

UNIT

11

Simple Random Sampling

→

1. Define simple random sampling with replacement and without replacement from a finite population. Also explain various procedures of drawing samples.

The simplest and common used method of sampling is simple random sampling.

Simple random sampling is a technique of drawing a sample in such a way that each and every unit of the population has an equal and independent chance of being included in the sample.

To draw a sample of size n from the population of size N in simple random sampling we have two methods.

1. Simple Random sampling with replacement (SRSWR)
2. Simple Random sampling without replacement (SRSWOR)

Simple Random sampling with replacement (SRSWR)

The technique of drawing each and every unit of population has an equal and independent chance of being included into the sample in such a way that the selected unit should be replaced again into the population before going to select the next unit is called simple random sampling with replacement. Simply written as SRSWR.

Suppose we have to select a sample of size n from the population of size N in SRSWR we proceed as follows.

1. Probability of drawing a unit from the N population units at the first draw is $\frac{1}{N}$ and is same for all draws.

2. Number of possible samples in SRSWR are N^n
Probability of selecting a sample of size n from a population of size N in SRSWR is $\frac{1}{N^n}$

All the draws in SRSWR are independent and identical.

Simple Random Sampling without Replacement (SRSWOR)

The technique of drawing each and every unit of the population has an equal and independent chance of being included into the sample in such a way that the selected unit should not be replaced by any other subsequent draws. The next draw is taken into the population before the next draw or sampling without replacement. Simply written as SRSWOR.

Probability of drawing a unit from N units at the first draw is $\frac{1}{N}$

at the second draw is $\frac{1}{N-1}$ (since $N-1$ units remain)

at the third draw is $\frac{1}{N-2}$ and so on.

2. Number of possible samples in SRSWOR are $N C_n$

3. Probability of selecting a sample of size n from a population of size N in SRSWOR is $\frac{1}{N C_n}$

4. All draws in SRSWOR are independent but not identical.

• Explain selection of a Simple Random sample.

Procedures of selecting a random sample.

There are several methods to draw a sample of size n out of N units through simple random sampling.

Selection of a simple random sample depends on the size and nature of the population.

A simple random sample can be obtained by the following two methods.

1. Lottery method.

2. Mechanical: Randomisation or Random numbers method.

1. Lottery method: This is the simplest method of selecting a random sample.

These slips should be as homogeneous as possible in shape, size, colour etc. to avoid human bias. Put these in slips in a bag and thoroughly.

shuffled and then n slips are drawn one by one. The n units (items) corresponding to numbers on the slips drawn will constitute a random sample. The only advantage in this method is its simplicity. If the population is large, the lottery method is time consuming and practically difficult.

Random Numbers method of mechanical Randomization.
If the population is large, the lottery method is time consuming, expensive and very difficult to select a random sample.

The most inexpensive method of drawing a random sample is random numbers method.

In this method the random number tables have been identified from various records / blocks so that each of the digits 0, 1, 2, ... 9 appear with the same frequency and independent of each other.

If a random sample is to be selected from the population of size $(N) \leq 99$ then the numbers can be combined by two digits from 00 to 99.

Similarly if $N \leq 999$ then combine three digits from 000 to 999 and if $N \leq 9999$ then combine four digits from 0000 to 9999 and so on.

obtaining a random sample from the random numbers method can be explained in the following steps:

1. Identify the total number of units in the population N with the numbers from 1 to N .
2. select any page of random number tables and pick up the numbers in any row or column or diagonal at random.
3. The selected numbers then constitute the required random sample.

Some of the random number tables in common use are

1. Tippett's random number tables
2. Fisher & Yates tables

There are some modified procedures of selecting random numbers.

The method of selecting a sample with the help of random number table is always advisable.

3) Explain SRSWOR vs SRSWR

SRSWOR

SRSWR

1. If the unit selected in any draw is not replaced in the population before making the next draw then the sampling plan is called SRSWOR. If the unit selected in any draw is replaced back before making the next draw then the sampling plan is called SRSWR.

2. All draws are independent but not identical.

All draws are independent and identical.

3. Number of possible samples in SRSWOR are $N_c n$

Number of possible samples in SRSWR are N^n

4. The probability of selecting a sample of size n from a population of size N is $\frac{1}{N_c n}$

The probability of selecting a sample of size n from a population of size N is $\frac{1}{N^n}$

5. $E(\bar{x}) = \mu$

$E(\bar{x}) = \mu$

6. Sample mean square is an unbiased estimator of the population mean square

Sample mean square is an unbiased estimator of population variance i.e. $E(s^2) = \sigma^2$

i.e. $E(s^2) = \sigma^2$

7. Variance of the sample mean in SRSWOR is

Variance of the sample mean in SRSWR is $\frac{N-1}{N} \frac{\sigma^2}{n}$

$$\frac{N-n}{N} \frac{\sigma^2}{N}$$

Since variance is less efficiency is more

The efficiency of the sample mean is the lesser than in SRSWOR

8. Probability of selecting a unit at first draw is $\frac{1}{N}$ at second draw is $\frac{1}{N-1}$ and so on

9. Probability of selecting a unit in every draw is $\frac{1}{N}$ only.

SRSWR is suffered from the draw back having same unit two or more times in a sample.
In SRSWR we have distinct elements in a sample.
In many cases SRSWR is preferred to SRSWR.

4. Define merits and demerits of simple random sampling merits (advantages).

1. Since simple random sampling is probability sampling it eliminates personal bias, as such a simple random sample is more representative of the population as compared with the judgement or purposive sampling.
2. we can estimate the efficiency of the estimator through their standard error in simple random sampling method.

3. The calculations for estimating the population parameters are easy.

4. simple random sampling method specifies the three principles of a good sampling design.

5. simple random sampling method is more popular and widely used sampling method which reduces the cost, time and labour.

Demerits (Limitations) & Disadvantages.

1. To select a simple random sample, it requires an up-to-date frame i.e. list or guide of the population. This is not possible every time.
A simple random sampling suffers from this drawback severely.

2. The size of the sample in simple random sample should be more otherwise sample may not be a good representative of the population. Note that if the sample size increases, non-sampling errors also increase.

If the population is too large, simple random sampling method is not easy to apply.

4. If the population is heterogeneous the estimates by simple random sampling method are not reliable. Hence

Cases we prefer stratified Random sampling.

5. Simple Random sampling may give non-random looking samples in some cases.

6. In simple random sampling with replacement we may get the same unit more than once. So the number of distinct observations is reduced. In such cases there may be increase in sampling errors.

Notations and Terminology in Simple Random sampling.

