

Unit - 7

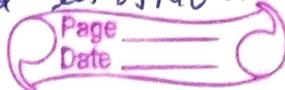
Analysis of Variance :- The analysis of variance is a powerful statistical tool for tests of significance. The test of significance based on t -distⁿ is an adequate procedure only for testing the significance of the difference betⁿ two sample means. In a situation when we have three or more samples to consider at a time an alternative procedure is needed for testing the hypothesis that all the samples are drawn from the same popⁿ. i.e. they have the same mean.

eg :- Five fertilizers are applied to four plots each of wheat and yields of wheat on each of the plot is given. Now we are interested in finding out whether the effect of these fertilizers on the yield is significantly different or in other words, whether the samples have come from the same normal popⁿ.

The answer to this problem is provided by the technique of analysis of variance. Thus basic purpose of the analysis of variance is to test the homogeneity of several means.

The term "Analysis of variance" was introduced by prof. R. A. Fisher in 1930's to deal with problem in the analysis of agronomical data.

Cochran's Thm: used to justify results relating to the prob. distⁿ of statistics that are used in ANOVA.



Variation is inherent in nature. The total variation in any set of numerical data is due to a no. of causes which may be classified as

- Assignable causes, and
- Chance causes.

The variation due to assignable causes can be detected and measured whereas the variation due to chance causes is beyond the control of human hand and cannot be traced separately.

Definition:

According to Prof. R.A. Fisher, Analysis of Variance (ANOVA) is the "separation of variance ~~for~~ ^{ascrivable} to one group of causes from the variance ~~for~~ ^{ascrivable} to other group."

The ANOVA consists in the estimation of the amount of variation due to each of the indept. factors (causes) separately and then comparing this estimates due to assignable factors (causes) with the estimates due to chance factor (causes).

Assumption :- For the validity of the F-test in ANOVA, the following assumption are made:-

- i) The obserⁿ are indept.
- ii) Parent popⁿ from which observation are taken ~~as~~ is normal and
- iii) Various treatment and environmental effects are additive in nature.

Cochran's Theorem: Let x_1, x_2, \dots, x_n denote a s.s. from normal popⁿ $N(\mu, \sigma^2)$. Let the sum of the squares of these values are

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

Where φ_j is a quadratic form of X_1, X_2, \dots, X_n with rank $(d.f.) r_j$; $j=1, 2, \dots, k$: Then the ~~l.v. $\varphi_1, \varphi_2, \dots, \varphi_k$~~ ^{Page Data} are mutually indept. and φ_j / σ^2 is a χ^2 -variante with r_j degrees of freedom iff $\sum_{j=1}^k r_j = n$.

ANOVA for One-Way Classification :- Let

(With one obsrⁿ per cell) be

Suppose that N obsrⁿ 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. ($N = \sum_{i=1}^k n_i$) as exhibited below

		Means	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
x_{k1}	$x_{k2} \dots x_{kn_k}$	\bar{x}_k	T_k
			G

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

i) The variation betwⁿ the classes or the variation due to different bases of classification commonly known as treatments.

ii) The variations within the classes i.e. the inherent variation of the l.v. within the obsrⁿ of a class.

The first type of variation is due to assignable causes which can be detected and controlled by human endeavour and the second type of variation is due to chance causes which are beyond the control of human hand.

The main object of analysis of variance technique is to ~~explore~~ estimate if there is significant difference between the class means in view of the due to chance causes i.e. variability within the ^{separate} classes.

e.g. In particular let us consider the effect of k different rations on the yield of milk of N cows (of same breed & stock). These N cows are divided into k classes of sizes n_1, n_2, \dots, n_k respectively.

$$N = \sum_{i=1}^k n_i$$

Here the sources of variation are

- (i) Effect of the rations (treatment) t_i $i=1, 2, k$
- (ii) error ϵ produced by numerous causes.

They are not detected and identified & produce a variation of random nature.

