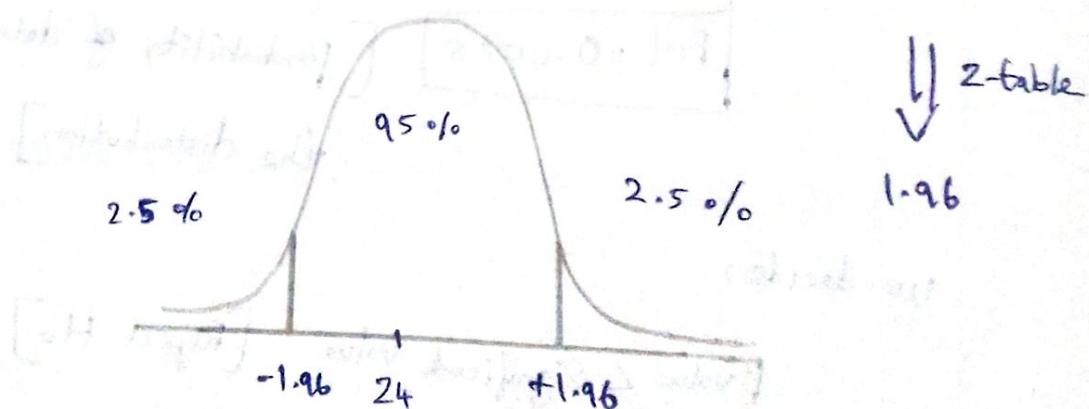
(i) Null hypothesis, $H_0 = \mu = 24$

(ii) Alternate hypothesis, $H_1 = \mu \neq 24$.

(iii) Significance Value = 0.05

(iv) State Decision boundary,

$$1 - 0.025 = 0.975$$



(V) Test statistics (Z-test),

$$\text{Z-Score} = \frac{\bar{x} - \mu}{\sigma}$$

$$\left(\frac{\sigma}{\sqrt{n}} \right) = \frac{25 - 24}{1.5} \times \sqrt{35}$$

$$= \frac{1 \times 6}{1.5} = 4$$

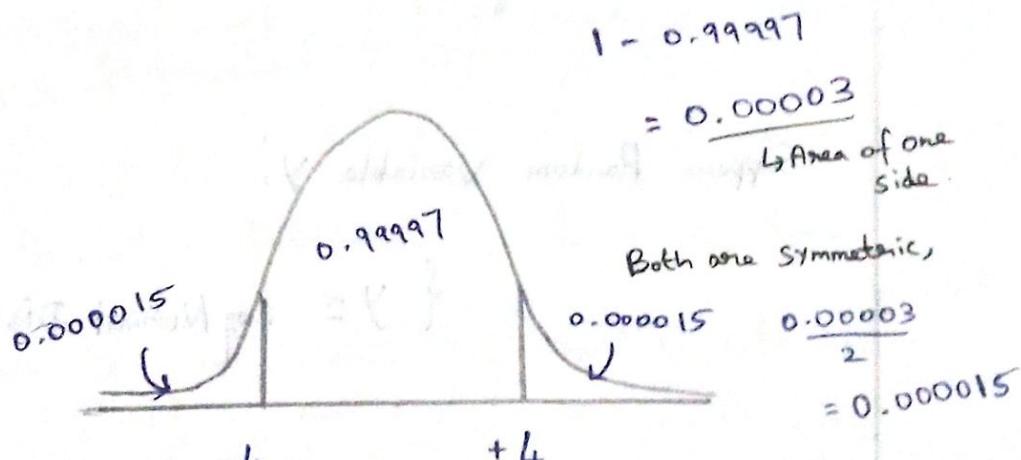
Z-score = 4

(vi) State Own decision,

$$Z = 4 > 1.96$$

So, reject the null hypothesis.

(vii) State Own decision with P-value,



↓ See Zscore value in Z-table

0.99997

(Area Under the Curve)

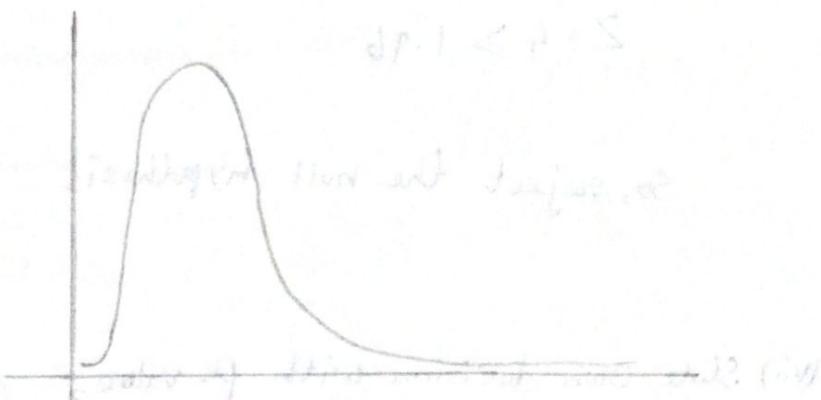
$$P\text{-Value} = 0.000015 + 0.000015$$

$$= 0.00003$$

$P \text{ value} \leq \text{Significance Value}$

So, reject the null hypothesis

Log Normal distribution:



Log Normal distribution

Suppose Random Variable Y ,

$$\{ Y \sim \text{Log Normal Distribution} \}$$

Using log function,

Apply: $\log(Y)$ of the Random Variable Y , It

Converted into Normal distribution.

Suppose we have log normal distribution, We need

to change the distribution to normal distribution for analysis and research.

$$\{y \approx \text{Log Normal Distribution}\}$$



$\log(y)$ $\xrightarrow{\text{Convert to}}$ Normal distribution.

Example:

→ Wealth distribution.

→ People writing big comments, etc.

Bernoulli's Distribution:

Bernoulli's distribution is focus on two outcomes One is Success and another one is failure in the term of probability.