Population

Sample

values of units Y_1, Y_2, \dots, Y_N

values of units y_1, y_2, \dots, y_n

size = N

size = n

Population mean = $\bar{Y}_N = \frac{1}{N} \sum Y_i$

Sample mean = $\bar{y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$

$$= \frac{1}{n} \sum_{i=1}^n a_i y_i$$

$a_i = \begin{cases} 1 & \text{if } i\text{th unit is included in the sample} \\ 0 & \text{if } i\text{th unit is not included in the sample} \end{cases}$

$$\text{variance} \sigma^2 = \frac{1}{N} \sum_{i=1}^N (Y_i - \bar{Y})^2$$

$$= \frac{1}{N} \left(\sum_{i=1}^N Y_i^2 - N\bar{Y}^2 \right)$$

$$s^2 = \sigma^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$= \frac{1}{n} \left(\sum_{i=1}^n y_i^2 - n\bar{y}^2 \right)$$

Population mean square

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2$$

$$= \frac{1}{N-1} \left[\sum_{i=1}^N Y_i^2 - N\bar{Y}^2 \right]$$

Sample mean square

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$= \frac{1}{n-1} \left(\sum_{i=1}^n y_i^2 - n\bar{y}^2 \right)$$

Let \therefore 1. Sample mean = $\bar{y}_n = \bar{y} = \frac{1}{n} \sum_{i=1}^n a_i y_i$

here $a_i = \begin{cases} 1 & \text{if } i\text{th unit is included in the sample} \\ 0 & \text{if } i\text{th unit is not included in the sample} \end{cases}$

$$E(a_i) = 1 \cdot P(a_i = 1) + 0 \cdot P(a_i = 0)$$

$E(a_i) = 1 \cdot P(i\text{th unit is included in the sample of size } n)$

$= 1 - P(i\text{th unit is not included in the sample of size } n)$

$$= P(a_i=1) \cdot P(a_j=1/a_{i-1})$$

$$= \frac{n}{N} \frac{n-1}{N-1}$$

$$E(a_i a_j) = \frac{n(n-1)}{N(N-1)}$$

$$4. \left(\sum_{i=1}^n y_i \right)^2 = \sum_{i=1}^n y_i^2 + \sum_{i \neq j=1}^n y_i y_j$$

Theorem 1 :- The probability that a specified unit of the population being selected at any given draw (say r th draw) is equal to the probability of it being selected at the first draw in Simple Random Sampling.

Proof :- Let A_i denote the event of selecting a specified unit at the i th draw $i=1, 2, \dots, r$. we are required to prove $P(A_r) = P(A_1)$

In SRSWR, clearly $P(A_r) = P(A_1) = \frac{1}{N}$

Where N is the size of the population.

In SRSWOR, an item is selected in the r th draw means that it is not selected in the previous $(r-1)$ draws

$$\begin{aligned} P(A_r) &= P(A_1' \cap A_2' \cap A_3' \cap \dots \cap A_r) \\ &= P(A_1') P(A_2' | A_1') P(A_3' | A_1' \cap A_2') \dots \\ &\quad \dots P(A_r | A_1' \cap A_2' \cap A_3' \dots \cap A_{r-1}') \end{aligned}$$

By the multiplication theorem of probability since draws are independent.

$$P(A_r) = P(A_1') P(A_2') P(A_3') \dots P(A_{r-1}') P(A_r)$$

$$P(A_r) = \left(1 - \frac{1}{N}\right) \left(1 - \frac{1}{N-1}\right) \left(1 - \frac{1}{N-2}\right) \dots \left(1 - \frac{1}{N-r+2}\right) \frac{1}{N-r+1}$$

$$P(A_r) = \left(\frac{N-1}{N}\right) \left(\frac{N-2}{N-1}\right) \left(\frac{N-3}{N-2}\right) \dots \frac{N-r+1}{N-r+2} \frac{1}{N-r+1}$$

$$P(A_r) = \frac{1}{N} = P(A_1)$$

$$P(A_r) = P(A_1) = \frac{1}{N}$$

Note :- The probability that a specified unit is selected in the r th draw is equal to the probability of it being selected in the first draw.

$$E(\bar{y}_n) = \bar{Y}_N$$

Proof :- Suppose the population consists of N units Y_1, Y_2, \dots, Y_N from which the sample of size y_1, y_2, \dots, y_n are drawn without replacement.

$$\text{Population mean } \bar{Y}_N = \frac{1}{N} \sum_{i=1}^N Y_i$$

$$\text{Sample mean } \bar{y}_n = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\bar{y}_n = \frac{1}{n} \sum_{i=1}^n a_i Y_i$$

Taking expectation on both sides

$$E(\bar{y}_n) = \frac{1}{n} \sum_{i=1}^n E(a_i) Y_i \quad \text{--- (1)}$$

We calculate $E(a_i)$, a_i takes the values

$$a_i = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ unit is included into the sample} \\ 0, & \text{if } i^{\text{th}} \text{ unit is not included into the sample} \end{cases}$$

$$E(a_i) = \sum_{i=0}^1 a_i P(a_i)$$

$$= 1 P(a_i=1) + 0 P(a_i=0)$$

$$E(a_i) = 1 P(i^{\text{th}} \text{ unit is included into the sample of size } n) + 0 P(i^{\text{th}} \text{ unit is not included into the sample of size } n)$$

$$E(a_i) = \frac{n}{N} \quad \text{--- (2)}$$

Substitute equation (2) in (1) we get

$$E(\bar{y}_n) = \frac{1}{n} \sum_{i=1}^n \frac{n}{N} Y_i$$

$$= \frac{1}{N} \sum_{i=1}^N Y_i = \bar{Y}_N$$

In SRSWOR sample mean is an unbiased estimate of population mean.

3) Prove that in simple Random Sampling with Replacement (SRSWR), sample mean (\bar{y}) is an unbiased estimate of the population mean (\bar{Y})

$$\text{i.e. } E(\bar{y}) = \bar{Y} \quad \& \quad E(\bar{y}_n) = \bar{Y}_N$$

If y_1, y_2, \dots, y_n units are drawn from Y_1, Y_2, \dots, Y_N units of population with the probability $\frac{1}{N}$ at every draw \therefore y_i gets from any one of

$\therefore E(y_i) = \bar{Y}_N \quad \forall i$ independent

$\therefore E(\bar{Y}_n) = E\left(\frac{1}{n} \sum_{i=1}^n y_i\right)$

$= \frac{1}{n} \sum_{i=1}^n E(y_i)$

$= \frac{1}{n} \sum_{i=1}^n \bar{Y}_N$

$= \frac{1}{n} n \bar{Y}_N = \bar{Y}_N$ [from eqn ①]

In SRSWR sample mean is an unbiased estimator of population mean.