Mathematical Model :- Let us consider the linear model will be

$$\begin{aligned} x_{ij} &= \mu + \alpha_i + \epsilon_{ij} \\ &= \mu + (\mu_i - \mu) + \epsilon_{ij} \\ &= \mu + \alpha_i + \epsilon_{ij} \quad ; \text{ where } (i=1, 2, \dots, k) \\ &\qquad\qquad\qquad (j=1, 2, \dots, n_i) \end{aligned} \tag{1}$$

- (i) x_{ij} is the yield from the j th cow, ($j=1, 2, \dots, n_i$) fed on the i th rations ($i=1, 2, \dots, k$)
- (ii) μ is the general mean effect given by

$$\bar{u} = \frac{\sum_{i=1}^k (n_i u_i)}{N} \quad - (2)$$

where u_i is the fixed effect due to the i^{th} ration, i.e. if there were no treatment differences and no chance causes then the yield of each cow will be \bar{u} .

(iii) α_i is the effect of the i^{th} ration (treatment) is given by

$$\alpha_i = u - \bar{u}_i \quad (i=1, 2, \dots, k) \quad - (3)$$

i.e. the i^{th} treatment (ration) increase or decrease the yield by an amount α_i .
using (2), we get

$$\begin{aligned} \sum_{i=1}^k n_i \alpha_i &= \sum_i n_i (\bar{u}_i - \bar{u}) = \sum_i n_i \bar{u}_i - N \bar{u} \\ &= N \bar{u} - N \bar{u} = 0 \end{aligned}$$

(iv) ε_{ij} is the error effect due to chance.

Assumption in the model

- (i) All the obsrⁿ x_{ij} are indept.
- (ii) Different effects are additive in nature
- (iii) ε_{ij} are i.i.d. $N(0, \sigma_e^2)$

under the (iii) assumption, the model (1) becomes

$$E(x_{ij}) = u_i + \varepsilon_{ij} = \bar{u} + \alpha_i + \varepsilon_{ij} ; \quad \begin{cases} i=1, 2, \dots, k \\ j=1, 2, \dots, n_i \end{cases}$$

Null hypothesis of one for ANOVA for one-way classification.

Here we want to test the equality of popⁿ
means i.e. the homogeneity of different
relations. Hence null hypothesis is given by

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

which from ③ reduces to

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_K = 0$$

Statistical Analysis of the Model :- Let us write

$\bar{x}_{i\cdot}$ = mean of the i^{th} class.

$$= \frac{\sum_{j=1}^{n_i} x_{ij}}{n_i}; \quad (i=1, 2, \dots, k)$$

$$\begin{aligned} \bar{x}_{..} &= \text{overall mean} = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{n_i} x_{ij} \\ &= \frac{1}{N} \sum_{i=1}^k n_i \bar{x}_{i\cdot} \end{aligned}$$

The parameters μ and α_i in ① are estimated by the principle of least squares on minimising the error (residual) sum of squares given by

$$E = \sum_{i=1}^k \sum_{j=1}^{n_i} \varepsilon_{ij}^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \mu - \alpha_i)^2$$

The normal eqⁿ for estimating μ and α_i are

$$\frac{\partial E}{\partial \mu} = -2 \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \mu - \alpha_i) = 0 \quad (*)$$

and $\frac{\partial E}{\partial \alpha_i} = -2 \sum_{j=1}^{n_i} (x_{ij} - \mu - \alpha_i) = 0 \quad (**)$

from $\textcircled{*}$, we get

$$\sum_{i,j} x_{ij} - N\bar{u} - \sum n_i \bar{x}_i = 0$$

$$\hat{u} = \frac{1}{N} \sum_{i,j} x_{ij} = \bar{x}_{..}$$

$[\because \sum n_i \bar{x}_i = 0]$

& from $\textcircled{**}$, we get

$$\sum_j x_{ij} - n_i \hat{u} + n_i \hat{\alpha} = 0$$

$$\Rightarrow \hat{\alpha} = \frac{1}{n_i} \sum_j x_{ij} - \hat{u} = \bar{x}_{i..} - \hat{u}$$

$$\hat{\alpha} = \bar{x}_{i..} - \bar{x}_{..}$$

Hence substituting the value of u & α in $\textcircled{1}$
the model becomes

$$x_{ij} = \bar{x}_{..} + (\bar{x}_{i..} - \bar{x}_{..}) + (x_{ij} - \bar{x}_{i..})$$

