

## Analysis of Variance (ANOVA):

It is the separation of variance ascribable to one group of causes from the variance ascribable to other group. ANOVA consists in the estimation of the amount of variation due to each of the independent factors (causes) separately and then comparing these estimates due to assignable factors with the estimates due to chance factors, the latter being known as experimental error.

1. ANOVA I way classification

2. ANOVA II way classification.

### 1. ANOVA I way classification:

Formulae:

$$RSS = \text{Raw sum of squares} = \sum_{i=1}^K \sum_{j=1}^n n_{ij}^2$$

$$G = \text{Grand Total} = \sum_{i=1}^K \sum_{j=1}^n n_{ij}$$

$$CF = \text{Correction factor} = \frac{G^2}{N}; N = \text{Total no. of observations.}$$

$$TSS = \text{Total sum of squares} = RSS - CF$$

$$SST = \text{Treatment sum of squares} = \sum_{i=1}^K \left( \frac{T_{i.}^2}{n_i} \right) - CF$$

$$SSE = \text{Error Sum of Squares} \quad (\text{Sum of squares due to error})$$

$$= TSS - SST$$

ANOVA 1 way Table

| Source of Variation | Degrees of Freedom (df) | Sum of Squares | Mean Sum of Squares     | Variance Ratio                            |
|---------------------|-------------------------|----------------|-------------------------|---|
| Treatments          | K-1                     | SST            | $MST = \frac{SST}{K-1}$ |   |
| Error               | n-K                     | SSE            | $MSE = \frac{SSE}{n-K}$ | $F = \frac{MST}{MSE} \sim F_{(K-1)(n-1)}$ |
| TSS                 | n-1                     |                |                         |   |

Problem:

The following table shows the time (in hours) of four batches of electric lamps.

| Batches (K) | Life of bulbs in hrs |      |      |      |      |      |      |      |
|-------------|----------------------|------|------|------|------|------|------|------|
| 1           | 1600                 | 1610 | 1650 | 1680 | 1700 | 1720 | 1800 | -    |
| 2           | 1580                 | 1640 | 1640 | 1700 | 1750 | -    | -    | -    |
| 3           | 1460                 | 1550 | 1600 | 1620 | 1640 | 1660 | 1740 | 1820 |
| 4           | 1510                 | 1520 | 1530 | 1570 | 1600 | 1680 | -    | -    |

Perform an analysis of variance of the data and show that a significance test does not reject their homogeneity.

$$H_0: M_1 = M_2 = M_3 = M_4 \rightarrow \text{Null Hypothesis}$$

$$H_1: M_1 \neq M_2 \neq M_3 \neq M_4 \rightarrow \text{Alternative Hypothesis}$$

Coding method: Subtract each number from 1600 and divide by 10  
 (mark repeated numbers)

|   | $T_i$ |    |    |     |     |     |     |     |                             |
|---|-------|----|----|-----|-----|-----|-----|-----|-----------------------------|
| 1 | 0     | -1 | -5 | -8  | -10 | -12 | -20 | -   | -56                         |
| 2 | 2     | -4 | -4 | -10 | -15 | -   | -   | -   | -31                         |
| 3 | 14    | +5 | 0  | -2  | -4  | -6  | -14 | -22 | -29                         |
| 4 | 9     | 8  | 7  | 3   | 0   | -8  | -   | -   | 19                          |
|   |       |    |    |     |     |     |     |     | <u><math>G = -97</math></u> |

|   | $x_{ij}^2$                 | $n_i$    |
|---|----------------------------|----------|
| 1 | 734                        | 7        |
| 2 | 361                        | 5        |
| 3 | 957                        | 8        |
| 4 | 267                        | 6        |
|   | $RSS = \sum \sum x_{ij}^2$ | $N = 26$ |

$$N = 26$$

$$RSS = 2319$$

$$G = \sum \sum x_{ij} = -97$$

$$\text{Correction factor} = \frac{G^2}{N}$$

$$C.F = \frac{(-97)^2}{26} = 361.8846$$

$$\approx 2319$$

$$\text{Total sum of squares} = TSS = RSS - CF$$

$$= 2319 - 361.8846$$

$$= 1957.1154$$

A)

Sum of squares due to treatments

$$\begin{aligned}
 SST &= \sum_{i=1}^4 \left( \frac{T_i^2}{n_i} \right) - CF \\
 &= \left[ \frac{(-56)^2}{7} + \frac{(-31)^2}{5} + \frac{(-29)^2}{8} + \frac{(19)^2}{6} \right] - 361.8846 \\
 &= 443.6071
 \end{aligned}$$

Sum of squares due to errors

$$\begin{aligned}
 SSE &= TSS - SST \\
 &= 1957.1154 - 443.6071 \\
 &= 1513.5083
 \end{aligned}$$

ANOVA 1 way Table:

| Source of variation | Degrees of freedom df            | Sum of Squares    | Mean sum of squares                                       | Variance ratio                             |
|---------------------|----------------------------------|-------------------|---|--|
| Treatments          | $K-1$<br>$\Rightarrow 4-1 = 3$   | $SST = 443.6071$  | $MST = \frac{SST}{K-1}$<br>$= 443.6071/3$<br>$= 147.8690$ | $F = \frac{MST}{MSE}$                      |
| Error               | $N-K$<br>$\Rightarrow 26-4 = 22$ | $SSE = 1513.5083$ | $MSE = \frac{SSE}{N-K}$<br>$= 68.7952$                    | $= \frac{147.8690}{68.7952}$<br>$= 2.1494$ |
| TSS                 | $N-1$<br>$\Rightarrow 26-1 = 25$ |                   |   |  |

Table value of F at 5% level of significance (5)  
at (3,22) df :  
is 3.0491

calculated value < Table value

2.1493      3.0491

Accept  $H_0$

Problem:

4 salesmen were posted in different areas by a company. The no of units of commodity X sold by them are as follows. Is there a significant difference in performance of salesmen?