Theorem ③ :- show that in simple Random Sampling without Replacement (SRSWOR) the sample mean square (sample variance s^2) is an unbiased estimator of the population mean square (S^2) i.e. $E(s^2) = S^2$

Proof :- Suppose the population consists of N units Y_1, Y_2, \dots, Y_N from which the sample of size n y_1, y_2, \dots, y_n are drawn without replacement

Sample mean $\bar{y}_n = \frac{\sum_{i=1}^n y_i}{n}$

We know that $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y}_n)^2$

$s^2 = \frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - n \bar{y}_n^2 \right]$

$= \frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - n \left(\frac{\sum_{i=1}^n y_i}{n} \right)^2 \right]$

$= \frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - \frac{n}{n^2} \left[\sum_{i=1}^n y_i \right]^2 \right]$

$= \frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - \frac{1}{n} \left[\sum_{i=1}^n y_i^2 + \sum_{i \neq j=1}^n y_i y_j \right] \right]$

$= \frac{1}{n-1} \sum_{i=1}^n y_i^2 - \frac{1}{n(n-1)} \sum_{i=1}^n y_i^2 - \frac{1}{n(n-1)} \sum_{i \neq j=1}^n y_i y_j$

[$\because \left(\sum_{i=1}^n y_i \right)^2 = \sum_{i=1}^n y_i^2 + \sum_{i \neq j=1}^n y_i y_j$]

$= \left(\frac{1}{n-1} - \frac{1}{n(n-1)} \right) \sum_{i=1}^n y_i^2 - \frac{1}{n(n-1)} \sum_{i \neq j=1}^n y_i y_j$

$= \frac{(n-1)}{n(n-1)} \sum_{i=1}^n y_i^2 - \frac{1}{n(n-1)} \sum_{i \neq j=1}^n y_i y_j$

$$= \frac{1}{n} \sum_{i=1}^n y_i^2 - \frac{1}{n(n-1)} \sum_{i \neq j=1}^n y_i y_j \rightarrow \textcircled{1}$$

Taking Expectation on both sides

$$E(S^2) = \frac{1}{n} E \left[\sum_{i=1}^n y_i^2 \right] - \frac{1}{n(n-1)} E \left[\sum_{i \neq j=1}^n y_i y_j \right] \rightarrow \textcircled{2}$$

$$\text{Consider } E \left[\sum_{i=1}^n y_i^2 \right] = E \left[\sum_{i=1}^N a_i y_i^2 \right]$$

where $a_i = \begin{cases} 1, & \text{if } i^{\text{th}} \text{ unit is included into the sample} \\ 0, & \text{if } i^{\text{th}} \text{ unit is not included into the sample} \end{cases}$

$$\begin{aligned} \therefore E \left[\sum_{i=1}^n y_i^2 \right] &= \sum_{i=1}^N E(a_i) y_i^2 \\ &= \sum_{i=1}^N \frac{n}{N} y_i^2 \\ &= \frac{n}{N} \sum_{i=1}^N y_i^2 \quad \text{--- } \textcircled{3} \end{aligned}$$

and now consider

$$\begin{aligned} E \left[\sum_{i \neq j=1}^n y_i y_j \right] &= E \left[\sum_{i \neq j=1}^N a_i a_j y_i y_j \right] \\ &= \sum_{i \neq j=1}^N E(a_i a_j) y_i y_j \end{aligned}$$

$$\begin{aligned} \text{Consider } E[a_i a_j] &= 1 P(a_i a_j = 1) + 0 P(a_i a_j = 0) \\ &= 1 P(a_i = 1 \cap a_j = 1) \\ &= P(a_i = 1) \cdot P(a_j = 1 / a_i = 1) \end{aligned}$$

$$E(a_i a_j) = \frac{n}{N} \frac{n-1}{N-1}$$

$$E \left(\sum_{i \neq j=1}^n y_i y_j \right) = \sum_{i \neq j=1}^N \frac{n(n-1)}{N(N-1)} y_i y_j$$

$$E \left(\sum_{i \neq j=1}^n y_i y_j \right) = \frac{n(n-1)}{N(N-1)} \sum_{i \neq j=1}^N y_i y_j \quad \text{--- } \textcircled{4}$$

Substitute the values of $\textcircled{3}$ and $\textcircled{4}$ in equation

we get

$$E(S^2) = \frac{1}{n} \frac{n}{N} \sum_{i=1}^N y_i^2 - \frac{1}{n(n-1)} \frac{n}{N} \frac{(n-1)}{(N-1)} \sum_{i \neq j=1}^N y_i y_j$$

$$E(S^2) = \frac{1}{N} \sum_{i=1}^N y_i^2 - \frac{1}{N(N-1)} \sum_{i \neq j=1}^N y_i y_j \quad \text{--- } \textcircled{5}$$

(as compared with eqn $\textcircled{1}$ and eqn $\textcircled{5}$ for small letters, just replacing capital letters we get the same as above)

Compare eqn $\textcircled{1}$ and $\textcircled{5}$

$$E(s^2) = S^2$$

∴ In SRSWOR the sample mean square is an unbiased estimate of population mean square (S^2)

Theorem (4) :- In SRSWOR variance of sample mean is given by $\text{var}(\bar{Y}_n) = \frac{N-n}{N} \frac{S^2}{n}$ or $(1-f) \frac{S^2}{n} = \left(\frac{1}{n} - \frac{1}{N}\right) S^2$

where $f = \frac{n}{N}$ = sampling fraction
where S^2 = population mean square

Proof :- Sample mean = $\bar{Y}_n = \frac{\sum y_i}{n}$

$$\text{we have } \text{var}(\bar{Y}_n) = E(\bar{Y}_n^2) - (E(\bar{Y}_n))^2$$

$$\text{var}(\bar{Y}_n) = E(\bar{Y}_n^2) - \bar{Y}_N^2 \rightarrow \textcircled{1} \quad (\because E(\bar{Y}_n) = \bar{Y}_N)$$

$$\bar{Y}_n = \frac{\sum_{i=1}^n y_i}{n}$$

Taking expectation and square on both sides we get

$$E(\bar{Y}_n^2) = E\left(\frac{1}{n} \sum_{i=1}^n y_i\right)^2$$

$$= \frac{1}{n^2} E\left[\sum_{i=1}^n y_i\right]^2$$

$$= \frac{1}{n^2} E\left[\sum_{i=1}^n y_i^2 + \sum_{i \neq j=1}^n y_i y_j\right]$$

$$\because \left(\sum_{i=1}^n y_i\right)^2 = \sum_{i=1}^n y_i^2 + \sum_{i \neq j=1}^n y_i y_j$$

