ANALYSIS OF VARIANCE (ANOVA)

INTRODUCTION

The test of significance based on t-distribution is an adequate procedure only for testing the significance of the difference between 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 population I.e., the means are equal.

For example, three types of fertilisers are applied to five plots each and their yields on each of the plot is given as follows

plots	Yield of wheat in tons		
	fertiliser	Fertiliser	fertiliser C
	А	В	
1	20	18	25
2	21	20	25
3	23	17	25
4	16	15	25
5	20	25	25
Mean	100/5=20	95/5=19	125/5=25

We have to study if the effect of these fertilisers on the yield is significantly different, or in other words the samples are from the same population. The answer to this is provided by the technique of analysis of variance.

The basic purpose of analysis of variance is to test the homogeneity of several means.

The term 'Analysis of variance' was introduced by Prof. R.A. Fisher in 1920's.

Variation is inherent in nature,

The total variation in any set of numerical data is due to a number of causes which may be calculated as i) assignable causes and ii) chance causes

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

Examples of assignable causes of variation

Inappropriate procedures, substandard raw materials, measurement errors, temperature etc.,

DEFINITION:

According to Prof. R.A. Fisher, Analysis of Variance (ANOVA) is the "Separation of variance ascribable to one group of causes from the variance ascribable to other group"

The 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(causes), with the estimate due to chance factor (causes). The later being known as experimental error.

ASSUMPTIONS FOR ANOVA TEST

ANOVA test is based on the test statistics F (variance Ratio)

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

- i) The observations are independent,
- ii) Parent population from which observations are taken is normal, and
- iii) Various treatment and environmental effects are additive in nature.

IMPORTANCE:

ANOVA technique enables us to compare several population means simultaneously and thus results in lot of savings in time and money

The origin of ANOVA technique lies in agricultural experiments but it finds its applications in almost all types of design of experiments in various diverse fields such as in industry, education, psychology, business etc.,

The ANOVA technique is not designed to test the equality of several population variances. It's objective is to test the equality of several population means or the homogeneity of several independent sample means.

In addition to testing homogeneity of several sample means, the ANOVA technique is now frequently applied in testing the linearity of the fitted regression line or the significance of the correlation ratio.

MATHEMATICAL MODEL

FIXED EFFECT MODEL AND RANDOM EFFECT MODELS

A fixed effects model is a statistical model in which the model parameters are fixed or non-random quantities.

In random effects model in which all or some of the model parameters are random variables.

Fixed effects model:

Suppose the k-levels of the factor (treatments) under consideration are the only levels of interest and all these are included in the experiment by the investigator or out of a large number of classes, the k classes (treatments) in the model have been specifically chosen by the experimenter. In such a case α_i 's the effect of the ith treatment [$\alpha_i = (\mu_i - \mu)$] are fixed constants (unknown) and the model is a fixed effect model.

In the fixed effect model, the conclusions about the test of hypothesis regarding the parameters α_i 's will apply only to k-treatments (factor levels) considered in the experiment.

These conclusions cannot be extended to other remaining treatments (factors) which are not considered in the experiment.

Random effects model:

Suppose we have a large number of classes (treatments) and we want to test, through an experiment if all these class effects are equal or not. Due to consideration of time, money or administrative convenience it may not be possible to include all the factor levels in the experiment. In such a situation, we take only a random sample of factor levels in the experiment and after studying and analysing the sample data, we draw conclusions which would be valid for all the factor levels whether included in the experiment or not. In such a situation the parameters α_i 's in the model will not be fixed constants but will be random samples and the model is random effect model.

In the random effect model if the null hypothesis of the homogeneity of class (treatment) effects is rejected, then to test to test the difference between two class(treatments) effects we cannot apply the t-test because all treatments are not included in the experiment.