TB = 3.4903

| (K) | A  | B  | C  | D  |
|-----|----|----|----|----|
| A   | 20 | 23 | 28 | 29 |
| B   | 25 | 32 | 30 | 21 |
| C   | 23 | 28 | 35 | 18 |
| D   | 15 | 21 | 19 | 25 |

Sol  $H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4$   $\rightarrow$  Depends on no of cases (i.e A,B,C,D)

$H_1 : \mu_1 \neq \mu_2 \neq \mu_3 \neq \mu_4$

b) Coding Method: Subtract each number from 28

|   | $T_i$      | $\sum n_{ij}^2$ | $n_i$                  |
|---|------------|-----------------|------------------------|
| A | +8 +5 0 -1 | 12              | 90                     |
| B | 3 -4 -2 7  | 4               | 78                     |
| C | 5 0 -7 10  | 8               | 174                    |
| D | 13 7 9 3   | 32              | 308                    |
|   |            | G = 56          | $\frac{N=16}{RSS=650}$ |

$$\text{Here } N=16, G=56, RSS=650$$

$$CF = \frac{G^2}{N} = \frac{(56)^2}{16} = 196$$

$$\text{Total sum of squares} = TSS = RSS - CF$$

$$TSS = 650 - 196$$

$$\rightarrow TSS = 454$$

Sum of squares due to treatments

$$SST = \sum_{i=1}^4 \left( \frac{T_i^2}{n_i} \right) - CF$$

$$= \left[ \frac{(12)^2}{4} + \frac{(4)^2}{4} + \frac{(8)^2}{4} + \frac{(32)^2}{4} \right] - 196$$

$$= 312 - 196$$

$$SST = 116$$

(7)

Sum of squares due to errors

$$\begin{aligned} SSE &= TSS - SST \\ &= 454 - 116 \end{aligned}$$

$$SSE = 338$$

ANOVA 1 way table:

| source of variation | Degrees of freedom df            | Sum of Squares | Mean Sum of Squares                                     | Variance Ratio                            |
|---------------------|----------------------------------|----------------|---|---|
| Treatments          | $K-1$<br>$\Rightarrow 4-1 = 3$   | $SST = 116$    | $MST = SST/K-1$<br>$= 116/3$<br>$= 38.6666$             | $F = \frac{MST}{MSE}$                     |
| Error               | $N-K$<br>$\Rightarrow 16-4 = 12$ | $SSE = 338$    | $MSE = \frac{SSE}{N-K}$<br>$= \frac{338}{12} = 28.1666$ | $= \frac{38.6666}{28.1666}$<br>$= 1.3727$ |
| Total               | $N-1$<br>$\Rightarrow 16-1 = 15$ | $TSS = 454$    |   |   |

Table value of F at 5% level of significance

at  $(3, 12)$  df is  $3.4903$

Calculated value < Table value

1.3727                    3.4903

Accept  $H_0$

Inference: There is no significant difference in the performance of the 4 salesman.

Q. 2. ANOVA II way classification:

Rows  $\rightarrow K$

Columns  $\rightarrow h$

Formulae:

$$CF = \frac{G^2}{N}, \quad G = \sum_{i=1}^K \sum_{j=1}^h y_{ij}$$

$$TSS = RSS - CF, \quad RSS = \sum_{i=1}^K \sum_{j=1}^h y_{ij}^2$$

$$SST = \frac{1}{h} \sum T_i^2 - CF; \quad T_i = \sum_{j=1}^h y_{ij}$$

Sum of squares due to varieties  $SSV$

$$SSV = \frac{1}{K} \sum T_j^2 - CF, \quad T_j = \sum_{i=1}^h y_{ij}$$

$$SSE = TSS - SST - SSV$$

ANOVA II way table:

| Source of Variation | Degrees of freedom(df) | Sum of Squares | Mean Sum of Squares            | Variance Ratio        |
|---------------------|------------------------|----------------|--------------------------------|-----------------------|
| Treatments          | K-1                    | SST            | $MST = \frac{SST}{K-1}$        | $F = \frac{MST}{MSE}$ |
| Varieties           | h-1                    | SSV            | $MSV = \frac{SSV}{h-1}$        | $F_{(K-1), (h-1)}$    |
| Errors              | (K-1) x (h-1)          | SSE            | $MSE = \frac{SSE}{(K-1)(h-1)}$ | $F = \frac{MSV}{MSE}$ |
| Total               | hk-1                   | -              | -                              | -                     |

problem:

(9)

A tea company appoints 4 salesmen A, B, C & D and observes that sales in 3 seasons - summer, winter and monsoon. The figures (in lakhs) are given in the following table.