$$E(\bar{Y}_n^2) = \frac{1}{n^2} \left[E\left(\sum_{i=1}^n y_i^2\right) + E\left(\sum_{i \neq j=1}^n y_i y_j\right) \right] \rightarrow \textcircled{2}$$

$$\text{consider } E\left[\sum_{i=1}^n y_i^2\right] = E\left[\sum_{i=1}^N a_i y_i^2\right] = \sum_{i=1}^N E(a_i) y_i^2$$

$$E\left[\sum_{i=1}^n y_i^2\right] = \sum_{i=1}^N \frac{n}{N} y_i^2 = \frac{n}{N} \sum_{i=1}^N y_i^2 \rightarrow \textcircled{3}$$

$$E(a_i) = \frac{n}{N}$$

$$\text{But } \sum_{i=1}^N (y_i - \bar{Y}_N)^2 = \sum_{i=1}^N y_i^2 - N \bar{Y}_N^2$$

$$\sum_{i=1}^N y_i^2 = \sum_{i=1}^N (y_i - \bar{Y}_N)^2 + N \bar{Y}_N^2$$

$$\sum_{i=1}^N y_i^2 = (N-1)S^2 + N \bar{Y}_N^2 \rightarrow \textcircled{4}$$

$$\because S^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y}_N)^2$$

substitute $\textcircled{4}$ and in $\textcircled{3}$ we get

$$E\left[\sum_{i=1}^n y_i^2\right] = \frac{n}{N} \left[(N-1)S^2 + N \bar{Y}_N^2 \right]$$

$$E\left[\sum_{i=1}^n y_i^2\right] = n\left[\left(\frac{N-1}{N}\right)S^2 + \bar{Y}_N^2\right] \rightarrow (5)$$

considered $E\left[\sum_{i \neq j=1}^n y_i y_j\right] = E\left[\sum_{i \neq j=1}^n a_i a_j y_i y_j\right]$

$$E\left[\sum_{i \neq j=1}^n y_i y_j\right] = \sum_{i \neq j=1}^n E(a_i a_j) y_i y_j$$

$$E\left[\sum_{i \neq j=1}^n y_i y_j\right] = \sum_{i \neq j=1}^n \frac{n(n-1)}{N(N-1)} y_i y_j$$

$$E\left[\sum_{i \neq j=1}^n y_i y_j\right] = \frac{n(n-1)}{N(N-1)} \sum_{i \neq j=1}^n y_i y_j \rightarrow (6)$$

$$\therefore \left(\sum_{i=1}^n y_i\right)^2 = \sum_{i=1}^n y_i^2 + \sum_{i \neq j=1}^n y_i y_j$$

$$\sum_{i \neq j=1}^n y_i y_j = \left(\sum_{i=1}^n y_i\right)^2 - \sum_{i=1}^n y_i^2$$

$$\sum_{i \neq j=1}^n y_i y_j = (N\bar{Y}_N)^2 - (N-1)S^2 - N\bar{Y}_N^2 \quad [\text{from (4)}]$$

$$\sum_{i \neq j=1}^n y_i y_j = N\bar{Y}_N^2 [N-1] - (N-1)S^2$$

$$= (N-1) [N\bar{Y}_N^2 - S^2] \rightarrow (7)$$

$$= \frac{N(N-1)}{N} \left(\bar{Y}_N^2 - \frac{S^2}{N}\right) \rightarrow (7)$$

Substitute eqn (7) in eqn (6)

$$E\left[\sum_{i \neq j=1}^n y_i y_j\right] = \frac{n(n-1)}{N(N-1)} \left[N(N-1) \left(\bar{Y}_N^2 - \frac{S^2}{N}\right)\right]$$

$$E\left[\sum_{i \neq j=1}^n y_i y_j\right] = n(n-1) \left[\bar{Y}_N^2 - \frac{S^2}{N}\right] \rightarrow (8)$$

Substitute eqn (5) and eqn (8) in eqn (2)

$$E(\bar{Y}_n^2) = \frac{1}{n^2} \left[E\left(\sum_{i=1}^n y_i^2\right) + E\left(\sum_{i \neq j=1}^n y_i y_j\right) \right]$$

$$E(\bar{Y}_n^2) = \frac{1}{n^2} \left[n \left(\frac{N-1}{N} S^2 + \bar{Y}_N^2\right) + n(n-1) \left(\bar{Y}_N^2 - \frac{S^2}{N}\right) \right]$$

$$E(\bar{Y}_n^2) = \frac{1}{n} \left[\frac{N-1}{N} S^2 + \bar{Y}_N^2 + (n-1) \bar{Y}_N^2 - (n-1) \frac{S^2}{N} \right]$$

$$= \frac{1}{n} \left[\left(\frac{N-1-n+1}{N}\right) S^2 + (1+n-1) \bar{Y}_N^2 \right]$$

$$= \frac{1}{n} \left(\frac{N-n}{N} S^2 + n \bar{Y}_N^2 \right)$$

$$E(\bar{Y}_n^2) = \frac{N-n}{Nn} S^2 + \bar{Y}_N^2 \rightarrow (9)$$

Substitute the value of eqn (9) in eqn (1) we get

$$\begin{aligned} \text{Var}(\bar{y}_n) &= C(\bar{y}_n^2) - \bar{y}_n^2 \\ &= \frac{N-n}{N} \frac{s^2}{n} + \bar{y}_n^2 - \bar{y}_n^2 \end{aligned}$$

$$\text{Var}(\bar{y}_n) = \frac{N-n}{N} \frac{s^2}{n} = \left(1 - \frac{n}{N}\right) s^2 = (1-F) s^2$$

$$\text{Var}(\bar{y}_n) = \left(\frac{1}{n} - \frac{1}{N}\right) s^2$$

Hence the proof.