Tossing a Coin,

$$P(H) = 0.5$$

So Consider,

$$P = 0.5 \quad [\text{Probability of Success}] \quad [\text{One chance of outcome}]$$

$$q = 1 - p$$

$$q = 1 - 0.5$$

$q = 0.5$ Probability of failure (or) other chance of outcome

Suppose, we do not a fair coin, [one-tail]

one-tail favor tail

$$P(H) = 0.3$$

$$P(T) = p_{tail} \text{ or } 1 - P(H)$$

$$= 1 - 0.3$$

$P(T) = 0.7$, These type of scenario, we use the Bernoulli's distribution.

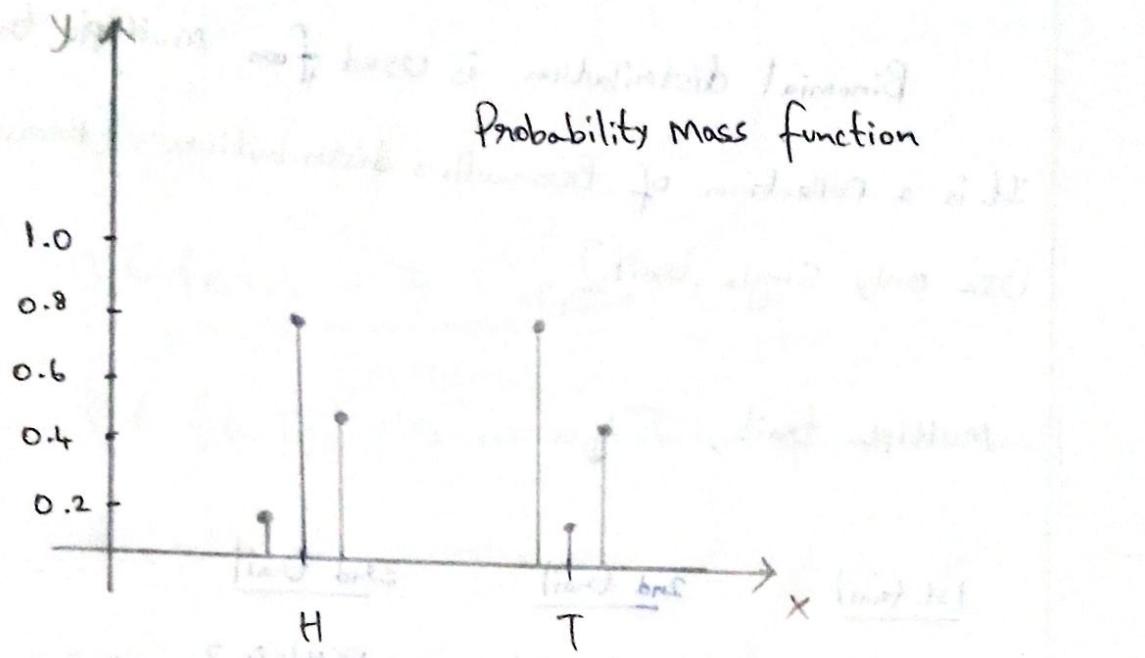
The Bernoulli distribution discovered by Jacob

Bernoulli and this distribution is for discrete variables.

So, it is called as discrete probability distribution of a random variable.

This distribution ensures the PMF (Probability Mass function) because it is discrete supported One and for Categorical Variable. But pdf are used in Continuous Variable. (i.e) Normal distribution.

Probability Mass function of Bernoulli's



$$P(x=0) = 0.2 \text{ and } P(x=1) = 0.8$$

$$P(x=0) = 0.8 \text{ and } P(x=1) = 0.2$$

$$P(x=0) = 0.5 \text{ and } P(x=1) = 0.5$$

Parameters, $0 \leq p \leq 1$

$$q = 1 - p$$

formula,

PMF,

$$\begin{cases} q = 1 - p & \text{if } k=0 \\ p & \text{if } k=1 \end{cases}$$

$$p^k (1-p)^{1-k}$$

CDF,

$$\begin{cases} 0 & \text{if } k \leq 0 \\ 1-p^k & \text{if } 0 \leq k < 1 \\ 1 & \text{if } k \geq 1 \end{cases}$$

Binomial Distribution:

Binomial distribution is used for multiple trials.

It is a collection of Bernoulli's distribution. [Because it uses only single trial].

Multiple trial, \rightarrow

1st trial

$$P(H) = 0.5$$

$$P(T) = 0.5$$

2nd trial

$$P(H) = 0.6$$

$$P(T) = 0.4$$

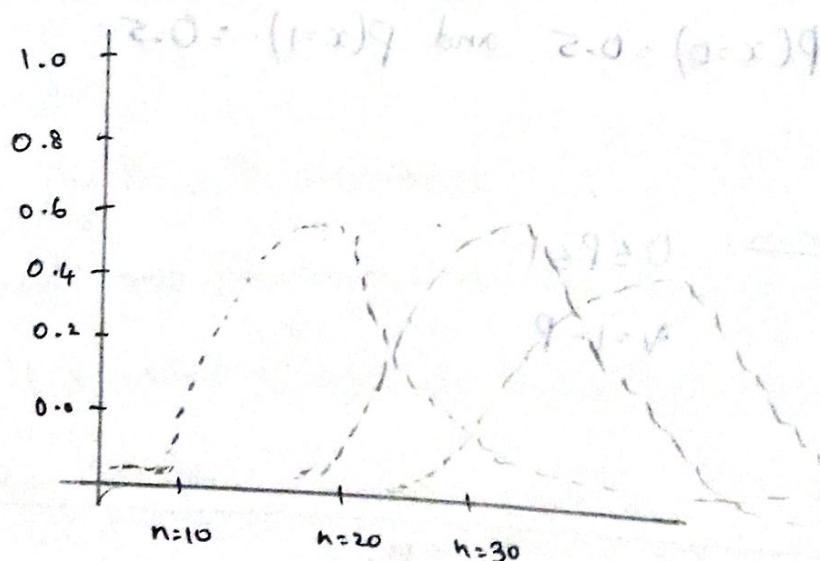
3rd trial

$$P(H) = 0.3$$

$$P(T) = 0.7$$

$$E.O = (1-x)^q \text{ and } E.S = (x-x)^q$$

Probability Mass function = q true $x.3 = (x-x)^q$



Ques 1

is also if $p = 0.5$ and $n = 10$

Ques 2

if $p = 0.5$ and $n = 20$

$p = 0.4$ and $n = 30$

$(q-1)^{n-1}$

$(q-1)^{n-1}$

Cumulative Mass function graph based on the Cumulative
Sum of the each Sample.