| Seasons       | A  | B  | C  | D  | Seasons Total |
|---------------|----|----|----|----|---------------|
| Summer        | 36 | 36 | 21 | 35 | 128           |
| Winter        | 28 | 29 | 31 | 32 | 120           |
| Monsoon       | 26 | 28 | 29 | 29 | 112           |
| Salemen Total | 90 | 93 | 81 | 96 | 360           |

Hypo

$$H_{01}: \mu_1 = \mu_2 = \mu_3$$

$$H_{02}: \mu_1' = \mu_2' = \mu_3' = \mu_4'$$

$$H_{11}: \mu_1 \neq \mu_2 \neq \mu_3$$

$$H_{12}: \mu_1' \neq \mu_2' \neq \mu_3' \neq \mu_4'$$

Coding Method: Subtract each value from 29

|                                 | A  | B  | C  | D  | T <sub>i</sub>                 | T <sub>i</sub> ' <sup>2</sup> | n <sub>i</sub> |
|---------------------------------|----|----|----|----|--------------------------------|-------------------------------|----------------|
| s                               | -7 | -7 | 8  | -6 | -12                            | 144                           | 4              |
| w                               | 1  | 0  | -2 | -3 | -4                             | 16                            | 4              |
| m                               | 3  | 1  | 0  | 0  | 4                              | 16                            | 4              |
| P <sub>j</sub>                  | -3 | -6 | 6  | 9  | G = -12 $\sum T_{i,j}^2 = 176$ | N = 12                        |                |
| T <sub>i,j</sub> ' <sup>2</sup> | 9  | 36 | 36 | 81 | $\sum T_{i,j}^2 = 162$         |                               |                |

Q

$$N = 12$$

$$\text{Rows} = k = 3$$

$$\text{Columns} = h = 4$$

$$G = -12$$

$$CF = \frac{G^2}{N} = \frac{(-12)^2}{12} = 12$$

$$RSS = (-7)^2 + (-7)^2 + 8^2 + (-6)^2 + (1)^2 + (0)^2 + (-2)^2 + \\ (-3)^2 + (3)^2 + 1^2 + 0^2 + 0^2$$

$$RSS = 222$$

$$SST = \frac{1}{h} \sum T_i^2 - CF = \frac{1}{4} (176) - 12 = 32$$

$$SSV = \frac{1}{k} \sum T_{ij}^2 - CF = \frac{1}{3} (162) - 12 = 42$$

$$TSS = RSS - CF \\ = 222 - 12 = 210$$

$$SSE = TSS - SST - SSV = 210 - 32 - 42 = 136$$

### ANOVA II way Table:

| Source of Variation | df                           | Sum of Squares | Mean Sum of Squares                                 | Variance (F) ratio  |
|---------------------|------------------------------|----------------|---|---|
| Treatments          | $k-1$<br>$3-1=2$             | $SST = 32$     | $MST = SST/(k-1)$<br>$= 32/2 = 16$                  | $F = \frac{MS\bar{E}}{MST} = \frac{22.667}{16}$<br>$= 1.4167$   |
| Varieties           | $h-1$<br>$4-1=3$             | $SSV = 42$     | $MSV = SSV/(h-1)$<br>$= 42/3 = 14$                  | $F = \frac{MS\bar{E}}{MSV} = \frac{22.667}{14}$   |
| Error               | $(k-1)(h-1)$<br>$(2)(3)=6$   | $SSE = 136$    | $MS\bar{E} = SSE/(h-1)$<br>$= 136/6$<br>$= 22.6667$ | $= 1.619$   |
| Total               | $hk-1$<br>$= 12-1$<br>$= 11$ |                |   | <p>Here numerator and denominator are reversed because it should be <math>\rightarrow N &gt; P</math></p> |

Here numerator and denominator are reversed because it should be  $\rightarrow N > P$

Since the F ratios are  $< 1$

Accept  $H_0$

Cal val  $<$  Tab val

$1.4167 < 19.33$  Accept  $H_{01}$

Cal val  $<$  Tab Val

$1.619 < 8.94$  Accept  $H_{02}$

Since Num and Deno  
are reversed.

Look up for Num  
F at (6, 2)  $\rightarrow$  Deno  
F at (6, 3) in  
Table values.

→ there is no significant difference between the seasons.

→ there is no significant difference between the performance of salesmen.

### Problem:

The following table represents the sales (Rs 1000)  
per month of 3 brands of soaps allocated among  
three cities.

| Brands | cities |    |    |
|--------|--------|----|----|
|        | A      | B  | C  |
| I      | 12     | 48 | 30 |
| II     | 42     | 54 | 57 |
| III    | 9      | 42 | 21 |

Test whether ① the mean sales of the three brands  
are equal.

② the mean sales of soaps in each city are equal.

$$\begin{array}{ll}
 H_0 : \mu_1 = \mu_2 = \mu_3 & H_{02} : \mu_1' = \mu_2' = \mu_3' \\
 H_{11} : \mu_1 \neq \mu_2 \neq \mu_3 & H_{12} : \mu_1' \neq \mu_2' \neq \mu_3'
 \end{array}$$