5.2. ONE-WAY CLASSIFICATION

Let us suppose that N observations y_{ij} , $(i = 1, 2, ..., k; j = 1, 2, ..., n_i)$ of a random variable Y are grouped, on some basis, into k classes of sizes $n_1, n_2, ..., n_k$ respectively, $\left(N = \sum_{i=1}^k n_i\right)$ as

exhibited	m	Table	9.1.	
				() =

TARIE	5.1 : ONE-WAY	CLASSIFIED	DATA
IABLE	DI UNE-WAI	OLAUOII ILD	P

Class	Sample Observations			ALL PROPERTY AND A REAL PR		Mean	
1	<i>y</i> ₁₁	y ₁₂	ffeet Me	y 1n1	R Tliebol	1 is y11 be	
er of side	y ₂₁	y ₂₂	the facto	<i>y</i> _{2<i>n</i>₂}	$ \frac{1}{2} T_2 = T_2$	\overline{y}_{2}	
the investi (5-1) have	he model	ents) in t	atieneern s (treatm	ie k-classe	of claises, th	arge number	
ants ⁱ unka	y _{i1}	Yi2	h a tase	y _{in,}	e expiriment move af th	$\begin{array}{c} y_i, y_i, y_i \\ y_i, y_i \\ y_i, y_i \\ y_i$	
	: 	: Y _{k2}	s. In the	: Y _{kn} ,	\dot{T}_{k}	$\frac{1}{\overline{y}_{k}}$	

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

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

(*ii*) The variation *within the classes, i.e.*, the 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 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 examine if there is significant difference between the class means in view of the inherent variability within the separate classes.

In particular, let us consider the effect of k different rations on the yield in milk of N cows (of the same breed and stock) divided into k classes of sizes $n_1, n_2, ..., n_k$ respectively,

 $N = \sum_{i=1}^{n} n_i$. Here the sources of variation are :

- (i) Effect of the ration (treatment): t_i ; i = 1, 2, ..., k.
- (ii) Error (ϵ) produced by numerous causes of such magnitude that they are not detected and identified with the knowledge that we have and they together produce a variation of random nature obeying Gaussian (Normal) law of errors.

Mathematical Model. In this case the linear mathematical model will be :

$$y_{ij} = \mu_i + \varepsilon_{ij} = \mu + (\mu_i - \mu) + \varepsilon_i$$

- $= \mu + \alpha_i + \varepsilon_{ij}$; where $(i = 1, 2, ..., k; j = 1, 2, ..., n_i)$ (5.1)
- (i) y_{ij} is the yield from the *j*th cow, $(j = 1, 2, ..., n_i)$ fed on the *i*th ration (i = 1, 2, ..., k). ...(5.2)

(*ii*) μ is the general mean effect given by :

$$\mu = \sum_{i=1}^{k} n_i \mu_i / N$$
 with our dominant of a matrix of bolic \dots (5.2a)

where μ_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 by μ_i ,

(*iii*) α_i is the effect of the *i*th ration given by : $\alpha_i = \mu_i - \mu$, (i = 1, 2, ..., k) ... (5.2b) *i.e.*, the *i*th ration increases (or decreases) the yield by an amount α_i . On using (5.2a) and (5.2b), we get

$$\sum_{i=1}^{\kappa} n_i \alpha_i = \sum_i n_i (\mu_i - \mu) = \sum_i n_i \mu_i - \mu \sum n_i = N, \ \mu - \mu, \ N = 0 \qquad \dots (5.2c)$$
(*iv*) ε_{ii} is the error effect due to chance. $\dots (5.2d)$

ANOVA FOR FIXED EFFECT MODEL

The fixed effect or parametric model used is

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

 $= \mu + \alpha_i + \varepsilon_{ij}$ (i=1,2...k; j=1,2...n_i) where $\alpha_{i} = \mu_{i} - \mu_{i}$

where y_{ij} is the yield from the ith row and jth column

 μ is the general mean effect given by $\mu = \sum_{i=1}^{k} ni\mu i/N$

 μ_i is the fixed effect due to ith treatment

 α_i is the effect of the ith treatment given by

 $\alpha_i=\mu_i-\mu$

 ε_{ij} is the error effect due to chance

ASSUMPTIONS IN THE MODEL

- i) all the observations $(y_{ij's})$ are independent and $y_{ij} \sim N(\mu_i, \sigma_e^2)$
- ii) different effects are additive in nature

iii) ϵ_{ij} are i.i.d., N(0, σ_e^2)

under the third assumption, the model becomes

E(y_{ij}) = μ_i = $\mu + \alpha_i$ (i=1,2...k; j=1,2...n_i)