Parameters,

$n \in \{0, 1, 2, \dots\}$ number of trials

$p \in \{0, 1\} \rightarrow$ Success probability for each trial

$$q = 1 - p$$

formula,

$$\binom{n}{k} p^k q^{n-k}$$

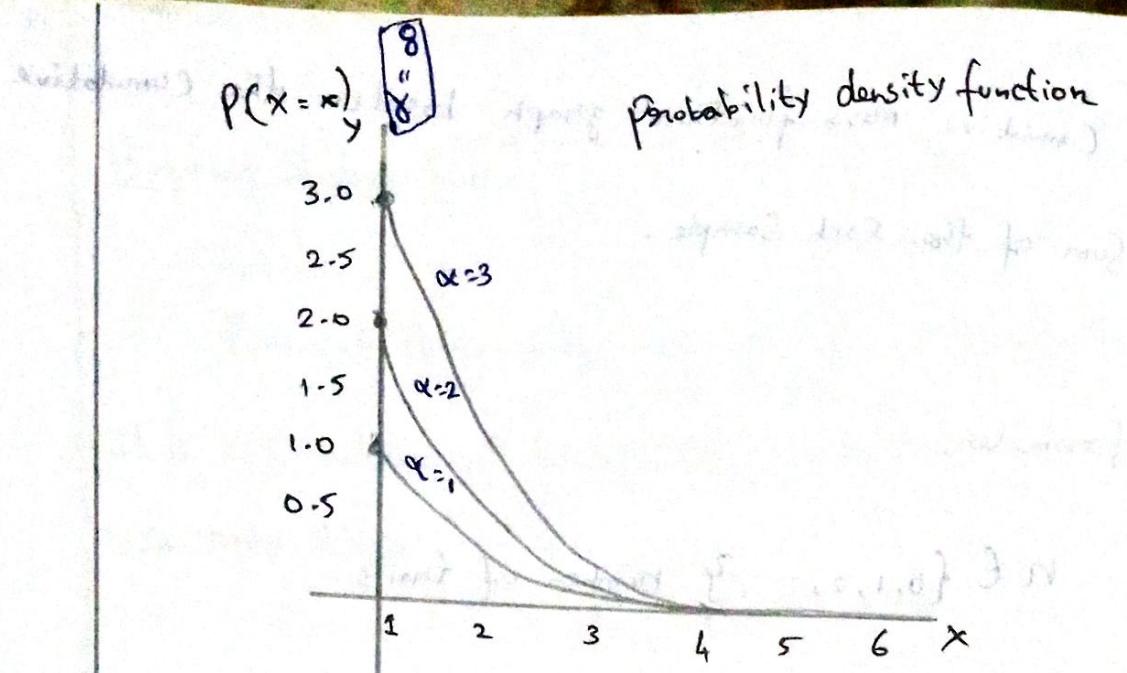
Pareto Distribution: [Non-Gaussian distribution]

→ Pareto distribution is a power-law probability

distribution that is used in description of social,
quality, control, scientific and many type of

Observable Phenomena.

The structure would be an,



This structure of distribution is called as the Pareto distribution.

The 80 - 20 rule is applicable for Pareto distribution.

Parameter,

$x_m > 0$ Scale (real)

$\alpha > 0$ Shape (real)

formula,

PDF

CDF

$$\frac{\alpha x_m^\alpha}{x^{\alpha+1}}$$

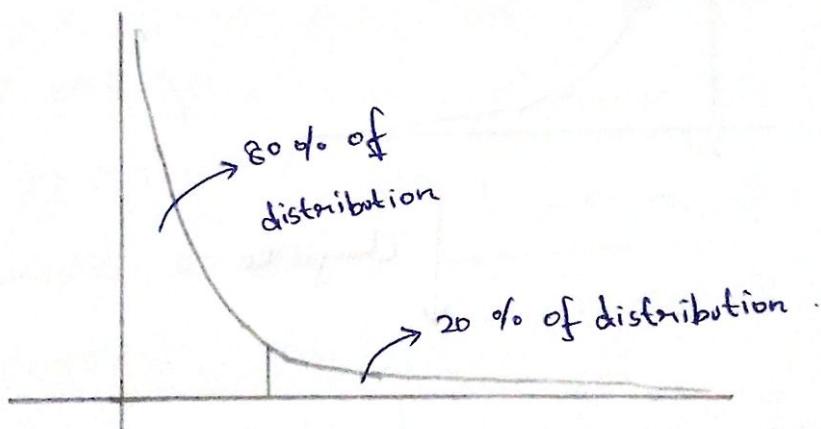
$$1 - \left(\frac{x_m}{x}\right)^\alpha$$

Power-law Distribution: [Application of Pareto distribution]

- Power law is a functional relationship between two quantities.
- It is the application of Pareto distribution.

The structure would be,

80 - 20 rule



Example,

→ 80% of the wealth is distributed with 20% of the people.