Coding method: Subtract each number from 42

|         | A    | B   | C   | $T_i$               | $T_i^2$             | $n_i$ |
|---------|------|-----|-----|---------------------|---------------------|-------|
| I       | 30   | -6  | 12  | 36                  | 1296                | 3     |
| II      | 0    | -12 | -15 | -27                 | 729                 | 3     |
| III     | 33   | 0   | 21  | 54                  | 2916                | 3     |
| $T_j$   | 63   | -18 | +18 | $G=63$              | $\sum T_i^2 = 4941$ | $N=9$ |
| $T_j^2$ | 3969 | 324 | 324 | $\sum T_j^2 = 4617$ |                     |       |

$$\text{Row} \rightarrow k=3 \quad N=9$$

$$\text{Columns} \rightarrow h=3 \quad G=63$$

$$\begin{aligned}
 RSS &= (30)^2 + (-6)^2 + (12)^2 + 0^2 + (-12)^2 + (-15)^2 + (33)^2 + \\ 
 &\quad 0^2 + (21)^2 \\ 
 &= 2979
 \end{aligned}$$

$$CF = \frac{G^2}{N} = \frac{(63)^2}{9} = 441$$

$$SST = \frac{1}{h} \sum T_i^2 - CF = \frac{1}{3} (4941) - 441 = 1206$$

$$SSV = \frac{1}{k} \sum T_j^2 - CF = \frac{1}{3} (4617) - 441 = 1098$$

$$\begin{aligned}
 TSS &= RSS - CF \\ 
 &= 2979 - 441 = 2538
 \end{aligned}$$

$$\begin{aligned}
 SSE &= TSS - SST - SSV \\
 &= 2538 - 1206 - 1098 \\
 &= 234
 \end{aligned}$$

ANOVA II Way Table:

| Source of variation | df                           | Sum of Squares | Mean sum of squares                             | Variance ratio (F)                                      |
|---------------------|------------------------------|----------------|---|---|
| Treatments          | $k-1$<br>$3-1 = 2$           | $SST = 1206$   | $MST = SST/(k-1)$<br>$= 1206/2 = 603$           | $F = \frac{MST}{MSE} = \frac{603}{58.5}$<br>$= 10.3077$ |
| Varieties           | $h-1$<br>$3-1 = 2$           | $SSV = 1098$   | $MSV = SSV/(h-1)$<br>$= 1098/2 = 549$           | $F = \frac{MSV}{MSE} = \frac{549}{58.5}$                |
| Error               | $(k-1)(h-1)$<br>$(2)(2) = 4$ | $SSE = 234$    | $MSE = SSE/(k-1)(h-1)$<br>$= 234/4$<br>$= 58.5$ | $= 9.3846$  |
| Total               | $hk-1$<br>$(3)(3)-1 = 8$     |                |   |   |

→ Table value of F at (2,4) df is at 5% level of significance is 6.94

$$\text{Cal val} > \text{Tab val}$$

Reject  $H_01$

$$\text{Cal val} > \text{Tab val}$$

$$9.3846 > 6.94$$

Reject  $H_02$

Inference: ① The mean sales of three brands significantly differ.

② The mean sales of soaps in each city are not equal.

Problem:

A trucking company wishes to test the average life of each of the four brands of tyres. The company uses all brands on randomly selected trucks. The records showing the lives (1000 of miles) of tyres are given in the table.

Test the hypothesis that the average life of each brand of tyres is the same. Assume  $\alpha = 0.01$  ( $1\cdot1$  table value level of significance)

| Brand I | II | III | IV |  |
|---------|----|-----|----|--|
| 20      | 19 | 21  | 15 |  |
| 23      | 15 | 19  | 17 |  |
| 18      | 17 | 20  | 16 |  |
| 14      | 20 | 17  | 18 |  |
| -       | 16 | 16  | -  |  |

$$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$$

$$H_1: \mu_1 \neq \mu_2 \neq \mu_3 \neq \mu_4$$

|           | ANOVA 1 way<br>method is used |    |    |    |
|-----------|-------------------------------|----|----|----|
| Brand I   | 20                            | 23 | 18 | 17 |
| Brand II  | 19                            | 15 | 17 | 20 |
| Brand III | 21                            | 19 | 20 | 17 |
| Brand IV  | 15                            | 17 | 16 | 18 |

Coding Method: Subtract each number from 17

|     |    |    |    |    | $T_i$     | $\sum n_{ij}^2$ | $n_i$    |
|-----|----|----|----|----|-----------|-----------------|----------|
| I   | -3 | -6 | -1 | 0  | -10       | 46              | 4        |
| II  | -2 | 2  | 0  | -3 | 1         | 18              | 5        |
| III | -4 | -2 | -3 | 0  | 1         | 30              | 5        |
| IV  | 2  | 0  | 1  | -1 | -         | 6               | 4        |
|     |    |    |    |    | $G = -18$ | $RSS = 100$     | $N = 18$ |

$$\text{Here } N = 18, G = -18, RSS = 100$$

$$CF = \frac{G^2}{N} = \frac{(-18)^2}{18} = 18$$

$$TSS = RSS - CF = 100 - 18 = 82$$

$$\begin{aligned}
 SST &= \sum_{i=1}^4 \left( \frac{T_i^2}{n_i} \right) - CF \\
 &= \left[ \frac{60^2}{4} + \frac{(-2)^2}{5} + \frac{(-8)^2}{5} + \frac{(2)^2}{4} \right] - 18 \\
 &= 39.6 - 18 \\
 &= 21.6
 \end{aligned}$$

$$\begin{aligned}
 SSE &= TSS - SST \\
 &= 82 - 21.6 \\
 &= 60.4
 \end{aligned}$$