STATISITICAL ANALYSIS OF THE MODEL

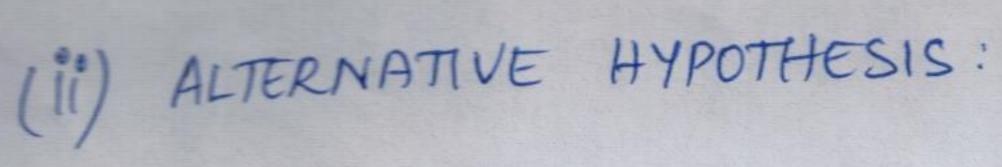
Null hypothesis: we have to test the equality of the population means. Hence the null hypothesis is given by

 $H_0 = \, \mu_1 = \mu_2 \, = \ldots = \mu_k = \mu$

which reduces to $H_0 = \alpha_1 = \alpha_2 = \ldots = \alpha_k = 0$ since $\alpha_{i=} \mu_{i-} \mu$

Alternate hypothesis:

 $H_1: \text{ at least two of the means } \mu_1, \mu_2, \ldots \mu_k \text{ are different.}$ let us write



H = M, # M2 # # MK

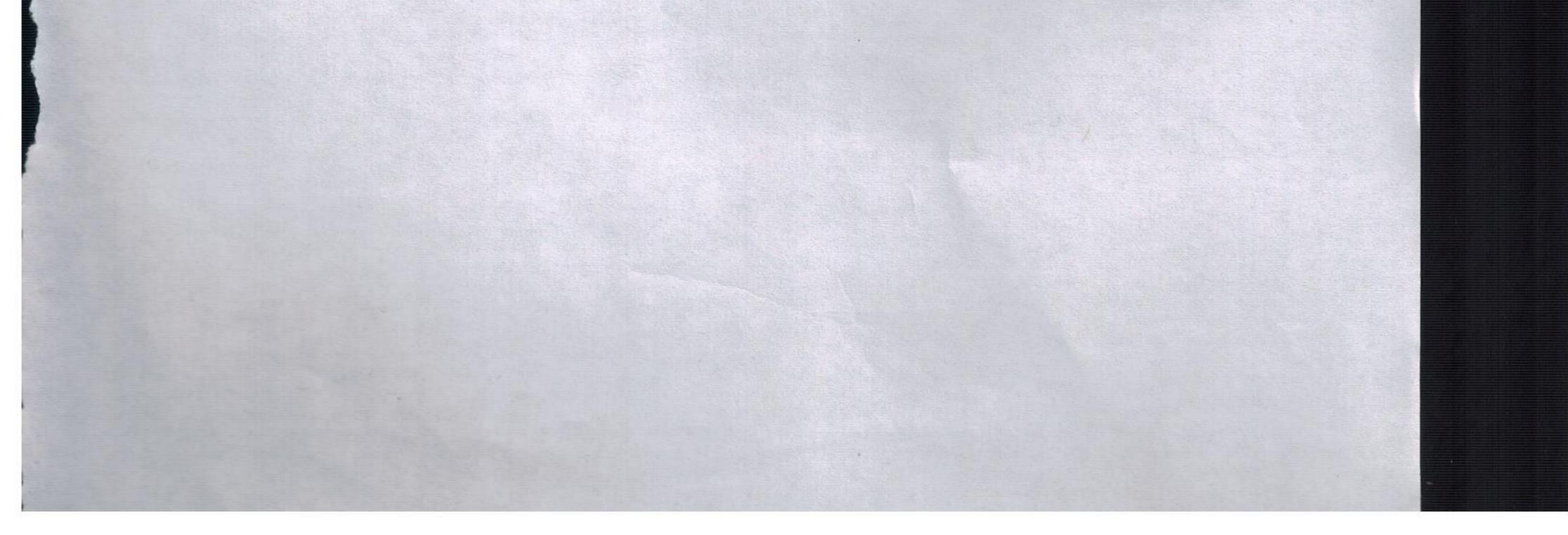
we know that,

Yi. = mean of the ith class $= \underbrace{z}_{j=1}^{n_{i}} y_{ij}^{i} \\ = \underbrace{y_{j=1}^{n_{i}}}_{n_{i}} i = 1, 2, 3 \cdots , k.$

 $\overline{y}_{..} = \text{overall mean} = \frac{1}{N} \stackrel{k}{\leq} \stackrel{n_i}{\leq} \stackrel{n_i}{(y_{ij})}$