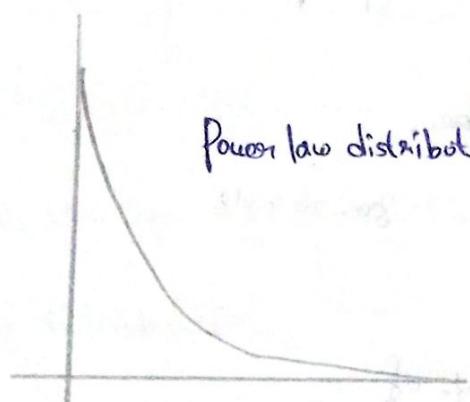
→ 80% of the company project are done by 20% of the people in a team.

→ 80% of sales is done by the 20% of the famous product.

This is also called as Power law (or) Pareto

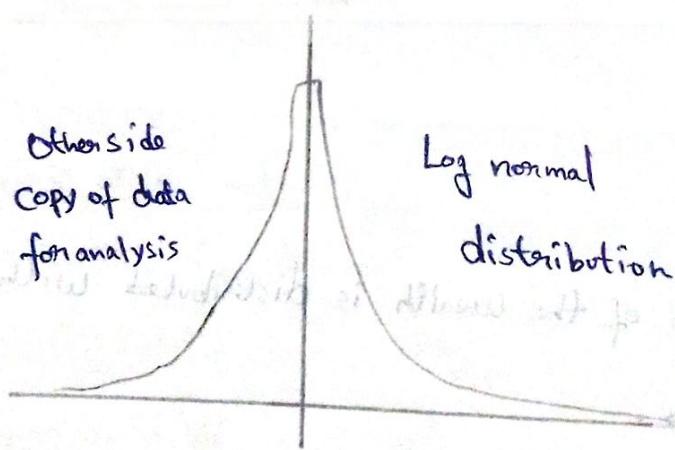
distribution.

One amazing thing is that we can change this distribution to Normal distribution for analysis.



Power law distribution

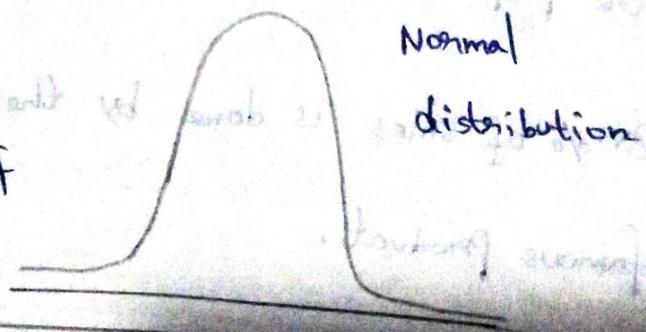
↓
Change to



Otherwise
copy of data
for analysis

Log normal
distribution

Using log
function ↓ Change to

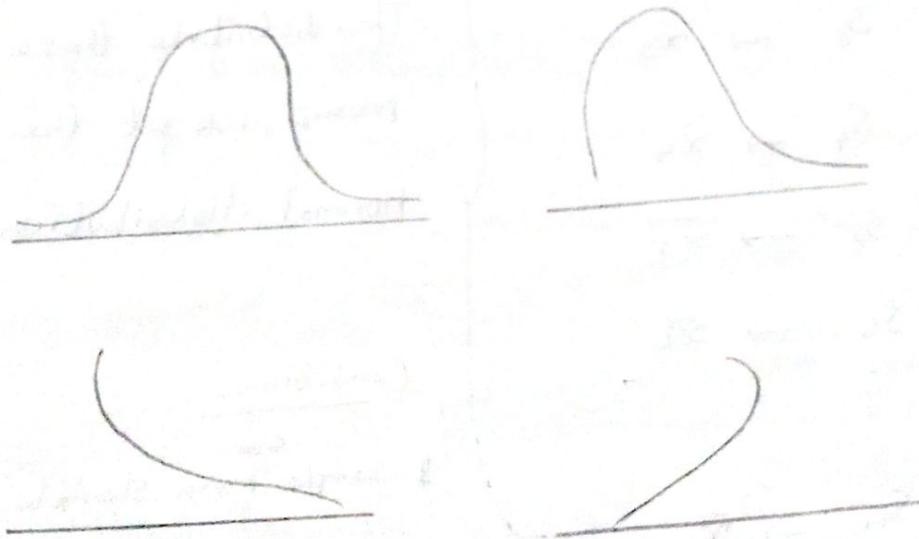


This is called as
Transformation of
data.

Normal
distribution

Central limit theorem:

Central limit theorem says that, In any distribution of data take ^{Mean} multiple samples ^{with Sample Size} and Convert is greater than 30. Then say that the and Convert to normal distribution for analysis.

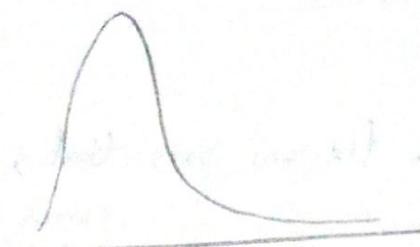


ANY DISTRIBUTION

↓
CONVERT TO [with the help of
taking samples]
Normal distribution

Example:

Consider the any type of distribution,



Take Some Samples

$S_1 \xrightarrow{\text{Find mean}} \bar{x}_1$

$S_2 \xrightarrow{\text{Find mean}} \bar{x}_2$

$S_3 \rightarrow \bar{x}_3$

$S_4 \rightarrow \bar{x}_4$

$S_5 \rightarrow \bar{x}_5$

$S_6 \rightarrow \bar{x}_6$

⋮
⋮

$S_m \rightarrow \bar{x}_m$

Then distribute these Sample

Means, we get the

normal distribution.

Conditions:

1. Sample Mean Should be $n \geq 30$

2. Samples are allowed ^{any} multiple

Samples. [Higher samples \rightarrow Good result].

→ Normal distribution.