## ANOVA 1 way Table:

| Source of Variation | df                   | Sum of Squares | Mean Sum of Squares                                 | Variation ratio (F)                  |
|---------------------|----------------------|----------------|---|--------------------------------------|
| Treatment           | $K-1$<br>$4-1 = 3$   | $SST = 21.6$   | $MST = \frac{SST}{K-1}$<br>$= \frac{21.6}{3} = 7.2$ | $F = \frac{MST}{MSE}$                |
| Error               | $N-K$<br>$18-4 = 14$ | $SSE = 60.4$   | $MSE = SSE/N-K$<br>$= 60.4/14$<br>$= 4.3143$        | $= \frac{7.2}{4.3143}$<br>$= 1.6689$ |
| Total               | $N-1$<br>$18-1 = 17$ | $TSS = 82$     |   |                                      |

Table value of F at 10% level of significance  
at (3, 14) df is 5.56

$$\text{Cal Val} < \text{Tab Val}$$

$$1.6689 < 5.56$$

Accept  $H_0$

Inference:

Avg life of each brand of tyres is the same.

## Estimation:

The theory of statistical inference can be divided into two major areas

1. Estimation of parameters.
2. Testing of Hypothesis.

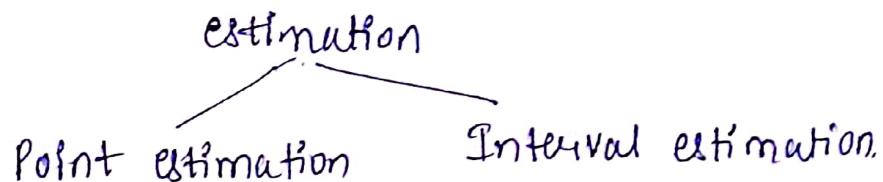
Parameter: A statistical constant derived from the population values is called 'parameter'.

ex: population mean,  $\mu$

population standard deviation,  $\sigma$

Statistic: A statistical constant derived from the sample values is called a 'statistic'

Estimation: It is a procedure of estimating a population parameter by using sample information or observations.



Point Estimate: A point estimate of some population parameter  $\theta$  is a single numerical value.

Point Estimator: A point estimator is a statistic for estimating a population parameter  $\theta$  and is denoted by  $\hat{\theta}$ .

Properties of estimator:

An estimator is said to be a good estimator if it is

- Unbiased
- Consistent
- Efficient
- Sufficient

Unbiased Estimator: A statistic  $\hat{\theta}$  is said to be an unbiased estimator or its value an unbiased estimate if and only if the mean of the sampling distribution of the estimator equals  $\theta$  i.e

$$\text{Mean of } \hat{\theta} = E(\hat{\theta}) = \theta$$

→ If the estimator is not unbiased, the difference  $E(\hat{\theta}) - \theta$  is called the bias of the estimate  $\hat{\theta}$ .

When the estimator is unbiased,

$$E(\hat{\theta}) - \theta = 0$$

i.e the bias is zero.

positive bias: If  $E(\hat{\theta}) > \theta$ , then  $\hat{\theta}$  is called positively biased.

Negative bias: If  $E(\hat{\theta}) < \theta$ , then  $\hat{\theta}$  is called negatively biased.

Problem:

If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are two unbiased estimators of  $\theta$ , then  $w_1\hat{\theta}_1 + w_2\hat{\theta}_2$  is an unbiased estimator of  $\theta$ . when  $w_1 + w_2 = 1$

Sol: Given that  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are unbiased estimator of  $\theta$ ,

$$\text{i.e } E(\hat{\theta}_1) = \theta \quad \text{--- (1)}$$

$$E(\hat{\theta}_2) = \theta \quad \text{--- (2)}$$

$$\begin{aligned} \text{Now } E(w_1\hat{\theta}_1 + w_2\hat{\theta}_2) &= w_1E(\hat{\theta}_1) + w_2E(\hat{\theta}_2) \\ &= w_1\theta + w_2\theta \quad (\text{from (1) \& (2)}) \\ &= (w_1 + w_2)\theta \\ &= \theta \quad (\because w_1 + w_2 = 1) \end{aligned}$$

$\therefore w_1\hat{\theta}_1 + w_2\hat{\theta}_2$  is an unbiased estimator of  $\theta$ .

Theorem: Let  $x_1, x_2, \dots, x_n$  be a random sample drawn from a given population with the mean  $\mu$  and variance  $\sigma^2$ , Show that the sample mean

$\bar{x}$  is an unbiased estimator of population mean  $M$   
i.e.  $E(\bar{x}) = M$

Proof:  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

taking expectations on both sides

$$\begin{aligned} E(\bar{x}) &= E\left(\frac{\sum_{i=1}^n x_i}{n}\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n x_i\right) = \frac{1}{n} [x_1 + x_2 + \dots + x_n] \\ &= \frac{1}{n} \{E(x_1) + E(x_2) + \dots + E(x_n)\} \\ &= \frac{1}{n} [M + M + \dots + M] \quad (\because E(x_i) = M) \\ &= \frac{1}{n} [nM] \\ &= M \end{aligned}$$

### Theorem:

for a random sample of size  $n$ ,  $x_1, x_2, \dots, x_n$  taken from a finite population  $s^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$  is not an unbiased estimator of  $\sigma^2$  but  $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$  is an unbiased estimator of  $\sigma^2$ .