Standard Error of the sample mean is SRSWOR

$$\text{SE}(\bar{y}_n) = \sqrt{V(\bar{y}_n)} = \sqrt{\frac{N-n}{N}} \frac{s}{\sqrt{n}}$$

ESTIMATION OF SE OF SAMPLE MEAN IN SRSWOR

If s^2 is not known, then we estimate s^2 by the sample mean square s^2

Since s^2 is an unbiased estimator of s^2

$$\text{Estimate of SE}(\bar{y}_n) = \sqrt{\frac{N-n}{N}} \frac{s}{\sqrt{n}}$$

Sampling fraction

$\frac{n}{N}$ is called sampling fraction and is denoted by f

$$\therefore f = \frac{n}{N}$$

finite population correlation (F.P.C)

$(1-f)$ is called finite population correlation and is usually written as f.p.c

$$\text{i.e. f.p.c} = 1-f$$

Note (1) If population is very large, then

$$f = \frac{n}{N} \rightarrow 0 \text{ and the f.p.c} \rightarrow 1$$

$$V(\bar{y}_n) = \frac{s^2}{n} \left(\because V(\bar{y}_n) = (1-f) \frac{s^2}{n} = (1-0) \frac{s^2}{n} = \frac{s^2}{n} \right)$$

$$(2) \text{ since } \hat{y} = N\bar{y}$$

$$V(\hat{y}) = V(N\bar{y})$$

$$= N^2 V(\bar{y}) = \frac{N^2(N-n)s^2}{Nn}$$

$$= N(N-n) \frac{s^2}{n}$$

$$= \frac{N^2 s^2}{n}$$

$$\therefore \frac{N-n}{N} = 1$$

$$\therefore SE(\bar{y}) = \frac{NS}{\sqrt{n}}$$

Theorem (5) Prove that in SRSWR variance of the sample mean is given by $V(\bar{y}) = \frac{\sigma^2}{n}$ or $\frac{N-1}{N} \frac{S^2}{n}$ where σ^2, n are population variance and sample size.

Proof:- Let the sample observations y_1, y_2, \dots, y_n be independent and identically drawn from the population with the same variance σ^2

$$\text{i.e. } \text{var}(y_i) = \sigma^2 \quad \forall i$$

$$\text{var}(\bar{y}_n) = \text{var}\left[\frac{1}{n} \sum_{i=1}^n y_i\right]$$

$$= \frac{1}{n^2} \sum_{i=1}^n V(y_i)$$

(covariance terms vanish since they are independent).

$$= \frac{1}{n^2} \sum_{i=1}^n \sigma^2$$

$$= \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

$$\text{i.e. } \text{var}(\bar{y}_n) = \frac{\sigma^2}{n} = \frac{N-1}{N} \frac{S^2}{n} \quad (\because (N-1)S^2 = N\sigma^2)$$

If N is large $\frac{1}{N} \rightarrow 0$ $\left(\frac{N-1}{N} = 1 - \frac{1}{N} = 1 - 0 = 1\right)$

$$\therefore V(\bar{y}_n) = \frac{\sigma^2}{n}$$

This is same as $V(\bar{y}_n)$ in SRSWOR

Note: (1) Comparison of variance in SRSWOR and SRSWR variance of the sample mean in SRSWOR is less than that in SRSWR

we know $\text{var}(\bar{y}_n) = \frac{N-n}{N} \frac{S^2}{n}$ in SRSWOR

$$\text{var}(\bar{y}_n) = \frac{N-1}{N} \frac{S^2}{n} \text{ in SRSWR}$$

Since sample size n is always greater than 1

$$N-n < N-1$$

$$\therefore \text{var}(\bar{y}_n)_{\text{SRSWOR}} < \text{var}(\bar{y}_n)_{\text{SRSWR}}$$

Theorem 6) Second method

Proof :- By definition variance of the sample mean \bar{y}

$$\begin{aligned} V(\bar{y}) &= E(\bar{y} - E(\bar{y}))^2 \\ &= E(\bar{y} - \bar{y})^2 \quad [\because E(\bar{y}) = \bar{y}] \\ &= E\left[\frac{\sum_{i=1}^n y_i}{n} - \bar{y}\right]^2 \\ &= E\left[\frac{\sum_{i=1}^n y_i - n\bar{y}}{n}\right]^2 \\ &= \frac{1}{n^2} E\left[\sum_{i=1}^n y_i - n\bar{y}\right]^2 = \frac{1}{n^2} E\left[\sum_{i=1}^n (y_i - \bar{y})\right]^2 \end{aligned}$$

$$\Rightarrow V(\bar{y}) = \frac{1}{n^2} E\left[\sum_{i=1}^n (y_i - \bar{y})^2 + \sum_{\substack{i=1 \\ i < j}}^n \sum_{j=1}^n (y_i - \bar{y})(y_j - \bar{y})\right]$$

$$\left(\because \left(\sum_{i=1}^n y_i\right)^2 = \sum_{i=1}^n y_i^2 + \sum_{i \neq j=1}^n y_i y_j\right)$$

$$V(\bar{y}) = \frac{1}{n^2} \left[\sum_{i=1}^n E(y_i - \bar{y})^2 + \sum_{i=1}^n \sum_{j=1}^n E(y_i - \bar{y})(y_j - \bar{y}) \right] \rightarrow \textcircled{1}$$

$$\left. \begin{aligned} E(y_i - \bar{y})^2 &= V(y_i) = \sigma^2 \\ E(y_i - \bar{y})(y_j - \bar{y}) &= \text{cov}(y_i, y_j) = 0 \end{aligned} \right\} \textcircled{2}$$

In SRSWR the units are drawn independent

$\therefore y_i$ and y_j are independent

Hence $\text{cov}(y_i, y_j) = 0$

Substitute $\textcircled{2}$ in $\textcircled{1}$

$$V(\bar{y}) = \frac{1}{n^2} [\sum \sigma^2 + 0] = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$ This is called the true variance

$$N\sigma^2 = \sum_{i=1}^N (y_i - \bar{y})^2$$

But the estimated variance $S^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2$

$$\Rightarrow (N-1)S^2 = \sum_{i=1}^N (y_i - \bar{y})^2$$

$$(N-1)S^2 = N\sigma^2$$

$$\sigma^2 = \frac{N-1}{N} S^2$$

$$V(\bar{y}) = \frac{\sigma^2}{n} = \frac{N-1}{N} \frac{S^2}{n}$$

Theorem 6:— In SRSWR the sample mean square ($s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$) is an unbiased estimator of the population variance (σ^2) i.e. $E(s^2) = \sigma^2$

Proof:— As we know in SRSWR

$$\left. \begin{aligned} E(y_i) &= \bar{Y}_N & \text{var}(y_i) &= \sigma^2 \\ E(\bar{y}_n) &= \bar{Y}_N & \text{var}(\bar{y}_n) &= \frac{\sigma^2}{n} \end{aligned} \right\} \forall i \quad \text{--- (1)}$$

$$\begin{aligned} \text{Now } E(s^2) &= E\left[\frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - n\bar{y}_n^2 \right]\right] \\ &= \frac{1}{n-1} \left[\sum_{i=1}^n E(y_i^2) - nE(\bar{y}_n^2) \right] \quad \text{--- (2)} \end{aligned}$$