Sample mean

Poisson distribution:

- Poisson distribution is a discrete probability distribution for the counts of events that occur randomly over a given period.
- Many experimental situations occur in which we observe the counts of events within the set units of time, area, volume etc.
- Discovered by French mathematician Simeon Denis Poisson in 1837.

Note,

Poisson distribution is used in cases where the chances of any individual event being a success is very small.

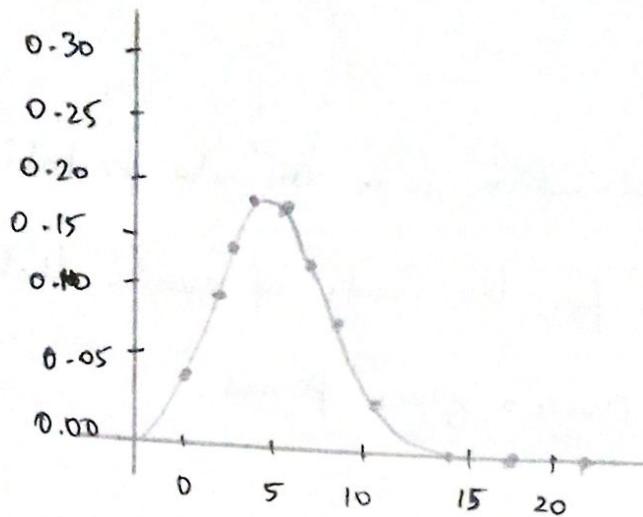
Eg:

The number of plane crash in India in one year.

$P(x=x)$

Probability Mass function

$$\lambda = 4$$



Parameter, $\lambda \in (0, \infty)$ (rate)

Formula,

The Probability of observing x events in a given interval is given by,

$$P(x=x) = \frac{e^{-\lambda} \lambda^x}{x!} \quad (x=0, 1, 2, \dots)$$

$$R \text{ Value} = 2.7182$$

It is also a kind of Pareto distribution.

Example Problem:

1. On average Cancer kills five people each year in India, $\lambda=5$. What is the probability that only one person is

killed this year?

Assuming these are independent random events, the number of people killed in a given year therefore has a Poisson distribution.

$$\lambda = 5$$

Let X be the number of people killed in a year,

$$P(X=k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

$$P(X=1) = \frac{e^{-5} 5^1}{1!} \approx 0.033$$

↑
One person
Probability

2. Suppose the trains arrive at the railway station with an average arrival rate of 4 trains per hour. What is the probability exactly 6 trains will arrive in a two hour period?

$$\lambda = 4 \quad k = 6$$

$$P(X=6) = \frac{2.7182^{(-4)} 4^6}{6!} = 0.104067$$

$$P(X=6) \approx 0.104067 \text{ (or) } 10.41\%$$

Uniform distribution:

A Uniform distribution refers to the probability distribution of discrete and continuous variables. It is applicable for Pdf and Pmf. It describes in which all values within a given range are equally likely to occur.

The same probability of being observed.

Example:

Rolling a fair Six-Sided Die

The die is Unbiased so all six sides have the same chance of landing face up. Each outcome (1, 2, 3, 4, 5, 6) is equally likely.

In this case,

$$P(x=1) = \frac{1}{6}$$

$$P(x=2) = \frac{1}{6}$$

$$P(x=3) = \frac{1}{6}$$

$$P(x=4) = \frac{1}{6}$$

$$P(x=5) = \frac{1}{6}$$

$$P(x=6) = \frac{1}{6}$$

Another Example:

Random Number Generation.

Random Number Generation has equally chance to create
(0 to 1) - i.e $P(0 \leq x \leq 1) = 1$.

Formula,

$$\text{PMF, } P_{\text{mf}} = \frac{1}{n}$$

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

Exponential Distribution:

→ The exponential distribution is a Continuous Probability

distribution that models the time between events

in a Poisson process.

→ It is commonly used to describe the time between

occurrence of continuous and memory less event.

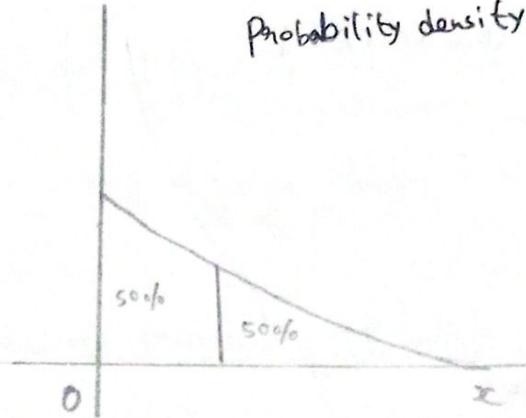
Pdf,

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

Cdf,

$$F(x/\lambda) = 1 - e^{-\frac{x}{\lambda}}$$

probability density function



Interarrival time:

It refers to the time elapsed between consecutive occurrence of events. It is the time gap or duration between one event happening and next event occurring.

Example: One customer arrive at 10 A.M, Another Customer from 10 to 11 AM like that.

Memoryless Events:

The Memoryless event means that the past does not affect the future events.

Example: Bus Waiting, Probability of head in fair coin.

Two types of problem in Exponential distribution:

Method 1: [with Interarrival time]

Time between Website Pages requests.

Suppose you are an IT analyst monitoring the website traffic for a ~~popular~~ Online store. You record the time stamps of user page requests and find the following time Interval (Interarrival time) between consecutive page requests. [0.2, 0.3, 0.1, 0.4, 0.2, 0.3, 0.5, 0.2, 0.3, 0.4]. So, what is the probability of next page request will occur within 5 seconds?

Soln:

Step 1: Calculate the average interarrival time.

$$\text{Average Interarrival time} = \frac{0.2 + 0.3 + 0.1 + 0.4 + 0.2 + 0.3 + 0.5 + 0.2 + 0.3 + 0.4}{10}$$

$$= 0.3 \text{ Seconds.}$$

Step 2: fit the exponential distribution.

fit the exp dist, 0.3 seconds as the rate Parameter.