 $= \frac{1}{N} \stackrel{k}{\leq} n_i \overline{y_i}.$ (iii) Least square Estimates of parameteres. The parameters je & a' given & above model are given by above of the model Yij = µ + x q + Eij are estimated by the principle of least squares, minimising the everos sum of squares. The mathematical models is $y_{ij} = \mu + \alpha_i + \epsilon_{ij} - 0$ $\Rightarrow E_{ij} = y_{ij} - \mu - \alpha_i - \emptyset$

squaring ξ taking $\xi \leq 0$ both sides we get i $E = \leq \leq E_{j}^{2} = \leq \leq (y_{j} - \mu - \alpha_{j})^{2}$ Estimation of pe. $E = \underset{i}{\leq} \underset{j}{\leq} (9ij - \mu - \alpha_i)^2$ $\frac{\partial E}{\partial E} = -2 \neq \leq \leq (y_{ij} - \mu - \alpha_{i}) = 0$ ope = ミミリッジョーション・ション・= 0 $= \underbrace{\leq}_{j} \underbrace{\leq}_{j} \underbrace{}_{j} - N\mu - \underbrace{\leq}_{j} n_{j} \alpha_{j} = 0$ $\Rightarrow \mu = \stackrel{\leq}{\Rightarrow} \stackrel{\leq}{\Rightarrow} \stackrel{\leq}{y} \stackrel{y}{y} = \stackrel{=}{y}$ $\therefore \leq n_i \alpha_i = 0$ from note 1 je = J..



Estimate of x; $E = \xi \xi \left(y_{ij} - \mu - \alpha_i \right)^2$ $\frac{\partial E}{\partial \alpha_i} = -2 \leq (y_{ij} - \mu - \alpha_i) = 0$ $= \leq y_{ij} - \leq \mu - \leq \alpha_i^{\circ} = 0$ $= \underbrace{\xi}_{j} \underbrace{\xi}_{j} \underbrace{\xi}_{j} - n_{i} \widehat{\mu} - n_{i} \alpha_{i} = 0$ $\Rightarrow n_i \alpha_i^{\alpha} = \leq y_{ij}^{\alpha} - n_i^{\alpha} \mu$ $\alpha_i^{\circ} = \frac{1}{n_i^{\circ}} \stackrel{i}{\equiv} y_{ij}^{\circ} - \frac{n_i^{\circ} \hat{\mu}_i}{n_i^{\circ}}$ $= \int = \int \frac{1}{j} \frac{1}$ = y: - y.. hence the model becomes $y_{ij} = \mu + \alpha_i^\circ + \epsilon_{ij}^\circ$ $=\overline{y}_{\cdot,\cdot}+(\overline{y}_{\cdot,\cdot}-\overline{y}_{\cdot,\cdot})+(\overline{y}_{\cdot,\cdot}^{\circ}-\overline{y}_{\cdot})$

TOTAL SUM OF SQUARE: T.S.S = S.S.E + S.S.T $T.S.S = \underbrace{\stackrel{k}{\leq}}_{i=1} \underbrace{\stackrel{n}{\leq}}_{j=1} (y_{ij} - \overline{y}_{..})^2$

 $= \underbrace{\xi}_{i=1}^{k} \underbrace{(y_{ij} - \overline{y_{i}}_{i} + \overline{y_{i}}_{i} - \overline{y_{i}}_{i}}_{i=1} \underbrace{(y_{ij} - \overline{y_{i}}_{i} + \overline{y_{i}}_{i} - \overline{y_{i}}_{i})^{2}}_{i=1}$ $= = = (y_{ij} - \overline{y_{i}})^2 + = n^{\circ} (\overline{y_i} - \overline{y_{i}})^2 +$ $2\left[= \{(\bar{y}_{i}, -\bar{y}_{i}) \geq (y_{i}, -\bar{y}_{i})\} \right]$ $= \underbrace{\neq \leq}_{j} (\underbrace{y_{j}}_{j} - \underbrace{y_{i}}_{j})^{2} + \underbrace{\leq ni}_{i} (\underbrace{y_{i}}_{i} - \underbrace{y_{i}}_{j})^{2}$

From 3, we get. $S_E^2 = S \cdot S \cdot E = \underset{i}{\leq} \underset{j}{\leq} (y_{ij} - \overline{y}_{i})^2$ $S_t^2 = S \cdot S \cdot T = \leq n \left(\overline{y_i} - \overline{y_i} \right)^2$

DEGREES OF FREEDOM!