We know that $\text{var}(y_i) = E(y_i^2) - (E(y_i))^2$ (from eqn (1))

$$E(y_i^2) = \text{var}(y_i) + (E(y_i))^2 = \sigma^2 + \bar{Y}_N^2 \quad \text{--- (3)}$$

$$\text{and } E(\bar{y}_n^2) = \text{var}(\bar{y}_n) + (E(\bar{y}_n))^2 = \frac{\sigma^2}{n} + \bar{Y}_N^2 \quad \text{--- (4)}$$

Sub the values in eqn (3) and (4) in eqn (2) from (1) we get

$$E(s^2) = \frac{1}{n-1} \left[\sum_{i=1}^n (\sigma^2 + \bar{Y}_N^2) - n \left(\frac{\sigma^2}{n} + \bar{Y}_N^2 \right) \right]$$

$$E(s^2) = \frac{1}{n-1} \left[n\sigma^2 + n\bar{Y}_N^2 - n\frac{\sigma^2}{n} - n\bar{Y}_N^2 \right]$$

$$E(s^2) = \frac{1}{n-1} (n-1)\sigma^2$$

$$E(s^2) = \sigma^2 \quad \text{where } s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

Problem

1. Draw a sample of size 2 from the population consisting of 3, 5, 7 by SRSWR and SRSWOR methods compare the two methods through your results.

Sol Given population units are 3, 5, 7

Population size $N=3$

Sample size $n=2$

$$\text{Population mean } \bar{y} = \frac{\sum y}{N} = \frac{3+5+7}{3}$$

$$= \frac{15}{3} = 5$$

Population mean square

$$s^2 = \frac{1}{N-1} \sum (y_i - \bar{y})^2 = \frac{1}{3-1} \left[(3-5)^2 + (5-5)^2 + (7-5)^2 \right]$$

$$s^2 = \frac{1}{2} [4+0+4] = \frac{8}{2} = 4$$

$$\sigma^2 = \text{Population variance} = \frac{1}{N} \sum (y_i - \bar{y})^2$$

$$\sigma^2 = \frac{1}{3} [(3-5)^2 + (5-5)^2 + (7-5)^2] = \frac{1}{3} [4+0+4] = 8$$

no	sample values	sample mean \bar{y}_n	y_n	deviation $s_i^2 = \frac{1}{n-1} \sum (y_i - \bar{y}_i)^2$
1	(3, 3)	3	9	$s_1^2 = \frac{1}{2-1} [(3-3)^2 + (3-3)^2] = 0$
2	(3, 5)	4	16	$s_2^2 = \frac{1}{2-1} [(3-4)^2 + (5-4)^2] = 2$
3	(3, 7)	5	25	$s_3^2 = 8$
4	(5, 3)	4	16	$s_4^2 = 2$
5	(5, 5)	5	25	$s_5^2 = 0$
6	(5, 7)	6	36	$s_6^2 = 2$
7	(7, 3)	5	25	$s_7^2 = 8$
8	(7, 5)	6	36	$s_8^2 = 2$
9	(7, 7)	7	49	$s_9^2 = 0$
		<u>45</u>	<u>237</u>	<u>24</u>

$$\sum \bar{y}_n = 45 \quad \sum \bar{y}_n^2 = 237 \quad \sum s_i^2 = 24$$

$$v(\bar{y}_n) = \frac{\sum \bar{y}_n^2}{N^n} - \left(\frac{\sum \bar{y}_n}{N^n} \right)^2 = \frac{237}{9} - \left(\frac{45}{9} \right)^2 = 26.3333 - 25 = 1.3333$$

$$v(\bar{y}_n)_{SRSWR} = \frac{\sigma^2}{n} = \frac{2.6667}{2} = 1.3333 \quad \text{The formula is verified.}$$

$$E(s^2) = \frac{1}{N^n} \sum_{i=1}^{N^n} s_i^2 = \frac{24}{9} = 2.6667 = \sigma^2$$

$$E(s^4) = \sigma^4$$

$$E(\bar{y}_n) = \frac{\sum \bar{y}_n}{N^n} = \frac{45}{9} = 5 = \bar{Y}_N$$

$$\text{ie } E(\bar{y}_n) = \bar{Y}_N$$

SRSWOR possible samples = $Nc_n = 3c_2 = 3$

Sample no	sample values	\bar{y}_n	\bar{y}_n^2	$\sum x_i^2$
1	(3,5)	4	16	2
2	(3,7)	5	25	8
3	(5,7)	6	36	2
		<u>15</u>	<u>77</u>	<u>12</u>

$$\sum \bar{y}_n = 15 \quad \sum \bar{y}_n^2 = 77 \quad \sum x_i^2 = 12 \quad Nc_n = 3$$

$$E(\bar{y}_n) = \frac{\sum \bar{y}_n}{Nc_n} = \frac{15}{3} = 5 = \bar{Y}_N$$

$$E(\bar{y}_n) = \bar{Y}_N$$

$$E(s^2) = \frac{\sum x_i^2}{Nc_n}$$

$$E(s^2) = \frac{12}{3} = 4$$

$$E(s^2) = s^2$$

$$V(\bar{y}_n) = \frac{\sum \bar{y}_n^2}{Nc_n} - \left(\frac{\sum \bar{y}_n}{Nc_n} \right)^2$$

$$= \frac{77}{3} - \left(\frac{15}{3} \right)^2$$

$$= 25.6667 - 25$$

$$= 0.6667$$

$$V(\bar{y}_n)_{\text{SRSWOR}} = \frac{N-n}{Nn} s^2 = \frac{3-2}{3 \times 2} \times 4$$

$$= \frac{4}{6} = \frac{2}{3} = 0.6667$$

The formula is verified.

Estimation of population mean, population total and variance of these estimators by SRSWR and SRSWOR
 (a) Estimation of population parameters.

Estimation of population mean, total and proportion is same in with and without replacement techniques whereas it will be different for population variance.

① Estimation of population mean

In simple Random sampling without replacement

$\frac{N-1}{N} \sigma^2$ is an unbiased estimator of the population

Variance in SRSWOR

i.e. Population Variance is estimated by $\frac{N-1}{N} s^2$

$$\text{i.e. } \hat{\sigma}^2 = \frac{N-1}{N} s^2$$

Estimation of Population Variance in SRSWOR

In SRSWOR we know that the sample mean square is an unbiased estimator of the population variance i.e. $E(s^2) = \sigma^2$

Therefore the population variance is estimated by sample mean square in SRSWOR i.e. $\hat{\sigma}^2 = s^2$

Question bank :-

1. Define simple random sampling. Write its merits and demerits.
2. Distinguish between SRSWR and SRSWOR
3. Explain the selection procedure of simple Random sampling (what is simple random sampling? mention the one of the method of drawing a random sample.)
Explain (a) Lottery method (b) Random numbers method in selecting SRS (explain methods of selecting simple random sample).