Step 3: Probability Calculation

$$\text{find } P(X \leq 0.5) = 1 - e^{(-\lambda t)}$$

$$= 1 - e^{(-0.3 \times 0.5)}$$

$$= 1 - e^{-0.15}$$

$$\approx 0.8607$$

$$\approx 0.1393$$

$$P(X \leq 0.5) \approx 0.1393 \text{ (or) } 13.93\%$$

why use Cdf formula?
because the value of λt is less than 1.
 (0.3×0.5)

So, use Cdf formula,
Otherwise use Pdf
formula.

Method 2: [without interarrival time]

Suppose you are IT analyst monitoring traffic for popular Online store. You observe the time between consecutive page requests made by users on the website, and you find that the time between page requests follows an exponential distribution with an average rate of 0.1 requests per second. ($\lambda = 0.1$ request per second).

What is the probability that the next page request will occur within 5 seconds?

Directly give $\lambda = 0.1$

Step 1: Convert the time to same unit of rate parameter.

$$\lambda = 0.1 \text{ (Rate parameter)}$$

$$x = 0.5 \text{ seconds}$$

Step 2: Calculate probability,

$$P(x \leq 5) = 1 - e^{(-0.1 + 0.5)}$$

$$= 1 - e^{(-0.05)}$$

$$= 1 - (0.9512294245)$$

$$= 0.0488$$

$$P(x \leq 5) = 0.0488 \text{ (or) } 4.8\%$$

Method 1 Question Suppose the interarrival time is 0.6 and

Say probability of next page request within 0.5 seconds,

[Use pdf]

So,

$$\lambda = 0.6 \text{ [Rate parameter]}$$

$$x = 0.5 \text{ Sec.}$$

Calculate,

$$\text{P.d.f. } f(x|\lambda) = \lambda e^{(-\lambda x)}$$

$$f(x|0.6) = 0.6 * e^{(-0.6x)}$$

$$P(x \leq 0.5) = \int_{[0, 0.5]} f(x|0.6) dx$$

$$P(x \leq 0.5) = \int_0^{0.5} (0.6 * e^{(-0.6x)}) dx$$

$$P(x \leq 0.5) = \left[-e^{(-0.6x)} \right] \text{ from 0 to } 0.5$$

$$P(x \leq 0.5) = \left[-e^{(-0.6 + 0.5)} \right] - \left[-e^0 \right]$$

$$P(x \leq 0.5) = [-0.30119] - [-1]$$

$$P(x \leq 0.5) \approx 0.69881$$

The probability of the next page request will occur within 0.5 seconds is 0.69881 (or)
69.88%.

- (Q8) This problem also can be solved with the help of Complementarity of Cdf. (ccop).

ANOVA TEST: [F - Test]

Analysis of Variance (ANOVA) is a statistical test for detecting difference in the group means when there is One parametric dependent variable and One or more independent Variable.

Types:

One-way Anova:

One way Anova is used when you have one independent Variable (factor) that divides your data into two or more groups.

Two way Anova:

Two way Anova is used when you have two independent variables (factors) that divides your data into different groups.

In ANOVA, F-statistics will be used.

Formula, [ONE WAY ANOVA]

$$F = \text{MSB} / \text{MSW}$$

MSB → Mean Square Between, Calculated as the Sum of Squares

between groups (SSB) divided by the degree of freedom.

$$\text{MSB} = \text{SSB} / \text{df}_B$$

SSB → Sum of Squares between, the variability between
the means of the group.

df_B → Degree of freedom between, the degrees of freedom
associated with between - group variability.

$$\text{MSW} = \text{SSW} / \text{df}_W$$

SSW → Sum of squares within, the variability
within each group.

df_W → Degrees of freedom within, the degrees of
freedom associated with group variability.

$$SSB = \sum (Group\ Mean - Grand\ Mean)^2 * (\text{Number of Observations in Group})$$

$$df_B = \text{Number of Groups} - 1$$

$$SSW = \sum (\text{Observation} - \text{Group Mean})^2$$

$$df_W = \text{Total number of Observations} - \text{Number of Groups}$$

Formula, [Two way ANOVA]

$$F = (SS_{AB} / df_{AB}) / (SS_W / df_W)$$

↓
↳ Both combine factor formula.

And

$$F = (SS_A / df_A) / (SS_W / df_W)$$

↓
↳ One factor

Example for One way and Two way Anova:

One-way Anova:

1. Scenario:

Imagine you work as a quality control manager in a factory that produces three different types of cookies: A, B and C. You want to determine if there are any significant differences in the average diameter of the cookies produced by these three types.

Data: you randomly sample 15 cookies from each type (A, B and C) and measure their diameters. The diameter measurements are within 8.0 to 9.0.

Soln:

Collect Samples,

Type A: 8.3, 8.4, 8.7, 8.9, 8.2, 8.8, 8.6, 8.4, 8.5, 8.7, 8.1, 8.3, 8.6, 8.4

Type B: 8.5, 8.7, 8.3, 8.6, 8.4, 8.2, 8.8, 8.4, 8.9, 8.6, 8.7, 8.5, 8.3, 8.4, 8.6

Type C: 8.8, 8.4, 8.9, 8.5, 8.7, 8.3, 8.6, 8.5, 8.7, 8.1, 8.2, 8.4, 8.6, 8.4, 8.5

Null Hypothesis (H₀):

There is no significant difference in the average diameter of cookies produced by type A, B and C.

Alternate Hypothesis (H₁):

There is a significant difference in the average diameter of cookies produced by at least one of the types A, B and C.

