

III B.Sc. Statistics

Subject name: Design of Experiment

Subject code :CST62

Unit :2

Analysis of variance [ANOVA]

Introduction:

Analysis of variance is a powerful statistical tool for test of significance. The basic purpose of the analysis of variance is to test the homogeneity of several means.

Definition:

Analysis of variance is the separation of variance ascribable to one group of causes from the variance ascribable to other group.

(X) Assumption:

- 1) The observations are independent.
- 2) Parent population from which observations are taken is normal.
- 3) Various treatments and environment effective are additive in nature.

Applications:

ANOVA techniques is now frequently apply in testing the linearity of the fitted regression line or the significance of the correlation ratio r .

Types of ANOVA:

- 1) One-way classification
- 2) Two-way classification

Cochran's theorem: (statement)

Let x_1, x_2, \dots, x_n denote a random sample from normal population $N(0, \sigma^2)$

Let the sum of the squares of these values be written in the form,

$$\sum_{i=1}^n x_i^2 = Q_1 + Q_2 + \dots + Q_k$$

where Q_j is a quadratic form in x_1, x_2, \dots, x_n with rank r_j , ($j=1, 2, \dots, k$) then the random variables Q_1, Q_2, \dots, Q_k are mutually independent and Q_j/σ^2 is chi-square variate with r_j degrees of freedom if and only if $\sum_{j=1}^k r_j = n$

One-way classification:

Let us suppose that N observations x_{ij} ($i=1, 2, \dots, K$) ($j=1, 2, \dots, n_i$) of a random variable X are grouped on some basis into K classes of sizes n_1, n_2, \dots, n_K respectively

				Mean	Total
x_{11}	x_{12}	\dots	x_{1n_1}	\bar{x}_1	T_1
x_{21}	x_{22}	\dots	x_{2n_2}	\bar{x}_2	T_2
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_{i1}	x_{i2}	\dots	x_{in_i}	\bar{x}_i	T_i
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_{k1}	x_{k2}	\dots	x_{kn_K}	\bar{x}_K	T_K
					G

The total variation in the observation x_{ij} can be split into the following two components

1) The variation between the classes or the variation due to different basis of classification commonly known

2) The variation within the classes inherent variation of the random variable within the observations of a class.

The first type of variation is due to assignable causes which can be detected and controlled by human endeavour.

The second type of variation is due to chance causes which are beyond the control of human hand.

2/2 Null hypothesis (H_0)

To test the equality of the population means i.e. all mean values are homogeneity

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k = \mu$$

Alternative hypothesis (H_1)

All mean values are not homogeneity (Heterogeneity)

$$H_1: \mu_1 \neq \mu_2 \neq \dots \neq \mu_k \neq \mu$$

Mathematical model:

The linear mathematical model is

$$x_{ij} = \mu_i + \epsilon_{ij}$$

$$\mu_i - \mu = \alpha_i$$

$$= \mu_i + \epsilon_{ij} + \mu - \mu$$

$$= \mu + (\mu_i - \mu) + \epsilon_{ij}$$

$$= \mu + \alpha_i + \varepsilon_{ij}$$

Where,

$x_{ij} \Rightarrow$ represent i^{th} row and j^{th} column

$\mu \Rightarrow$ general mean effect

$\alpha_i \Rightarrow$ effect of i^{th} row

$\varepsilon_{ij} \Rightarrow$ error effect due to chance

Statistical analysis

We consider

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{..})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.} + \bar{x}_{i.} - \bar{x}_{..})^2$$

$$= \sum_{i=1}^k \sum_{j=1}^{n_i} \left[(x_{ij} - \bar{x}_{i.})^2 + (\bar{x}_{i.} - \bar{x}_{..})^2 \right.$$

$$\left. + 2(x_{ij} - \bar{x}_{i.})(\bar{x}_{i.} - \bar{x}_{..}) \right]$$

$$= \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.})^2 + \sum_{i=1}^k \sum_{j=1}^{n_i} (\bar{x}_{i.} - \bar{x}_{..})^2$$

$$+ 2 \sum_{i=1}^k \sum_{j=1}^{n_i} [(x_{ij} - \bar{x}_{i.})(\bar{x}_{i.} - \bar{x}_{..})]$$

Since the product term is vanished

$$2 \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.})(\bar{x}_{i.} - \bar{x}_{..}) = 0$$

The algebraic sum of deviation of the i^{th} row from their mean is zero.