### \* MLE (Minimum Variance Unbiased Estimate)

or most efficient unbiased estimator:

If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are two unbiased estimators

of  $\theta$  and if  $\sigma_1^2$  and  $\sigma_2^2$  are variances of their sampling distributions and  $\sigma_1^2 < \sigma_2^2$ ,  $\hat{\theta}_1$  is said to be more unbiased estimator of  $\theta$ .

Note: If  $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$  are unbiased estimators of the parameter  $\theta$ ,  $\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2$  are variances of their sampling distributions, then  $\hat{\theta}_1$  is the most efficient estimator of  $\theta$  if  $\sigma_1^2 < \sigma_2^2 < \dots < \sigma_k^2$

e = efficiency of an estimator

$$= \frac{\text{var } \theta_i}{\text{var } \hat{\theta}_i} \quad \text{where } i = 2, 3, \dots, k$$

$$e < 1$$

• Consistency: An estimator  $T_n = T(x_1, x_2, \dots, x_n)$  based on a random sample of size  $n$ , is said to be a consistent estimator of  $r(\theta)$ ,  $\theta \in \Theta$  the parameter space, if  $T_n$  converges to  $r(\theta)$  in probability

i.e if  $T_n \xrightarrow[\text{converges}]{} r(\theta)$  as  $n \rightarrow \infty$

In other words,  $T_n$  is a consistent estimator of  $r(\theta)$  if for every  $\epsilon > 0$ ,  $m > 0$ , there exists a positive integer  $n \geq m(\epsilon, m)$  such that

$$P\{|T_n - r(\theta)| < \epsilon\} \rightarrow 1 \text{ as } n \rightarrow \infty$$

$$\Rightarrow P\{|T_n - r(\theta)| < \epsilon\} > 1 - \eta \text{ for } n > m$$

where  $m$  is some very large value of  $n$

. Sufficiency: An estimator is said to be sufficient for a parameter, if it contains all the information in the sample regarding the parameter.

If  $T = t(n_1, n_2, \dots, n_n)$  is an estimator of a parameter  $\theta$  based on a sample  $n_1, n_2, \dots, n_n$  of size  $n$  from the population with density  $f(n, \theta)$  such that the conditional distribution of  $n_1, n_2, \dots, n_n$  given  $T$ , is independent of  $\theta$ , then  $T$  is sufficient estimator for  $\theta$ .

### Methods of estimation:

1. Maximum likelihood method
2. Method of moments
3. Method of minimum chi-square
4. Method of minimum variance
5. Method of least squares
6. Method of Inverse Probability.

Maximum Likelihood Estimation: It was formulated by C.F. Gauss but as a general method of estimation was first introduced by Professor R.A. Fisher.

Likelihood function: let  $x_1, x_2, \dots, x_n$  be a random sample of size  $n$  from a population with density function  $f(x, \theta)$ . Then the likelihood function of the sample values  $x_1, x_2, \dots, x_n$  usually denoted by  $L = L(\theta)$  is their joint density function.

$$L = f(x_1, \theta) f(x_2, \theta) \cdots f(x_n, \theta) = \prod_{i=1}^n f(x_i, \theta)$$

$L$  gives the relative likelihood that the random variables assume a particular set of values  $x_1, x_2, \dots, x_n$ . For a given sample  $x_1, x_2, \dots, x_n$ ,  $L$  becomes a function of the variable  $\theta$ , the parameter.

The principle of maximum likelihood estimation is finding an estimator for the unknown parameter  $\theta = (\theta_1, \theta_2, \dots, \theta_k)$  which maximises the likelihood function  $L(\theta)$  for variations in parameters.

If there exists a function  $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$  of the sample values which maximises  $L$  for variations in  $\theta$ , then  $\hat{\theta}$  is to be taken as an estimator of  $\theta$ .  $\hat{\theta}$  is usually called maximum likelihood estimator.

thus  $\hat{\theta}$  is the solution of

$$\frac{\partial L}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L}{\partial \theta^2} < 0 \quad -\textcircled{2}$$

$L_{\textcircled{1}}$

since  $L > 0$  and  $\log L$  is non decreasing function of  $L$ ,  $L$  and  $\log L$  attain their extreme values (maxima or minima) at the same value of  $\hat{\theta}$

① can be rewritten as

$$\frac{1}{L} \frac{\partial L}{\partial \theta} = 0 \Rightarrow \frac{\partial \log L}{\partial \theta} = 0$$

this is a convenient form from practical view. If  $\theta$  is a vector valued parameter, then  $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3, \dots, \hat{\theta}_k)$  is given by the solution of simultaneous equations

$$\frac{\partial}{\partial \theta_i} \log L = \frac{\partial}{\partial \theta_i} \log L(\theta_1, \theta_2, \dots, \theta_k) = 0 \quad -\textcircled{3}$$

$i = 1, 2, \dots, k$

eqns ①, ② & ③ are usually referred to as likelihood eqns for estimating the parameters

Date:  
17/12/2020

## UNIT 2 : Testing of Hypothesis

Hypothesis is a statement regarding the population parameter. Hypothesis is of two types

1. Null Hypothesis
2. Alternative Hypothesis

Null Hypothesis: It is a no difference statement. It is denoted by  $H_0$ .

$$\text{ex: } H_0: \mu = 3$$

$$H_0: \mu_1 = \mu_2$$

$$H_0: \sigma^2 = 2$$

Alternative Hypothesis: It is complementary to Null Hypothesis. It is denoted by  $H_1$ . It is of two types. ex:  $\mu_1 \neq \mu_2$ ;  $\sigma_1 \neq \sigma_2$ ;  $\mu > 2 \Rightarrow H_1$