- In SRSWOR show that $E(\bar{Y}_n) = \bar{Y}_N$

i. In SRSWOR show that $E(s^2) = \sigma^2$

ii. In SRSWOR obtain $V(\bar{Y}_n) = \frac{N-n}{N} \frac{\sigma^2}{n}$

Show that the sample mean is an unbiased estimator of population mean in SRSWOR and find its variance in SRSWOR $V(\bar{Y}) = \frac{\sigma^2}{n}$

Define SRS. Show that sample mean is an unbiased estimator of population mean in SRSWOR

Show that sample mean square is an unbiased estimator of population mean square in SRSWOR.

Stratified Random Sampling

① Explain the need for stratification of the population in sample surveys and the procedure of stratified sampling.

All the methods of sampling, the most commonly used procedure in surveys is stratified random sampling when the population is heterogeneous with respect to the variable or characteristic under study then the technique of stratified random sampling is used to obtain more efficient results.

If the population is heterogeneous, then we divide the entire population into relatively homogeneous sub-groups called strata.

In every stratum (sub-group) we apply simple random sampling so that all the different groups should be represented proportionately equal.

Definition :- stratification means division into layers. Auxiliary information ie past data or some other information related to the population characteristic under study may be used to divide the population into various groups such that

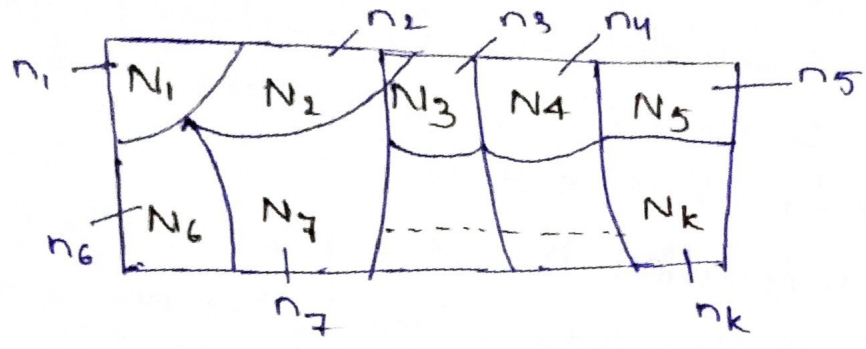
- i) units within each group are as homogeneous as possible
 - ii) The group means are as heterogeneous as possible.
- The population consist of N units is divided into k relatively homogeneous and mutually disjoint (means, non-overlapping) sub-groups are known as strata of sizes N_1, N_2, \dots, N_k such that $\sum_{i=1}^k N_i = N$.

From each stratum, a sample is drawn using simple random sampling without replacement (in general) of size n_i ($i=1, 2, \dots, k$) respectively such that $\sum_{i=1}^k n_i = n$, the sample is known as stratified random sample of size n and the technique of drawing such sample is called stratified random sampling.

The criterion used for the stratification of the universe into various strata is known as stratifying

ii) Selecting units at random from each stratum.

Diagrammatic representation of stratified Random sampling



The stratified random sampling is to be conducted keeping in mind that the population is to be divided into strata properly and a suitable sample size is to be selected from each stratum. Otherwise stratified random sample may not be trustworthy.

Example :- To study the cost of living of the people living in the state AP. In this case, the sampling data is to be collected from all income groups of people. Since population units are heterogeneous, one cannot assure that all the income groups like high income group, middle income group, low income group etc. are equally represented into the sample of the sampling is done through the simple random sampling. This is the by the technique stratified random.

Advantages of stratified Random sampling

• More Representative

Usually some sampling techniques like simple random sampling, some groups may be over represented, some may under represented and some may not represented at all. Stratified random sampling technique will be useful to represent all the sub-groups (strata) in the population. Hence stratified random sampling provide more representative cross-section of the population data.

• Greater Accuracy :-

Stratified random sampling provide more accurate

increasing precision since it enables us to obtain the precision for each stratum.

3. Administrative Convenience

In stratified random sampling, the samples will be concentrated more geographically as compared to simple random sampling.

Hence the time and money involved in the collection of the data and interviewing the individuals may be reduced to a great extent.

The supervision of field work also may be done with greater ease and convenience.

4. Sometimes the sampling problems will be different in different strata.

For example, the population may consist of literates and illiterates or people living in hostels, hospitals, jails and those living in ordinary homes.

In such situations, each of the above cases can be considered as a stratum.

This is a special advantage of stratified random sampling.

DISADVANTAGES

1. Stratified Random sampling may give highly biased estimates.

2. The selection of a stratified random sampling requires an up-to-date sampling frame. In practice up-to-date sampling frame is not available then it is impossible to identify the sampling units. Hence stratified Random sampling cannot be used.

Notations :-

Let N be the size of population

Let the population is divided into k strata of sizes N_1, N_2, \dots, N_k so that $N = \sum_{i=1}^k N_i$

Let us suppose that a simple random sample of size n_i (using SRSWOR) is drawn from the stratum N_i , $i = 1, 2, \dots, k$ such that $n = \sum_{i=1}^k n_i$

N = Total number of population units

N_i = The number of population units in the i^{th} stratum

n_i = The number of units selected in RRSWOR from the i^{th} stratum

$n = \sum_{i=1}^k n_i$: Total sample size drawn from all strata

Let Y_{ij} ($i=1, 2, \dots, k, j=1, 2, \dots, N_i$) be the value of the j^{th} unit in the i^{th} stratum.

Population mean of i^{th} stratum = $\bar{Y}_{N_i} = \frac{1}{N_i} \sum_{j=1}^{N_i} Y_{ij}$

Population mean = $\bar{Y}_N = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{N_i} Y_{ij}$

$$= \frac{1}{N} \sum_{i=1}^k N_i \bar{Y}_{N_i} \quad \left(\sum_{j=1}^{N_i} Y_{ij} = N_i \bar{Y}_{N_i} \right)$$

$$\bar{Y}_N = \sum_{i=1}^k P_i \bar{Y}_{N_i}$$

Where $P_i = \frac{N_i}{N}$ is called the weight of i^{th} stratum.

-m.

Population mean square of the i^{th} stratum

$$S_i^2 = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (Y_{ij} - \bar{Y}_{N_i})^2 \quad i=1, 2, \dots, k$$

Population mean square = $S^2 = \frac{1}{N-1} \sum_{i=1}^k \sum_{j=1}^{N_i} (Y_{ij} - \bar{Y}_N)^2$

Let y_{ij} be the value of j^{th} sampled unit drawn from the stratum.