Then we have

$$\sum_{i=1}^K \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{..}) = \sum_{i=1}^K \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.})^2 + \sum_{i=1}^K \sum_{j=1}^{n_i} (\bar{x}_{i.} - \bar{x}_{..})^2$$

where

$$\sum_{i=1}^K \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{..})^2 \Rightarrow \text{total sum of squares.}$$

$$= TSS \text{ or } ST^2$$

$$\sum_{i=1}^K \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.})^2 \Rightarrow \text{error sum of square}$$

$$= SSE \text{ or } SE^2$$

$$\sum_{i=1}^K n_i (\bar{x}_{i.} - \bar{x}_{..})^2 = \text{column sum of squares (or)}$$

sum of square due to treatment

$$= SSC \text{ (or)} SC^2$$

$$TSS / \text{Total sum of square} = SSC + SSE$$

$$ST^2 = SC^2 + SE^2$$

Degrees of freedom:

Degrees of freedom for
column $c-1$

Degrees of freedom for
error $N-c$.

Degrees of freedom for total
 $N-1$

Mean sum of square

The sum of squares divided by
its degrees of freedom is known
as mean sum of squares.

1) Mean sum of squares for column

$$SC^2 = \frac{SSC}{c-1} = \frac{SC^2}{c-1}$$

2) Mean sum of squares for error

$$SE^2 = \frac{SSE}{N-c} = \frac{SE^2}{N-c}$$

ANOVA Table

Source of Variation	Sum of Squares	Degrees of freedom	Mean sum of square	Variance Ratio
Between column	$SSC (or)$ SC^2 $\sum n(\bar{x}_i - \bar{x})^2$	$r - C - 1$	$MSSc (or) = \frac{SC^2}{SC^2}$ $C - 1$	$F_c = \frac{SC^2}{SE^2}$
Error	$SSE (or)$ SE^2 $\sum (x_{ij} - \bar{x}_i)^2$	$N - C - r$	$MSSE (or) = \frac{SE^2}{N - C}$ SE^2	
Total	$TSS (or)$ ST^2 $\sum \sum (x_{ij} - \bar{x})^2$	$N - 1 - r$		

If an observed value of $F >$ the tabulated value of F_α at specified level of significance that is 5% or 1% the H_0 is rejected.

Problem

- 1) Random samples of same size drawn from air conditioners from three manufactures showed the following life expectation (month)

Manufacturer I 34 28 32 29 31

Manufacturer II 23 34 37 32

Manufacturer III 30 39 35 37 32 31

use the ANOVA at 5% L.O.S and indicate whether the average life expectations of 3 different brands of room air conditioner differ significantly

calculation:

Null hypothesis (H_0):

The average life expectation of three brands of room air conditioners are same

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

Alternative hypothesis (H_1):

The average life expectation of three brands of room air conditioners are not same.

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

Level of significance

$$\alpha = 5\%$$

$$\alpha = 0.05$$

Test statistic:

X_1	X_2	X_3	X_1^2	X_2^2	X_3^2
34	23	30	1156	529	900
28	34	39	784	1156	1521
32	37	35	1024	1369	1225
29	32	37	841	1024	1369
31		32	961		1024
		31			961
$\Sigma X_1 =$	$\Sigma X_2 =$	$\Sigma X_3 =$	$\Sigma X_1^2 =$	$\Sigma X_2^2 =$	$\Sigma X_3^2 =$
154	126	204	4766	4078	7000