(i) The degrees of freedom for $T.S.S(S_T^2)$ is N-1.

(ii) The degrees of freedom for treatment S. Stillis K-1 (iii) The def of error sum of square (S_E^2) is N-k.

MEAN SUM OF SQUARES (M.S.S)

(i) The M.S.S due to treatments

 $= \frac{S_t^2}{k-1} = \Re_t^2$

(ii) M.S.S due to Emor $= SE^2 = \$E^2$ N-K

ANOVA TABLE !

The ANOVA table for one way

d.f

classified data

Sources of variation

sum of Squares mean sum of squares

variance ratio

Treatment

 S_t^2

 $k-1 \quad \varphi_t^2 = \frac{St^2}{k-1}$

 $F = \frac{\varphi_{E}^{2}}{\varphi_{E}^{2}}$ ~ F_K-1, N-K

Error

Total

SE²

 $N-k \quad \not \in \vec{E} = \frac{S\vec{E}^2}{N-k}$

ST2

N-1

Conclusion.

 $\int_{\mathcal{F}} F = \frac{\mathscr{F}t^2}{\mathscr{F}\varepsilon^2} \text{ is less than the table}$ value of F_{k-19}N-k. , we accept Ho otherwise reject Ho.

VARIANCE OF THE ESTIMATES :

The lineare model for one way classification is given by $y_{ij}^{\circ} = \mu + \alpha_{i} + \epsilon_{ij}^{\circ} j = 1, 2, ..., k$ $j = 1, 2, ..., n_{i}^{\circ}$

The least square estimates of μ & α ?

are given by



 $\hat{\mu} = \overline{y_{\cdot \cdot}} \quad \underbrace{\varepsilon}_{i} \quad \underbrace{\widetilde{z}_{i}}_{i} = \underbrace{\overline{y}_{i}}_{i} - \underbrace{\overline{y}_{\cdot \cdot}}_{i} - \underbrace{O}$ Then, Var $(\hat{\mu}) = E[\hat{\mu} - E(\hat{\mu})]^2$ $= E\left[\overline{y}:=E(\overline{y}:)\right]^{2}$ $= E(\overline{g}, -\mu)^2$ $= E(\overline{e}^{2}) = Var(\overline{e}) = \underbrace{\nabla \overline{e}}_{N}^{2}$ But $\hat{\alpha}_{i} = \hat{y}_{i} - \hat{y}_{i}$ we know that $\overline{y_i} = \mu + \alpha_i + \overline{E_i}$ and $\overline{y}_{\cdot \cdot} = \mu + \overline{\epsilon}_{\cdot \cdot}$

Then we have $\hat{a_i} = \mu + \hat{a_i} + \hat{\epsilon_i} - (r + \hat{\epsilon}_{-})$ $= d_i + E_i - E_i$ $E(\hat{\alpha}_i) = d_i \quad E_i \in \mathcal{N}(0, \sigma_e^2)$ $\therefore \hat{a}_i - E(\hat{a}_i) = \overline{E}_i - \overline{E}_i$ $V(\hat{a}_i) = E[\hat{a}_i - E(\hat{a}_i)]$ $= E \left[E_i - E_i \right)^{-1}$ $= E(\bar{e}_{i}^{2}) + E(\bar{e}_{i}^{2}) - 2E(\bar{e}_{i}^{2},\bar{e}_{i}) - 2E(\bar{e}_{i}^{2},\bar{e}) - 2E(\bar{e}_{i}^{2$

But $E[\overline{e_i}, \overline{e_i}] = E[(\frac{1}{n_i} \underbrace{e_i})(-\underbrace{e_i})(-\underbrace{e_i})]$ $= \frac{1}{n_{1}^{\circ}N} E \left[\left(e_{11}^{\circ} + e_{12}^{\circ} + \dots + e_{n_{1}}^{\circ} \right) \right] + \frac{e_{11} + e_{12} + \dots + e_{n_{1}}}{e_{k_{1}} + e_{k_{2}} + \dots + e_{n_{1}}} + \frac{e_{k_{1}} + e_{k_{2}} + \dots + e_{n_{1}}}{e_{k_{1}} + e_{k_{2}} + \dots + e_{n_{1}}}$ $=\frac{1}{h_{i}^{\circ}N}\left[E\left(e_{i}\right)+e_{i}+e_{i}+e_{i}\right)^{2}\right]$

all co-variance terms vanish, since Eij's

are un conclated.