We introduce the error component e_{ij} so
that both the sides are equal. This is the
deviation within the class which is due to
randomisation. Transposing $\bar{x}_{..}$ to the left
squaring both sides and summing over i, j
we get

$$\sum_{i,j} [x_{ij} - \bar{x}_{i..} + \bar{x}_{i..} - \bar{x}_{..}]^2$$

$$\begin{aligned} \sum_{i,j} (x_{ij} - \bar{x}_{..})^2 &= \sum_{i,j} \left[(\bar{x}_{i..} - \bar{x}_{..}) + (x_{ij} - \bar{x}_{i..}) \right]^2 \\ &= \sum_{i,j} (\bar{x}_{i..} - \bar{x}_{..})^2 + \sum_{i,j} (x_{ij} - \bar{x}_{i..})^2 \\ &\quad + 2 \sum_i \bar{x}_{i..} \sum_j (x_{ij} - \bar{x}_{i..}) \end{aligned}$$

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

But $\sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.}) = 0$ [algebraic sum of the deviation of the ratios from their mean is zero]

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

$S_T^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{..})^2$ is known as total sum of squares (T.S.S.)

$S_E^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_{i.})^2$ is called within sum of squares or error sum of squares (S.S.E.) and

$S_T^2 = \sum_i n_i (\bar{x}_{i.} - \bar{x}_{..})^2$ is called S.S. due to treatments (S.S.T.)

Then $T.S.S. = S.S.E. + S.S.T.$

Degrees of freedom for various S.S. :-

S_T^2 , the total S.S. will carry $(N-1)$ d.f. (one d.f. lost because of the linear constraint)

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

similarly the treatment sum of squares

S_t^2 will have $(k-1)$ d.f.

& S_E^2 the error s.s. will have $(N-k)$ d.f. .

Hence we see that the

$$N-1 = (N-k) + (k-1)$$

Mean Sum of Squares (M.S.S.)

~~Defn~~ ~~SS~~

The sum of squares divided by its degrees of freedom gives the corresponding variance or the mean sum of squares (M.S.S.). Thus

$$\frac{S_t^2}{(k-1)} = \frac{S.S.T.}{(k-1)} - S_E^2 \text{ is M.S.S. due to treatment}$$

$$\frac{S_E^2}{N-k} = \frac{S.S.E.}{(N-k)} = S_E^2 \text{ is M.S.S. due to error.}$$

Hence the test statistic for H_0 provided by the variance ratio

$$F = \frac{S_t^2}{S_E^2}$$

XXX

If H_0 is true the F - should take the value 1, otherwise it should be greater than unity. In order to find out if an observed value of F is significantly greater than unity. we have

to obtain the sampling distⁿ of the F-statistic defined as

$$F = \frac{S_t^2}{S_E^2}$$

Under H_0 by Cochran's thm., $\frac{S_t^2}{\sigma_e^2} \sim \chi^2_{k-1}$ and $\frac{S_E^2}{\sigma_e^2} \sim \chi^2_{N-k}$ are independently distributed as χ^2 variates with $(k-1)$ and $(N-k)$ d.f. respectively.
Hence the statistic

$$F = \left[\frac{\frac{S_t^2}{\sigma_e^2} \cdot \frac{1}{k-1}}{\frac{S_E^2}{\sigma_e^2} \cdot \frac{1}{N-k}} \right] = \frac{S_t^2}{S_E^2}$$

follows Snedecor's F (central) dist^m with $[(k-1), (N-k)]$ d.f.

Thus if an observed value of F obtained from xxx is greater than tabulated value of F for $(k-1, N-k)$ d.f. at specified level of significance (usually 5% or 1%) then H_0 is refuted at that level otherwise H_0 may be retained.