1. Two Tailed Alternative
2. Single Tailed alternative

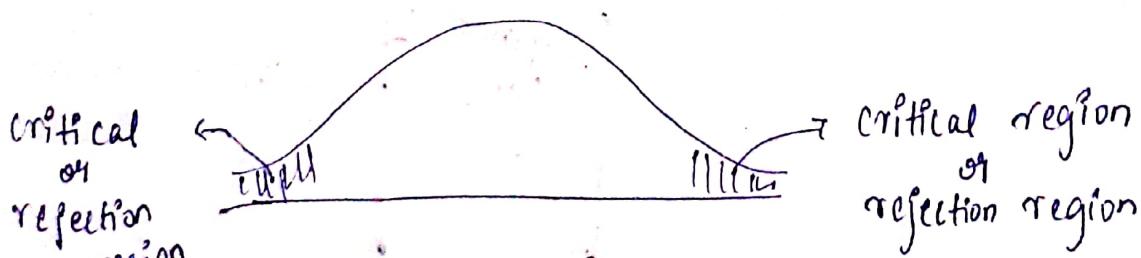
Two Tailed Alternative:

$$\text{ex: } H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

$$\text{ex: } H_0: \sigma^2 = 30$$

$$H_1: \sigma^2 \neq 30$$



## Single Tailed Alternative:

It is divided into two types. They are

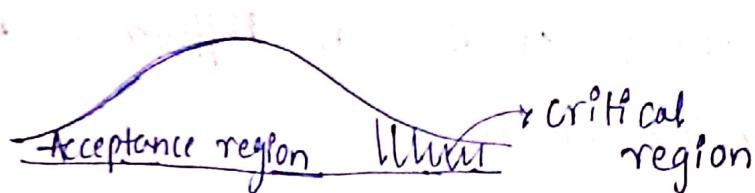
### 1. Right Tailed Alternative

$$\text{ex: } H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 > \mu_2$$

$$\text{ex: } H_0: \sigma^2 = 30$$

$$H_1: \sigma^2 > 30$$



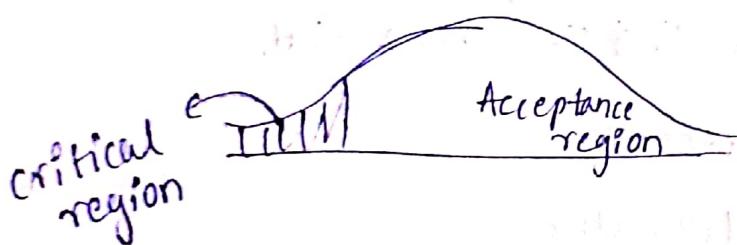
### 2. Left Tailed Alternative

$$\text{ex: } H_0: \mu_1 = \mu_2$$

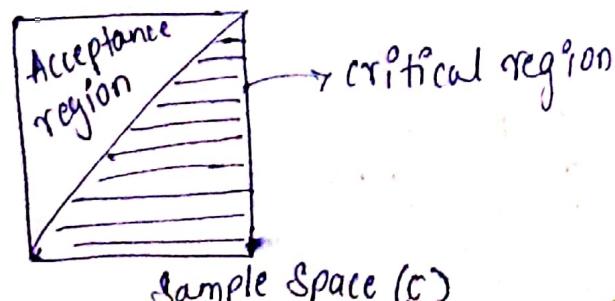
$$H_1: \mu_1 < \mu_2$$

$$\text{ex: } H_0: \sigma^2 = 30$$

$$H_1: \sigma^2 < 30$$



- Critical region: A region in the sample space ' $S$ ' which amounts to rejection of  $H_0$  is called the critical region or the rejection region. It is denoted by  $W$ .



Level of Significance: If  $\omega$  is the critical region and  $t = f(n_1, n_2, n_3 \dots n_n)$  is the value of the statistic based on random sample of size 'n', then

$$P(t \in \omega / H_0) = \alpha$$

Here  $\alpha$  is the level of significance and  $\omega$  is the critical region and

$$P(t \in \bar{\omega} / H_1) = \beta$$

$$\omega \cup \bar{\omega} = S$$

$$\omega \cap \bar{\omega} = \emptyset$$

Critical value or Significant value: The value of the test statistic which separates the critical region and the acceptance region is called critical value or significant value. It depends upon:

1. Level of Significance
2. Alternative Hypothesis



\* Types of Errors: There are two types of errors

1. Type I error
2. Type II error

1. Type I error: Probability of rejecting a good lot is called Type I error.

$$P(\text{rejecting a good lot}) = \alpha$$

This is called producers risk

2. Type II error: Probability of accepting a bad lot is Type II error.

$1 - \beta$  is called power of the test.