 $E(\overline{E_{i}}, \overline{E_{i}}) = \frac{1}{n_{i}N} \left[\operatorname{var}(\overline{E_{i}}) + \operatorname{var}(\overline{E_{i}}) + \operatorname{var}(\overline{E_{i}}) + \operatorname{var}(\overline{E_{i}}) \right]$ = 1 nd Je² Min

substituting 3 foora in 2 jareget. $: \operatorname{var}(\hat{a}_i) = E(\overline{e}_i^2) + E(\overline{e}_i^2) - 2E(\overline{e}_i, \overline{e}_i)$ $= \frac{\sigma_e^2}{n_i^2} + \frac{\sigma_e^2}{N} - 2 \frac{\sigma_e^2}{N}$ $= \frac{\sigma e^2}{ni} - \frac{\sigma e^2}{N}$ $= \sigma e^2 \left(\frac{1}{n_i} - \frac{1}{N} \right)$

Two WAY CLASSIFICATION: Suppose n obsuevations are classified into K categories for classes), say A, 1, A2, ..., AK according to some caterion A; and into h categories say B, B2...BL according to some criterion B, having Kh Combinations Ai, Bj; i=1,2...k, j=1, 2...h; ofter called cells.

This Scheme of classification according to two factors or criteria is called two-way classification and its analysis is called two-way analysis of variance. The number of observations in each cell may le equal or different, but we shall consider about in por cell so that

So that n=hk, se, the total number of cells is n=hk In the two way classification, the values of the response Variables are affected by too factors.

For example. The yield of milk may be affected by differences in treatment, (i.e) rations as well as the differences in variety, (ie) breed and atock of the cours. Let us now duppose that the n Cours are divided into n different groups or classes according to their breed and about, each groups containing k cous and then let us consider the affect of K treatment (ie) rations given at random to caps in each geoup on the yield of milk. Let Yij = Yield of milk from the cow jth breed or stock,

fed on the ration i

i=1,2,... k; j=1,2,...h

H₁: Albast two
$$\mu_{j}$$
's are different
H₁₂: Albast two μ_{j} 's are different
Least adquare Estimates of parameters.
The linear model is given by,
 $J_{ij} = \mu + \alpha_{i} + \mu_{j} + \epsilon_{ij}$
 $\epsilon_{ij} = J_{ij} - \mu - \alpha_{i} - \mu_{j}$
Why applying the principal q hast adquare we get.
 $E = \sum_{i,j} \epsilon_{ij}^{2} = \sum_{j,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})^{2}$
The normal equation for estimating
 $\frac{d\epsilon}{d\mu} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$
 $\frac{dE}{dki} = 0 \Rightarrow -\vartheta \sum_{i,j} (y_{ij} - \mu - \alpha_{i} - \mu_{j})$

$$\partial i = \mu + \alpha_i + \beta_j + c_{ij}$$

Where
$$\mu = \sum_{i} \sum_{j} \mu_{ij}$$

 $\mu_{i} = \frac{1}{n} \sum_{j=1}^{h} \mu_{ij}$
 $\mu_{j} = \frac{1}{k} \sum_{i=1}^{k} \mu_{ij}$
 $\alpha_{i} = \mu_{i} - \mu; \sum_{i=1}^{k} \alpha_{i} = 0$
 $\beta_{j} = \mu_{j} - \mu; \sum_{j=1}^{h} \beta_{j} = 0$

Statistical Analysis of a fixed effect model.