The above statistical analysis is very elegantly presented in the following table

ANOVA Table for ONE-WAY CLASSIFIED DATA

Sources of Variation	Sum of Squares	d.f.	Mean sum of squares	Variance Ratio
Treatment (actions)	S_t^2	$k-1$	$S_t^2 = \frac{S_t^2}{k-1}$	$\frac{S_t^2}{S_E^2} = F_{k-1, N-k}$
Error	S_E^2	$N-k$	$S_E^2 = \frac{S_E^2}{N-k}$	
Total	S_T^2	$N-1$		

Remark: For practical point of view; various sum of squares reduced to a great extent by using following simplified formulae

$$\begin{aligned}\text{Total S.S.} &= \sum_{i,j} \sum (x_{ij} - \bar{x}_{..})^2 \\ &= \sum \sum x_{ij}^2 - \frac{(\sum \sum x_{ij})^2}{N} \\ &= \sum \sum x_{ij}^2 - \frac{G^2}{N}\end{aligned}$$

where G is the grand total of all the observations and $N = n_1 + n_2 + \dots + n_k = \sum n_i$

The expression $\sum \sum x_{ij}^2$ the sum of squares of all the observations is known as raw sum of squares (R.S.S.) and the expression $\frac{G^2}{N}$ is called the correction factor (C.F.).

$$\text{Total S.S.} = \text{R.S.S.} - \text{C.F.}$$

$$\begin{aligned}\text{Error S.S.} &= \sum \sum (x_{ij} - \bar{x}_{i.})^2 = \sum_i \left[\sum_j (x_{ij} - \bar{x}_{i.})^2 \right] \\ &= \sum_i \left[\sum_j x_{ij}^2 - \frac{(\sum x_{ij})^2}{n_i} \right] = \sum_i \left[\sum_j x_{ij}^2 - \frac{T_{i.}^2}{n_i} \right] \\ S_E^2 &= \sum \sum x_{ij}^2 - \sum_i \left(\frac{T_{i.}^2}{n_i} \right)\end{aligned}$$

where $T_{i.}$ is the total yield from the units receiving the i th treatment

$$\begin{aligned}\text{Treatment S.S.} &= S_T^2 = \text{T.S.S.} - \text{S.S.E.} \\ &= \sum \left(\frac{T_{i.}^2}{n_i} \right) - \frac{G^2}{N} \\ &= \sum \left(\frac{T_{i.}^2}{n_i} \right) - \text{C.F.}\end{aligned}$$

ANOVA for Two-way Classification (with one observation per cell)

Let us consider the case when there are two factors which may affect the yield of milk.

e.g. the yield of milk may be affected by difference of treatments i.e. rations as well as the difference in variety i.e. breed and stock of the cows. Let us now suppose that the N cows are divided into h different groups or classes according to their breed and stock, each group containing k cows and then let us consider the effect of K treatments (i.e. rations given at random to cow in each group) on the yield of milk.

Let us suffix i refer to the treatments (rations) and suffix j refer to the varieties (breed of the cow). Then the yield of milk x_{ij} ($i=1, 2, \dots, K$, $j=1, 2, \dots, h$) of $N = h \times k$ cows furnish the data for the comparison of the treatments. The yield may be expressed as variate values in the following $K \times h$ two-way table

				mean	total		
x_{11}	x_{12}	\dots	x_{1j}	\dots	x_{1h}	$\bar{x}_{1\cdot}$	T_1
x_{21}	x_{22}	\dots	x_{2j}	\dots	x_{2h}	$\bar{x}_{2\cdot}$	T_2
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_{i1}	x_{i2}	\dots	x_{ij}	\dots	x_{ih}	$\bar{x}_{i\cdot}$	T_i
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_{K1}	x_{K2}	\dots	x_{Kj}	\dots	x_{Kh}	$\bar{x}_{K\cdot}$	T_K
Mean	$\bar{x}_{1\cdot}$	$\bar{x}_{2\cdot}$	$\bar{x}_{i\cdot}$	$\bar{x}_{K\cdot}$	Total		

total | T.1 T.2 T.j T.h \rightarrow 61

Mathematical Model :- let x_{ij} be the yield from the cow of j th variety fed on the i th ration ($i=1, 2, \dots, k$, $j=1, 2, \dots, h$).