N = No. of observation

ΣT = Total No. of observation

Grand total

correction factor:

$$\begin{aligned} C.F. &= \frac{[\Sigma T]^2}{N} \\ &= \frac{[1184]^2}{15} \\ &= \frac{234256}{15} \end{aligned}$$

$$\boxed{C.F. = 15617}$$

Total sum of square:

$$\begin{aligned} TSS &= \Sigma X_1^2 + \Sigma X_2^2 + \Sigma X_3^2 - C.F. \\ &= 4766 + 4078 + 7000 - 15617 \\ &= 15844 - 15617 \end{aligned}$$

$$TSS = 227$$

column sum of square:

$$SSC = \frac{(\Sigma X_1)^2}{N_1} + \frac{(\Sigma X_2)^2}{N_2} + \frac{(\Sigma X_3)^2}{N_3} - C.F.$$

$$= \frac{(154)^2}{5} + \frac{(126)^2}{4} + \frac{(204)^2}{6} - 15617$$

$$= \frac{23716}{5} + \frac{15876}{4} + \frac{41616}{6} - 15617$$

$$= 4743.2 + 3969 + 6936 - 15617$$

$$= 15648.2 - 15617$$

$$= 31.2$$

Error sum of square

$$SSE = SST - SSC$$

$$= 227 - 31.2$$

$$SSE = 195.8$$

ANOVA table

Source of Variation	Sum of Squares	Degrees of freedom	Mean sum of Square	Variance ratio
Between column	$SSC = 31.2$	$C-1 = 3-1 = 2$	$MSC = \frac{SSC}{C-1} = \frac{31.2}{2} = 15.6$	$F_c = \frac{MSE}{MSC} = \frac{16.3}{15.6} = 1.04$
Error	$SSE = 195.8$	$N-C = 15-3 = 12$	$MSE = \frac{SSE}{N-C} = \frac{195.8}{12} = 16.316$	
Total	$TSS = 227$	$N-1 = 15-1 = 14$		

Degrees of freedom:

$$f = (N-C, C-1)$$

$$= (15-3, 3-1)$$

$$= (12, 2)$$

Table value:

At 5% level of significance for (12, 2) degrees of freedom the table value is $F_{\alpha} = 19.16$

Conclusion :

Calculated value $<$ table value

$$1.045 < 19.41$$

$$|F| < F_{\alpha}$$

H_0 is accepted

Two-way classification

Let us consider the case when there are two factors which may affect the variate values x_{ij}

Example:

The yield of milk may be affected by differences in treatments that is locations as well as the differences in variety.

Let us now suppose that N cows are divided into h different groups, each group containing k cows.

Let us consider the effect of k treatments on the yield of milk.

	α_{11}	α_{12}	\dots	α_{1j}	\dots	α_{1h}	Mean $\bar{\alpha}_{1.}$
	α_{21}	α_{22}	\dots	α_{2j}	\dots	α_{2h}	$\bar{\alpha}_{2.}$
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
	α_{i1}	α_{i2}	\dots	α_{ij}	\dots	α_{ih}	$\bar{\alpha}_{i.}$
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
	α_{k1}	α_{k2}	\dots	α_{kj}	\dots	α_{kh}	$\bar{\alpha}_{k.}$
Mean	$\bar{\alpha}_{.1}$	$\bar{\alpha}_{.2}$	\dots	$\bar{\alpha}_{.j}$	\dots	$\bar{\alpha}_{.h}$	$\bar{\alpha}_{..}$
Total	$T_{.1}$	$T_{.2}$	\dots	$T_{.j}$	\dots	$T_{.h}$	

Null hypothesis (H_0)

The treatment as well as varieties are homogeneous

$$H_0(t): \mu_{1.} = \mu_{2.} = \dots = \mu_{k.} = \mu$$

$$H_0(v): \mu_{.1} = \mu_{.2} = \dots = \mu_{.h} = \mu$$

these equivalent

$$H_t: \alpha_1 = \alpha_2 = \dots = \alpha_k = 0$$

$$H_v: \beta_1 = \beta_2 = \dots = \beta_h = 0$$

Mathematical model:

Let us suppose that α_{ij} ($i = 1, 2, \dots, k$, $j = 1, 2, \dots, h$) are independent the mathematical model

$$\alpha_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}$$

- $x_{ij} \Rightarrow$ Effect of i th row and j th column
 $\mu \Rightarrow$ General mean effect
 $\alpha_i \Rightarrow$ effect due to i th row
 $\beta_j \Rightarrow$ effect due to j th column
 $\epsilon_{ij} \Rightarrow$ Error effect due to chance