Calculations:

Step 1: Calculate mean for each type.

$$\text{Mean of Type A } (\bar{x}_A) = \frac{(8.3 + 8.4 + 8.7 + \dots + 8.6 + 8.4)}{15}$$

$$\boxed{\bar{x}_A = 8.48}$$

$$\text{Mean of Type B } (\bar{x}_B) = \frac{(8.5 + 8.7 + 8.3 + \dots + 8.4 + 8.6)}{15}$$

$$\boxed{\bar{x}_B = 8.48}$$

$$\text{Mean of Type C } (\bar{x}_C) = \frac{(8.8 + 8.4 + 8.9 + \dots + 8.4 + 8.5)}{15}$$

$$\boxed{\bar{x}_C = 8.52}$$

Step 2: Calculate Overall mean (\bar{x} -overall) and total

Sum of Squares (SS-total).

$$\bar{x}_{\text{overall}} = \frac{(\bar{x}_A + \bar{x}_B + \bar{x}_C)}{3} = \frac{8.48 + 8.48 + 8.52}{3}$$

$$\boxed{\bar{x}_{\text{overall}} = 8.4933}$$

$$SS_{\text{total}} = \sum (x_i - \bar{x}_{\text{overall}})^2$$

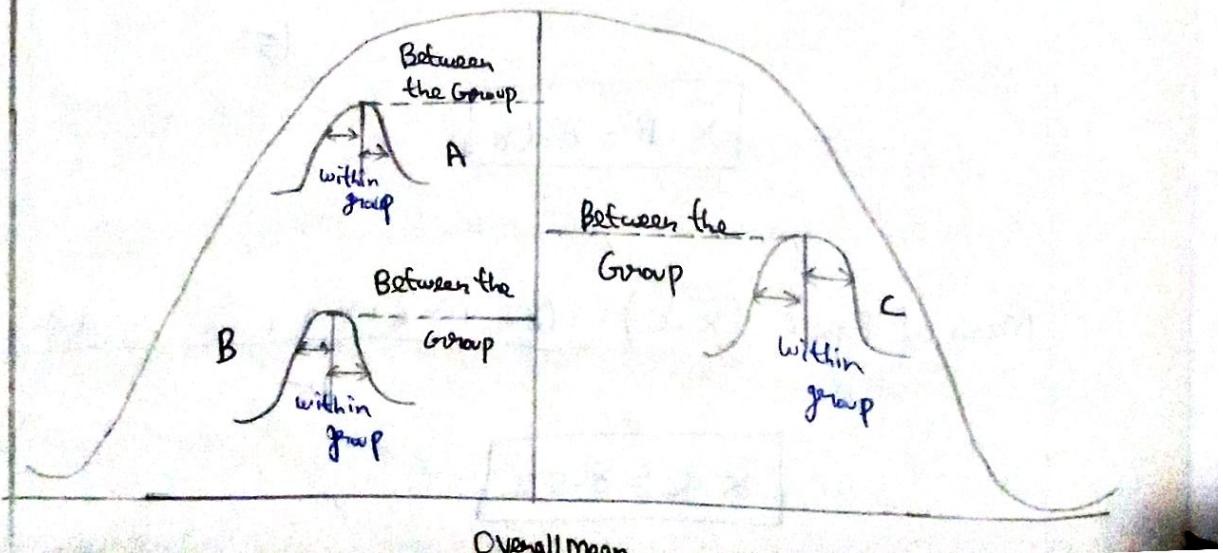
↳ Each data in the three groups.

$$SS_{\text{total}} = (8.3 - 8.4933)^2 + (8.4 - 8.4933)^2 + \dots + (8.4 - 8.4933)^2 + (8.5 - 8.4933)^2$$

Step 3:

Calculate the Sum of Squares between groups (SS-between).

↓ Nothing but
Variability between the group.



$$SS_{\text{between}} = \sum (\text{no of observation}) * (\text{group mean} - \text{grand mean})$$

$$SS_{\text{between}} = n_A * (\bar{x}_A - \bar{x}_{\text{overall}})^2 + n_B * (\bar{x}_B - \bar{x}_{\text{overall}})^2 + n_C * (\bar{x}_C - \bar{x}_{\text{overall}})^2$$

$n_A, n_B, n_C \rightarrow \text{no of observation}$

Sample ($n=15$)

$$= 15 * (8.48 - 8.4933)^2 + 15 * (8.48 - 8.4933)^2 +$$

$$15 * (8.52 - 8.4933)^2$$

$$= 0.0544$$

$$\boxed{SS_{\text{between}} = 0.0544}$$

Step 4:

Calculate the Sum of Squares within groups (SS_{within})

↳ variability within the group.

$$SS_{\text{within}} = \sum (x_i - \bar{x}_i)^2$$

$$SS_{\text{within}} = \sum (\text{for each group data points} - \text{for each respective Group mean})^2$$

For Type A,

$$\bar{x}_A = 8.48$$

$$SS_{\text{within_A}} = (8.3 - 8.48)^2 + (8.4 - 8.48)^2 + \dots + (8.6 - 8.48)^2 + (8.4 - 8.48)^2$$

$$SS_{\text{within-}A} = 0.0328$$

Same like,

$$SS_{\text{within-}B} = 0.1064$$

$$SS_{\text{within-}C} = 0.2144$$

Step 5: Calculate degree of freedom.

$df_{\text{between}} = k - 1$, where k is the no of groups.

$df_{\text{within}} = N - k$, where N is the total no of measurements.