Let us suppose that x_{ij} ($i=1, 2, \dots, k$, $j=1, 2, \dots, h$) are indept., normally distributed as $N(\mu_{ij}, \sigma^2)$.
The linear mathematical model becomes

$$E(x_{ij}) = \mu_{ij}$$

$$\text{or } x_{ij} = \mu_{ij} + \varepsilon_{ij} \quad \text{--- (1)}$$

where ε_{ij} are i.i.d. $N(0, \sigma^2)$

μ_{ij} is further split into the following parts:

(i) The general mean effect μ given by

$$\mu = \frac{\sum \sum \mu_{ij}}{N} \quad \text{--- (2)}$$

(ii) The effect α_i ($i=1, 2, \dots, k$) due to the i th ration given by

$$\alpha_i = \mu_{i\cdot} - \mu$$

$$\text{where } \mu_{i\cdot} = \frac{1}{h} \sum_{j=1}^h \mu_{ij} \quad ; \quad (i=1, 2, \dots, k)$$

$$\text{obviously } \sum_{i=1}^k \alpha_i = 0 \quad \text{--- (3)}$$

(iii) The effect β_j , ($j=1, 2, \dots, h$) due to j th variety (breed of cow) given by

$$\beta_j = \bar{u}_{ij} - u$$

$$u_{ij} = \frac{1}{k} \sum_{i=1}^k \bar{u}_{ij}; (j=1, 2, \dots, h) \quad (4)$$

Obviously $\sum_{j=1}^h \beta_j = 0$

- (iv) The interaction effect v_{ij} , when the i^{th} level of first factor (varieties) and j^{th} level of second factor (breed of cow) occur simultaneously and given by

$$v_{ij} = \bar{u}_{ij} - \bar{u}_{i\cdot} - \bar{u}_{\cdot j} + u$$

$$\sum_j v_{ij} = 0, \quad i = 1, 2, \dots, k$$

$$\sum_i v_{ij} = 0, \quad j = 1, 2, \dots, h$$

and thus we have

$$\bar{u}_{ij} = u + (\bar{u}_{i\cdot} - u) + (\bar{u}_{\cdot j} - u) + (v_{ij} - \bar{u}_{ij})$$

(5)

and consequently (1) becomes

$$x_{ij} = u + d_i + \beta_j + v_{ij} + \epsilon_{ij} \quad (6)$$

where ϵ_{ij} is the error effect due to chance &

$$\sum_{i=1}^k d_i = 0 = \sum_{j=1}^h \beta_j$$

$$\therefore \sum_i v_{ij} = \sum_j v_{ij} = 0$$

As there is only one obsrⁿ in each cell, the obsrⁿ corresponding to the i^{th} level of ration & j^{th} level of breed of cow is only one v.e.

x_{ij} : But we cannot estimate by one value alone. Hence in this case of one obsrⁿ-per cell the interaction effect $\gamma_{ij}=0$ and the model (6) reduces to

$$x_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij} \quad \text{--- (7)}$$

Statistical Analysis of the Model

Let us write

$$\begin{aligned} \bar{x}_{i\cdot} &= \text{Mean yield of } i^{\text{th}} \text{ ration} \\ &= \frac{1}{h} \sum_{j=1}^h x_{ij} \quad (i = 1, 2, \dots, k) \end{aligned}$$

$$\begin{aligned} \bar{x}_{\cdot j} &= \text{mean yield of } j^{\text{th}} \text{ variety} \\ &= \frac{1}{k} \sum_{i=1}^k x_{ij} \quad (j = 1, 2, \dots, h) \end{aligned}$$

$$\bar{x}_{\cdot \cdot} = \text{overall mean} = \frac{1}{hk} \sum_{i=1}^k \sum_{j=1}^h x_{ij}$$

$$= \frac{1}{h} \sum_j \left(\frac{1}{k} \sum_{i=1}^k x_{ij} \right) = \frac{1}{h} \sum_j \bar{x}_{\cdot j}$$

$$= \frac{1}{k} \sum_i \left(\frac{1}{h} \sum_{j=1}^h x_{ij} \right) = \frac{1}{k} \sum_i \bar{x}_{i\cdot}$$