Statistical Analysis:

We can consider

$$\begin{aligned}
 \sum_{i=1}^k \sum_{j=1}^h (x_{ij} - \bar{x}_{..})^2 &= \sum \sum [x_{ij} - \bar{x}_{i.} + \bar{x}_{i.} - \bar{x}_{.j} \\
 &\quad + \bar{x}_{.j} - \bar{x}_{..} + \bar{x}_{..} - \bar{x}_{..}]^2 \\
 &= \sum \sum [(\bar{x}_{i.} - \bar{x}_{..}) + (\bar{x}_{.j} - \bar{x}_{..})]^2 \\
 &= \sum \sum (\bar{x}_{i.} - \bar{x}_{..})^2 + \sum \sum (\bar{x}_{.j} - \bar{x}_{..})^2 \\
 &\quad + 2 \sum \sum (\bar{x}_{i.} - \bar{x}_{..})(\bar{x}_{.j} - \bar{x}_{..}) \\
 &= \text{SSR} + \text{SSC} + 2 \sum \sum (\bar{x}_{i.} - \bar{x}_{..})(\bar{x}_{.j} - \bar{x}_{..})
 \end{aligned}$$

Since algebraic sum of deviation set of observation about their mean is zero that is all the product term is vanished. So, we have

$$\sum_{i=1}^k \sum_{j=1}^h (\alpha_{ij} - \bar{\alpha}_{i.})^2 = \sum_{i=1}^k \sum_{j=1}^h (\alpha_{ij} - \bar{\alpha}_{i.} - \bar{\alpha}_{.j} + \bar{\alpha}_{..})^2$$

$$+ \bar{\alpha}_{..})^2 + h \sum_{i=1}^k (\bar{\alpha}_{i.} - \bar{\alpha}_{..})^2$$

$$+ k \sum_{j=1}^h (\bar{\alpha}_{.j} - \bar{\alpha}_{..})^2$$

$$TSS = SSR + SSC + SSE$$

$$ST^2 = SR^2 + SC^2 + SE^2$$

Where,

* $\sum_{i=1}^k \sum_{j=1}^h (\alpha_{ij} - \bar{\alpha}_{i.})^2$ is the total sum of squares

= TSS or ST^2

* $h \sum_{i=1}^k (\bar{\alpha}_{i.} - \bar{\alpha}_{..})^2$ is row sum of square

= SSR or SR^2

* $k \sum_{j=1}^h (\bar{\alpha}_{.j} - \bar{\alpha}_{..})^2$ is column sum of square

= SSC or SC^2

* $\sum_{i=1}^k \sum_{j=1}^h (\alpha_{ij} - \bar{\alpha}_{i.} - \bar{\alpha}_{.j} + \bar{\alpha}_{..})^2$ is error

sum of square

= SSE or SE^2

Degrees of freedom

$$\text{For row } \gamma = r - 1$$

$$\text{For column } \gamma = c - 1$$

$$\text{For error } \gamma = (r - 1)(c - 1)$$

$$\text{For total } \gamma = N - 1 \text{ (or) } rc - 1$$

Mean sum of square

The sum of square divided by its degrees of freedom is known as mean sum of square.