The treatments and the varieties of homogenius. Ho, : The treatment are homogeneous Ho2 : The variety are homogeneous i.e), Ho1 : H1.= H2.=...= Mi.=M Ho2 : M.)= H.2=...=Mj.=M

Alternative Hypothesis:

Hit : At least two of the Mi.'s are different. Hiv : At least two of the Mij's are different Or their equalent.

$$\begin{split} \hat{\mu} &= \frac{1}{h\kappa} \sum_{j} \sum_{j} 9_{ij} = \widehat{9} .. \\ \hat{\alpha}_{i} &= \frac{1}{n} \sum_{j} 9_{ij} - \widehat{\mu} \\ &= \widehat{9}_{i} - \widehat{9}_{i} .. \\ \hat{\beta}_{j} &= \frac{1}{k} \sum_{j} 9_{ij} - \widehat{\mu} \\ &= \widehat{9}_{ij} - \widehat{9}_{i} .. \\ \hat{\beta}_{j} &= \frac{1}{k} \sum_{j} 9_{ij} - \widehat{\mu} \\ &= \widehat{9}_{ij} - \widehat{9}_{i} .. \\ \text{Thus the linear model becomes} \\ g_{ij} &= \widehat{y}_{..} + (\widehat{y}_{i} - \widehat{9}_{..}) + (\widehat{y}_{ij} - \widehat{9}_{..}) + (\widehat{y}_{ij} - \widehat{9}_{..} - \widehat{9}_{.j} + \widehat{9}_{..}) \\ \text{Poutitioning of the dums } q \ dquares. \\ \sum_{j} \sum_{j} (9_{ij} - \widehat{9}_{..})^{2} &= \sum_{j} \sum_{j} \left[(9_{ij} - \widehat{9}_{i} - \widehat{9}_{.j} + \widehat{9}_{..}) + (\widehat{9}_{i,} - \widehat{9}_{..})^{2} \\ &= \sum_{j} \sum_{j} (9_{ij} - \widehat{9}_{..})^{2} = \sum_{j} \sum_{j} (9_{ij} - \widehat{9}_{i,} - \widehat{9}_{.j} + \widehat{9}_{..}) + (\widehat{9}_{i,} - \widehat{9}_{..})^{2} \\ &= \sum_{j} \sum_{j} (9_{ij} - \widehat{9}_{i,} - \widehat{9}_{.j} + \widehat{9}_{..}) + 2\sum_{j} \sum_{j} (\widehat{9}_{.j} - \widehat{9}_{..})^{2} \\ &= \sum_{j} \sum_{j} (9_{ij} - \widehat{9}_{i,} - \widehat{9}_{.j} + \widehat{9}_{..}) + 4\sum_{j} \sum_{j} (\widehat{9}_{.j} - \widehat{9}_{..}) \\ &(\widehat{9}_{ij} - \widehat{9}_{i,} - \widehat{9}_{.j} + \widehat{9}_{..}) + 2\sum_{j} \sum_{j} (\widehat{9}_{.i} - \widehat{9}_{..}) (\widehat{9}_{.j} - \widehat{9}_{..}) \\ &(\widehat{9}_{ij} - \widehat{9}_{i,} - \widehat{9}_{.j} + \widehat{9}_{..}) + 2\sum_{j} \sum_{j} (\widehat{9}_{.i} - \widehat{9}_{..}) (\widehat{9}_{.j} - \widehat{9}_{..}) \\ \end{aligned}$$

•

Now.

$$\begin{split} \sum_{i,j}^{W} \left(\overline{y}_{i}, -\overline{y}_{\cdot} \right) \left(\underline{y}_{ij} - \overline{y}_{i}, -\overline{y}_{\cdot j} + \overline{y}_{\cdot} \right) &= \sum_{i}^{\infty} \left[\left(\overline{y}_{i}, -\overline{y}_{\cdot} \right) \sum_{j}^{\infty} \left(\underline{y}_{ij} - \overline{y}_{i}, -\overline{y}_{\cdot j} + \overline{y}_{\cdot} \right) \right] \\ &= \sum_{i}^{\infty} \left[\left(\overline{y}_{i}, -\overline{y}_{\cdot} \right) \sum_{j}^{0} \sum_{j}^{0} \left(\underline{y}_{ij} - \overline{y}_{i}, \right) - \sum_{j}^{0} \left(\overline{y}_{ij} - \overline{y}_{\cdot}, \right) \sum_{j}^{0} \right] = 0, \\ \therefore \sum_{i}^{\infty} \sum_{j}^{\infty} \left(\underline{y}_{ij} - \overline{y}_{\cdot}, \right)^{2} &= h \sum_{i}^{\infty} \left(\overline{y}_{i}, -\overline{y}_{\cdot}, \right)^{2} + K \sum_{j}^{\infty} \left(\overline{y}_{ij} - \overline{y}_{\cdot}, \right)^{2} + \sum_{i,j}^{\infty} \left(\underline{y}_{ij} - \overline{y}_{i}, -\overline{y}_{i} + \overline{y}_{j} + \overline{y}_{j} \right) \\ O^{n}, \qquad S_{T}^{2} = S_{T}^{2} + S_{V}^{2} + S_{V}^{2} \\ Mhere, \qquad S_{T}^{2} = \sum_{i,j}^{\infty} \sum_{j}^{0} \left(\underline{y}_{ij} - \overline{y}_{\cdot}, \right)^{2} \text{ is the total s.s} \\ S_{T}^{2} &= h \sum_{i,j}^{\infty} \left(\overline{y}_{ii} - \overline{y}_{\cdot}, \right)^{2} \text{ is s.s due to treatment} \end{split}$$