The least square estimates of the parameters μ , α_i and β_j are obtained on minimizing

the error sum of squares

$$E = \sum_{i=1}^k \sum_{j=1}^n e_{ij}^2 = \sum_{ij} (x_{ij} - \mu - \alpha_i - \beta_j)^2$$

The normal eqⁿ: for eq estimating μ , α_i & β_j are respectively

$$\frac{\partial E}{\partial \mu} = 0 = -2 \sum_{ij} (x_{ij} - \mu - \alpha_i - \beta_j)$$

$$\frac{\partial E}{\partial \alpha_i} = 0 = -2 \sum_j (x_{ij} - \mu - \alpha_i - \beta_j)$$

$$\frac{\partial E}{\partial \beta_j} = 0 = -2 \sum_i (x_{ij} - \mu - \alpha_i - \beta_j)$$

since $\sum_i \alpha_i = 0 = \sum_j \beta_j$; we get from the above eqⁿ:

$$\hat{\mu} = \frac{1}{nk} \sum_{ij} x_{ij} = \bar{x}_{..}$$

$$\hat{\alpha}_i = \frac{1}{n} \sum_j x_{ij} - \hat{\mu} = \bar{x}_{i..} - \bar{x}_{..}$$

$$\hat{\beta}_j = \frac{1}{k} \sum_i x_{ij} - \hat{\mu} = \bar{x}_{.j} - \bar{x}_{..}$$

Thus linear model (7) becomes

$$x_{ij} = \bar{x}_{..} + (\bar{x}_{i..} - \bar{x}_{..}) + (\bar{x}_{.j} - \bar{x}_{..}) \\ + (x_{ij} - \bar{x}_{i..} - \bar{x}_{.j} + \bar{x}_{..})$$

The error term e_{ij} being so chosen that both sides are equal

Transposing $\bar{x}_{..}$ to the left side, squaring

and summing both sides over $i \& j$, we get

$$\begin{aligned} \sum_{i,j} (x_{ij} - \bar{x}_{..})^2 &= \sum_{i,j} [(x_{ij} - \bar{x}_{i.} - \bar{x}_{.j} + \bar{x}_{..}) + (\bar{x}_{i.} - \bar{x}_{..}) \\ &\quad + (\bar{x}_{.j} - \bar{x}_{..})]^2 \\ &= \sum_{i,j} (x_{ij} - \bar{x}_{i.} - \bar{x}_{.j} + \bar{x}_{..})^2 + \sum_{i,j} (\bar{x}_{i.} - \bar{x}_{..})^2 \\ &\quad + \sum_{i,j} (\bar{x}_{.j} - \bar{x}_{..})^2 + 2 \sum_{i,j} (x_{ij} - \bar{x}_{i.} - \bar{x}_{.j} \\ &\quad + \bar{x}_{..}) \\ &\quad (\bar{x}_{i.} - \bar{x}_{..}) + 2 \sum_{i,j} (x_{ij} - \bar{x}_{i.} - \bar{x}_{.j} + \bar{x}_{..}) \\ &\quad (\bar{x}_{.j} - \bar{x}_{..}) + 2 \sum_{i,j} (\bar{x}_{i.} - \bar{x}_{..}) (\bar{x}_{.j} - \bar{x}_{..}) \end{aligned}$$

since algebraic sum of deviation of a set of observations about their mean is zero.

similarly other product terms also vanishes.
Hence

$$\begin{aligned} \sum_{i,j} (x_{ij} - \bar{x}_{..})^2 &= \sum_{i,j} (x_{ij} - \bar{x}_{i.} - \bar{x}_{.j} + \bar{x}_{..})^2 + h \sum_i (\bar{x}_{i.} - \bar{x}_{..})^2 \\ &\quad + k \sum_j (\bar{x}_{.j} - \bar{x}_{..})^2 \end{aligned}$$

$$S_T^2 = S_E^2 + S_T^2 + S_V^2$$

T.S.S. = S.S. due to error + S.S. due to treatments
+ S.S. due to varieties.