$$P(\text{accepting a bad lot}) = \beta$$

This is called consumer's risk

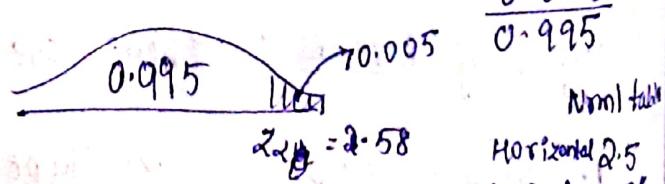
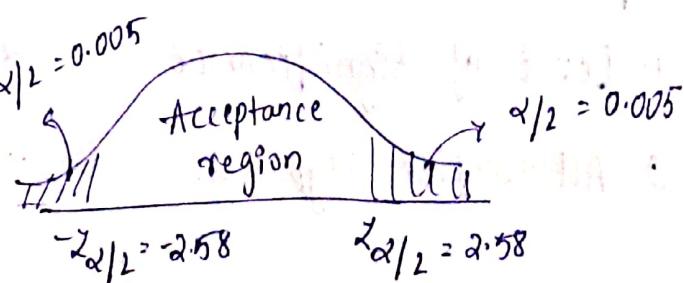
|                   | 1%                    | 5%                    | 10%                   |
|-------------------|-----------------------|-----------------------|-----------------------|
| Two Tailed test   | $Z_{\alpha/2} = 2.58$ | $Z_{\alpha/2} = 1.96$ | $Z_{\alpha/2} = 1.65$ |
| Right Tailed test | $Z_{\alpha} = 2.33$   | $Z_{\alpha} = 1.65$   | $Z_{\alpha} = 1.29$   |
| Left Tailed test  | $Z_{\alpha} = -2.33$  | $Z_{\alpha} = -1.65$  | $Z_{\alpha} = -1.29$  |

. Two tailed 1%:

$$\alpha = 1\% = 0.01$$

$$\frac{\alpha}{2} = \frac{0.01}{2} = 0.005$$

$$\alpha/2 = 0.005$$

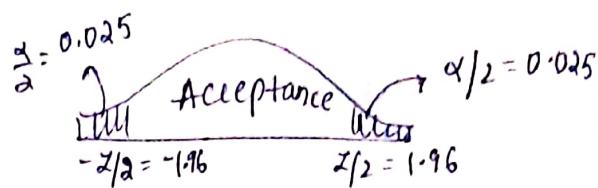


$$1\% \text{ two tailed } |Z_{\alpha/2}| = 2.58$$

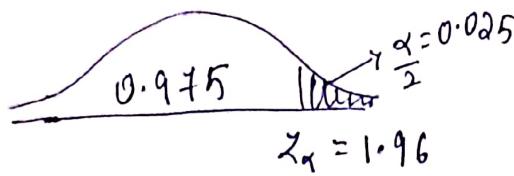
• 5.1. Two Tailed:

$$\alpha = 0.05 = 5\%$$

$$\frac{\alpha}{2} = \frac{0.05}{2} = 0.025$$



Total area under  
Normal curve = 1



From  $-z_\alpha$  upto  $z_\alpha$  area  
is  

$$1.000 - \frac{0.025}{0.975}$$

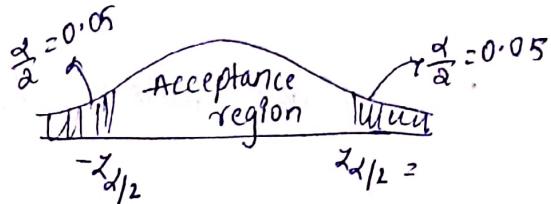
So two tailed test 5.1. los,

$$|z_{\alpha/2}| = 1.96$$

• 10.1. Two Tailed:

$$\alpha = 10\% = \frac{10}{100} = 0.1$$

$$\frac{\alpha}{2} = \frac{0.1}{2} = 0.05$$



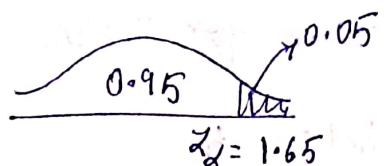
Nrm1 ttable

1.6 - horizontal

$\frac{0.05}{1.65}$  - vertical

10.1. two tailed

$$|z_{\alpha/2}| = 1.65$$



$$1.00 - \frac{0.05}{0.95}$$

• 5.1. Right tailed:

$$\alpha = 5\% = 0.05$$



$$1.00 - \frac{0.05}{0.95}$$

$$z_\alpha = 1.65$$

5.1. right tailed  $z_\alpha = 1.65$

1.1. right tailed:

$$\alpha = 1\% = 0.01$$



$$Z_\alpha = 2.33$$

$$\begin{array}{r}
 1.00 \\
 -0.01 \\
 \hline
 0.99
 \end{array}$$

1.1. right tailed,  $Z_\alpha = 2.33$

10.1. right tailed:

$$\alpha = \frac{10}{100} = 0.1$$



$$Z_\alpha = 1.29$$

$$\begin{array}{r}
 1.0 \\
 -0.1 \\
 \hline
 0.9
 \end{array}
 \quad \begin{array}{l}
 1.2 \text{ hori} \\
 \hline
 0.09 \text{ vert} \\
 \hline
 1.29
 \end{array}$$

10.1. right tailed,  $Z_\alpha = 1.29$

Procedure for Testing of Hypothesis:

1. Setup  $H_0$
2. Setup  $H_1$
3. Choose the level of significance,  $\alpha$
4. Calculate the test statistic under  $H_0$ .
5. Compare the calculated value of the test statistic with table value.

For two tailed and right tailed test

calculated value < Table value , Accept  $H_0$

Calculated value  $>$  Table value , Reject  $H_0$

Left Tailed test

calculated value < Table value , Reject  $H_0$

calculated value  $>$  Table value , Accept  $H_0$