* Mean sum of square for row

$$SR^2 = \frac{SSR}{r - 1} \text{ (or) } \frac{SR^2}{r - 1}$$

* Mean sum of square for column

$$SC^2 = \frac{SSC}{c - 1} \text{ (or) } \frac{SC^2}{c - 1}$$

* Mean sum of square for error

$$SE^2 = \frac{SSE}{(r - 1)(c - 1)} \text{ (or) } \frac{SE^2}{(r - 1)(c - 1)}$$

Test Statistic

For row

$$FR = \frac{SR^2}{SE^2}$$

For column

$$FC = \frac{SC^2}{SE^2}$$

ANOVA Table

Source of variation	Sum of squares	Degrees of freedom	Mean sum of squares	Variable ratio
Between row	$SSR = h \sum (\bar{x}_{i.} - \bar{x}_{..})^2$	$(r-1)$	$SR^2 = \frac{SR^2}{r-1}$	$F_R = \frac{MSSR}{MSSE}$
Between column	$SSC = K \sum (\bar{x}_{.j} - \bar{x}_{..})^2$	$(c-1)$	$SC^2 = \frac{SC^2}{c-1}$	
Error	$SSE = \sum \sum (x_{ij} - \bar{x}_{i.} - \bar{x}_{.j} + \bar{x}_{..})^2$	$(r-1)(c-1)$	$SE^2 = \frac{SE^2}{(r-1)(c-1)}$	$F_C = \frac{MSSC}{MSSE}$
Total	$TSS = \sum_{i=1}^K \sum_{j=1}^h (x_{ij} - \bar{x}_{..})^2$	$N-1$		

If an observed value of $F >$ the tabulated value of F_α at specified level of significance that is 5%. or 1%. the H_0 is rejected.

Multiple Range Test [MRT]

1) Student - Newman Keul's test:

The Newman Keul's method is a step-wise multiple comparison. This procedure used to identify sample means that are significantly different from each other. The procedure for Newman Keul's test are

i) The employers step-wise approach when comparing sample means.

ii) Prior to any mean comparison, all sample means are rank ordered in ascending or descending order, thereby producing an ordered range (P) of sample means.

iii) A comparison is made between the largest and smallest sample means within the largest range.

iv) Assuming that the largest range is H means ($P=H$). A significant difference between the largest and smallest mean as revealed by the "Newman Keul's" method could result in deflection of the null hypothesis for that specific range of means.

v) The next largest comparison of 2 sample means would be made within a smaller range 3 means ($p=3$).

vi) continue this process until a final comparison is made.

vii) If there is no significant difference between the 2 sample means, all the null hypothesis within that range would be retained no further comparison within smaller ranges are necessary.

18/2 Duncan's Multiple Range Test [DMRT]

Duncan's Multiple Range test is a post hoc test to measure specific differences between pairs of means.

This test is commonly used in agronomy and other agricultural research. The result of the test is a set of subsets of means, where in each subset means have been found not to be significantly different from one other.

This test is summarized the way in finding several significant differences with increasing values (descending order) which depending on the extent of the distance between the treatment means after arranged.

Treatment	Main yield (kg)	Rank
T ₂	2678	1
T ₃	2552	2
T ₄	2128	3
T ₁	2127	4
T ₅	1796	5
T ₆	1681	6
T ₇	1316	7

$$SE = \sqrt{\frac{2S^2}{r}}$$

$$= 217.68$$

(t-1) value of the smallest significant range are computed

$$R_p = \frac{(r_p)(SE)}{\sqrt{2}} \quad \text{for } p = 2, 3, \dots$$

from the largest mean subtract

R_p for largest p , then declare as significantly different from the largest mean.

For the remaining treatment whose values are larger, then the difference

R_p . continue this process till all the treatment of T .

Present the results by using the (T_1, T_2, \dots, T_n) to indicate which treatment pairs are significantly different from each other.

Tukey's Range Test [TRT]

Tukey's honest significance difference test can be used on raw data to find means that are significantly different from each other.

This test compares the means of every treatment to the means of every other treatment.

This test is based on a formula very similar to that of the t-test. Tukey's test is essentially a t-test except that it corrects for experiment wise error rate.

$$q_s = \frac{Y_A - Y_B}{SE}$$

where,

$Y_A \rightarrow$ is a larger of the two means being compared.

$Y_B \rightarrow$ is the smaller of the two means being compared.

$SE \rightarrow$ Standard error.

This q_s value can be compared to a q value from the range distribution.

If the q_c value is larger than the q critical value, the two means are said to be significantly different.

$$q = \frac{(\bar{Y}_{\max} - \bar{Y}_{\min})}{\sqrt{\frac{s^2}{n}}}$$

$$\sqrt{\frac{s^2}{n}}$$