Null hypothesis: — We set up the null hypothesis that treatment as well as varieties are homogeneous. In other words, the null hypothesis for treatment and varieties are respectively

$$H_T: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_K = \mu$$

$$H_V: \mu_1 = \mu_2 = \dots = \mu_h = \mu$$

Q1

$$H_t : d_1 = d_2 = \dots = d_k = 0$$

$$H_r \Rightarrow \beta_1 = \beta_2 = \dots = \beta_k = 0$$

Degrees of freedom for various S.S. :- The total S.S., S_T^2

being computed from $N = hk$ quantities, $(x_{ij} - \bar{x}_{..})$, which are subject to one linear constraint $\sum (x_{ij} - \bar{x}_{..}) = 0$, will carry $(N-1)$ d.f.

similarly for S_T^2 will be based on $(k-1)$ d.f.
since $\sum_i (\bar{x}_{ij} - \bar{x}_{..}) = 0$ +

S_V^2 will have $(h-1)$ d.f. since $\sum_j (\bar{x}_{..j} - \bar{x}_{..}) = 0$
will carry $(h-1)$ d.f.

and \sum & S_E^2 will carry $(N-1) - (k-1) - (h-1)$
 $= (h-1)(k-1)$ d.f.

Thus the partitioning of d.f. is as follows

$$(hk-1) = (k-1) + (h-1) + (h-1)(k-1)$$

which \Rightarrow that d.f. are additive.

Test statistic :- To obtain appropriate test statistic to test the hypothesis H_t & H_r , we need the expectation of the various mean sum of squares due to each independent factors.

$$\text{M.S.S. due to treatment} = \frac{S_T^2}{k-1} = S_T^2$$

$$M.S.S. \text{ due to variety} = \frac{S_V^2}{h-1} = \frac{\sigma_e^2}{h-1}$$

$$\text{error M.S.S.} = \frac{S_E^2}{(h-1)(k-1)} = \frac{\sigma_e^2}{(h-1)(k-1)}$$

Since various S.S. and d.f. are additive and since under H_T & H_V , each of S_T^2 , S_V^2 & S_E^2 provides an unbiased estimates of σ_e^2 , under the assumption of normal parent popn., we get by Cochran's thm

$\frac{S_T^2}{\sigma_e^2}$, $\frac{S_V^2}{\sigma_e^2}$ & $\frac{S_E^2}{\sigma_e^2}$ are mutually indept.

χ^2 with $(K-1)$, $(h-1)$ & $(h-1)(k-1)$ d.f. respectively
Hence under H_T and H_V respectively, we get

$$F_T = \frac{\frac{S_T^2}{\sigma_e^2}}{\frac{(h-1)(k-1)}{h-1, (h-1)(k-1)}} = \frac{\frac{S_T^2}{\sigma_e^2}}{\frac{S_E^2}{h-1, (h-1)(k-1)}} \sim F_{K-1, (h-1)(k-1)}$$

$$F_V = \frac{\frac{S_V^2}{\sigma_e^2}}{\frac{(h-1)(k-1)}{h-1, (h-1)(k-1)}} = \frac{\frac{S_V^2}{\sigma_e^2}}{\frac{S_E^2}{h-1, (h-1)(k-1)}} \sim F_{h-1, (h-1)(k-1)}$$

ANOVA for two-way classified Data

sources of variat-	Sum of squares	d.f.	M.S.S.	Variate ratio
Treatments	$S_T^2 = \sum_i h(\bar{x}_{i..} - \bar{x}_{..})^2$	$K-1$	$S_T^2 = \frac{S_T^2}{K-1}$	$F_T = \frac{S_T^2}{S_E^2}$ $a(h-1)(h-1)K-1$
varieties	$S_V^2 = \sum_j k(\bar{x}_{j..} - \bar{x}_{..})^2$	$h-1$	$S_V^2 = \frac{S_V^2}{h-1}$	$F_V = \frac{S_V^2}{S_E^2}$
residuals	$S_E^2 = \sum_{ij} (x_{ij} - \bar{x}_{i..} + \bar{x}_{j..} + \bar{x}_{..})^2$	$(h-1)(k-1)$	$S_E^2 = \frac{S_E^2}{(h-1)(k-1)}$	$(h-1), (h-1), (h-1)$
Total		$N-1$		