 $S_{v}^{2} = K \sum_{j} (y_{j} - \overline{y}_{o})^{2} is$ the s.s due to varieties,

and,

$$S_E^2 = \Sigma \Sigma (y_{ij} - \overline{y}_{i} - \overline{y}_{ij} + \overline{y}_{i})^2$$
 is the error or
 i_j
residuals

Degrees of fixedom for various s.s

$$S_{T}^{2} \text{ being computed in N=hk quantities(9ij-9.)}$$
Which are autient to one linear constraint $\sum_{i=j}^{r} (9ij-9..)=0$
Will carry (N-1)

$$S_{t}^{2} \longrightarrow (n-1) \longrightarrow \sum_{i}^{r} (9ij-9..)=0$$

$$S_{v}^{2} \longrightarrow (k-1) \longrightarrow \sum_{i}^{r} (9ij-9..)=0$$

$$S_{E}^{2} \longrightarrow (n-1)-(k-1)-(h-1) = (h-1)(k-1)$$
Thus the postitioning q dif is as follows
 $[h(k-1)] = (k-1)+(h-1)+(h-1)(k-1)$
Which implies that the dif are additive.
Test attaistic
Mean S.S due to

$$\frac{S_{t}^{2}}{k-1} = S_{t}^{2}$$
Mean s.s due to varieties. $= \frac{S_{v}^{2}}{k-1} = S_{v}^{2}$

Error Mean S.S =
$$\frac{S_E^2}{(h-1)(k-1)} = S_E^2$$

Over j from 1 to h and dividing by h, we get

$$\frac{1}{h} \sum_{j} y_{ij} = \frac{1}{h} \left[h_{M} + h\alpha_{i} + \sum_{j} B_{j} + \sum_{j} Z_{ij} \right] \Rightarrow$$

$$\overline{\mathcal{Y}_{i.}} = \mu + \alpha_j + \overline{z_{i.}}$$

Over i from 1 to 1 and dividing by k and using $\overline{z}\alpha_i = 0$ We get, $\overline{y}_{.j} = \mu + \beta_j + \overline{z}_{ij}$

Over i and j both and dividing by hk and using, We ashall get,

$$\overline{y}_{..} = \mu + \overline{\xi}_{..}$$

ANOVA TABLE FOR TWO-WAY CLASSIFLED

		1		
Sources of Variation			M-3.3	Valiance Ratio (F)
Factor A (ROWS)	S.S.A	(٢-١)	S.S.A =M.S.S.A 12-1	$F_{A} = \frac{S_{t}^{2}}{M \cdot s \cdot s \cdot F} - F(k-1)(k-1)(k-1)$ S_{t}^{2}
Factor B (COLUMINIS)	8.3°B	(h-1)	<u>S.s.B</u> h-1 -1	$F_{B} = \frac{M \cdot S \cdot S \cdot B}{M \cdot S \cdot S \cdot E} - F(h-1), (k-1)$ $\frac{M \cdot S \cdot S \cdot E}{S_{E}^{2}} \qquad (h-1)$
Error	S.S.E	(k-1) (h-1)	<u>S.s.F</u> = M.s.s.E (k-1)(h-1)	
Total	J.S.S	h12-)		

Conclusion :

 P_{f} $F_{t} \leq F_{(k-1),(h-1)(k-1)}$, we accept Hor

Othonvise,

If Fr & Fch-1), Ch-1) CK-1), We accept 1902 Otherwese

reject Hoz.