[i.e group A=15, group B=15, group C=15]

$$N = 45$$

$$df_{\text{between}} = 3 - 1 = 2$$

$$df_{\text{within}} = 45 - 3 = 42$$

Step 6: Calculate $\frac{MS_B}{MS_W}$

\downarrow
Mean square between

$$MS_{\text{between}} = \frac{SS_B}{df_B} = \frac{SS_{\text{between}}}{\text{degree of freedom between}}$$

$$= \frac{0.0544}{2} = 0.0272$$

$$MS_{\text{between}} = 0.0272$$

$$MS_{\text{within}} = \frac{SSW}{dfw} = \frac{SS_{\text{within}}}{\text{degree of freedom within}}$$

$$MS_{\text{within_A}} = \frac{0.0328}{14}$$

$$MS_{\text{within_A}} \approx 0.002343$$

Same like,

$$MS_{\text{within_B}} \approx 0.007600$$

$$MS_{\text{within_C}} \approx 0.015314$$

Step 7: Calculate the F-statistic

$$F = \frac{MSB}{MSW}$$

Type A,

$$MSB = 0.0272$$

$$MSW = 0.002343$$

$$F_{\text{-stat_A}} = \frac{0.0272}{0.002343} \approx 11.588$$

$$F_{\text{-stat_A}} \approx 11.588$$

Type B,

$$MSB = 0.0272$$

$$MSW = 0.007600$$

$$F\text{-stat-B} \approx 3.579$$

Type C,

$$MSB = 0.0272$$

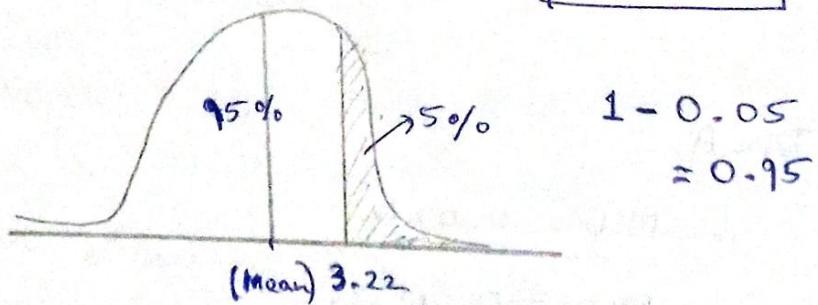
$$MSW = 0.015314$$

$$F\text{-stat-C} \approx 1.776$$

Step-8:

One-tail because [Is any significant difference or not]

$$\alpha = 0.05$$



See F-Test table with first and second degree of freedom,

$df_{\text{between}} = 2$ (see Second column in table)

$df_{\text{within}} = 42$ (see 42nd row in table)

$$(2,42) = 3.22$$

Step 9: [Decision]

~~Topper~~

$$F\text{-stat-A} \approx 11.588 > 3.22 \quad [\text{Reject}]$$

$$F\text{-stat-B} \approx 3.579 > 3.22 \quad [\text{Reject}]$$

$$F\text{-stat-C} \approx 1.776 < 3.22 \quad [\text{Accept}]$$

The Two groups are rejected and only one group is accepted. Therefore Null hypothesis are rejected and Alternate hypothesis are accepted.

Two-way Anova:

It is Quite large calculation. The calculation is much similar to one-way Anova but the factors difference.

Calculation step:

Scenario:

Imagine you are conducting experiment to investigate two factors, "Type of fertilizers" and "Amount of Water" on growth of plants. You have three types of fertilizers (A, B and C) and three levels of water (low, medium and

High). You measure the height of the plants after a certain period to determine if there are any significant difference due to two factors.

fertilizer/water Low Medium High

A 15 20 25

B 18 22 26

C 12 19 24

Step 1:

Mean of each type of Sample.

Type A - fertilizer mean, Type B - fertilizer mean, Type C - fertilizer mean

Type A - water mean, Type B - water mean, Type C - water mean

Step 2: Overall Mean

$$\bar{x}_{\text{grand}} = \frac{15 + 20 + 25 + 18 + 22 + 26 + 12 + 19 + 24}{9} = 20.4444$$

Step 3: SS-total = 107.111

Step 4: SS-between Using formula -

Step 5: SS-within Using formula -

Step 6: degree of freedom (calculate)

Step 7: MSB and MSW

Step 8: F-statistic

Note: Same like apply

water factor every step -

SS-between water

SS-between fertilizer like this

- for all .

[df - interaction] is some very -

[Refer ChatGPT]

CHEBYSHEV'S INEQUALITY:

Let Consider,

$x \approx$ follows Gaussian distribution (μ, σ) .
(Random Variable)

Let Say, Empirical formula for data distribution,

$$P(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$$

$$P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$$

$$P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.7\%$$

Suppose,

$y \approx$ Does not follow Gaussian distribution.
(Random Variable)

How to find distribution respect first, second and third
Sd? that why Chebyshov's inequality come.

formula,

$$P(\mu - k\sigma \leq \text{Random variable} \leq \mu + k\sigma) \geq \boxed{1 - \frac{1}{k^2}}$$

$k \rightarrow$ Range of Standard deviation.

We want to find Range = 2,

$$\boxed{k=2}$$

$$P(\mu - 2\sigma < y < \mu + 2\sigma) > 1 - \frac{1}{2^2}$$

$$= 1 - \frac{1}{4} = \frac{3}{4}$$

= 75%

The Second Standard deviation of Random Variable y Contains the 75% of total distribution of data.

Some Assumption for Each test: [JUST NOTES ALREADY DISCUSSED]

Z-TEST	T-TEST	CHI-SQUARE TEST	ANOVA TEST
<ul style="list-style-type: none"> → Random Sampling. → Normality (data follow G.D.). → Population SD is known. → Sample Size $n > 30$ is produce Super result. 	<ul style="list-style-type: none"> ONE Sample, → follow G.D. → Random Sample from population. → Data points are independent. Two Samples, → follow G.D. → Random Sample. → Homoscedasticity. 	<ul style="list-style-type: none"> → Independence [data is independent - nice to each other] → Sample Size at least 5. → Random Sampling. → Data should be Categorical, 	<ul style="list-style-type: none"> → Independence between each group. → Normality [data follow G.D.] → Homogeneity of variance (Homoscedasticity) If means two group variance is approx similar.