Mean of the sample selected from the i^{th} stratum

Mean of the sample selected from the i^{th} stratum

$$\bar{y}_{n_i} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}, \quad i=1, 2, \dots, k$$

Mean of the stratified sample

$$\bar{y}_n = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij} \quad (n_i \bar{y}_{n_i} = \sum_{j=1}^{n_i} y_{ij})$$

$$\text{or } \bar{y}_n = \frac{1}{n} \sum_{i=1}^k n_i \bar{y}_{n_i}$$

Sample mean square of the i^{th} stratum

$$s_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{n_i})^2 \quad i=1, 2, \dots, k$$

Sample mean square

$$s^2 = \frac{1}{n-1} \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_n)^2$$

We consider two estimates of the population mean

\bar{Y}_N which will be given below

$$\bar{y}_m = \frac{1}{n} \sum_{i=1}^k n_i \bar{y}_{n_i}, \text{ the stratified sample mean}$$

If $\frac{n_i}{n} = \frac{N_i}{N}$, then stratified sample mean.

considered as \bar{y}_{st} and is given below

$$\bar{y}_{st} = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_{n_i}$$

$$= \sum_{i=1}^k p_i \bar{y}_{n_i}$$

where p_i is weight of the i^{th} stratum

Note :- $\bar{y}_n = \frac{1}{n} \sum_{i=1}^k n_i \bar{y}_{n_i}$

$$\bar{y}_{st} = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_{n_i}$$

These two estimates of the population mean are identical if $\frac{n_i}{n} = \frac{N_i}{N}$

$$\frac{n_i}{N_i} = \frac{n}{N} \text{ (constant)}$$

$$n_i = \frac{n}{N} N_i$$

$$n_i = c N_i$$

$$n_i \propto N_i$$

This is known as proportional allocation.

Theorem 1: show that in stratified random sampling the mean of the stratified random sample \bar{y}_{st} is an unbiased estimator of the population mean \bar{y}_N i.e.

$$E(\bar{y}_{st}) = \bar{y}_N$$

Proof: We know that in SRSWOR $E(\bar{y}_n) = \bar{y}_N$

Since the sample of size n_i ($i=1, 2, \dots, k$) is drawn using SRSWOR from each of the stratum,

We know $E(\bar{y}_{n_i}) = \bar{y}_{N_i}$ — (1) $i=1, 2, \dots, k$

Consider the stratified random sample mean

$$\bar{y}_{st} = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_{n_i}$$

where \bar{y}_{n_i} is the mean of the sample drawn from the i^{th} stratum.

Now consider

$$E(\bar{y}_{st}) = E\left[\frac{1}{N} \sum_{i=1}^k N_i \bar{y}_{n_i}\right]$$

$$= \frac{1}{N} \sum_{i=1}^k N_i E(\bar{y}_{n_i})$$

$$= \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_{N_i} \quad (\because \text{from (1)})$$

$$= \bar{y}_N \quad (\because \bar{y}_N = \frac{1}{N} \sum_{i=1}^k N_i \bar{y}_{N_i})$$

$\therefore E(\bar{y}_{st}) = \bar{Y}_N$ in stratified Random sampling

Theorem 2 show that stratified Random sampling is more precise than simple random sampling of the sample mean (estimate of population mean) i.e.

$$V(\bar{y}_{st}) = \frac{1}{N^2} \sum_{i=1}^k N_i (N_i - n_i) \frac{s_i^2}{n_i} \quad \text{where } s_i^2 = \frac{1}{N_i} \sum_{j=1}^{N_i} (y_{ij} - \bar{y}_i)^2$$

$$= \sum_{i=1}^k \left(\frac{1}{n_i} - \frac{1}{N_i} \right) (P_i s_i)^2$$

Proof :- We know that in SRSWOR

$$V(\bar{y}_n) = \frac{N-n}{N} \frac{s^2}{n}$$

since the samples have drawn from each stratum by using SRSWOR we have

$$V(\bar{y}_{ni}) = \frac{N_i - n_i}{N_i} \frac{s_i^2}{n_i} \quad i=1, 2, \dots, k \quad \text{--- (1)}$$

where N_i = size of i^{th} stratum
 n_i = the sample size drawn from the i^{th} stratum by using SRSWOR
 s_i^2 = population mean square of i^{th} stratum

Now consider

$$V(\bar{y}_{st}) = V\left(\frac{1}{N} \sum_{i=1}^k N_i \bar{y}_{ni}\right)$$

$$= \frac{1}{N^2} \sum_{i=1}^k N_i^2 V(\bar{y}_{ni})$$

$$= \frac{1}{N^2} \sum_{i=1}^k N_i \left[\frac{N_i - n_i}{N_i} \frac{s_i^2}{n_i} \right] \quad \text{[from eqn (1)]}$$

$$V(\bar{y}_{st}) = \frac{1}{N^2} \sum_{i=1}^k N_i (N_i - n_i) \frac{s_i^2}{n_i}$$

$$\text{or } V(\bar{y}_{st}) = \frac{1}{N^2} \sum_{i=1}^k N_i^2 \left(\frac{1}{n_i} - \frac{1}{N_i} \right) s_i^2$$

$$= \sum_{i=1}^k \left(\frac{1}{n_i} - \frac{1}{N_i} \right) P_i^2 s_i^2 \quad \text{where } P_i = \frac{N_i}{N}$$

$$V(\bar{y}_{st}) = \sum_{i=1}^k \left(\frac{1}{n_i} - \frac{1}{N_i} \right) (P_i s_i)^2$$

Note 1 $V(\bar{y}_{st}) = \sum_{i=1}^k \frac{N_i^2}{N^2} \left(\frac{1}{n_i} - \frac{1}{N_i} \right) s_i^2$

$$V(\bar{y}_{st}) = \sum_{i=1}^k P_i^2 s_i^2 \left(\frac{1}{n_i} - \frac{1}{N_i} \right)$$

$$V(\bar{y}_{st}) = \sum_{i=1}^k \frac{P_i^2 s_i^2}{n_i} - \sum_{i=1}^k \frac{P_i^2 s_i^2}{N_i}$$

$$\textcircled{2} \quad V(\hat{Y}_{st}) = V(N \bar{y}_{st}) = N^2 V(\bar{y}_{st})$$

$$= N^2 \sum_{i=1}^k \frac{N_i}{N^2} (N_i - n_i) \frac{s_i^2}{n_i}$$

$$V(\hat{Y}_{st}) = \sum_{i=1}^k N_i (N_i - n_i) \frac{s_i^2}{n_i}$$

$$= \sum_{i=1}^k P_i^2 (1 - F_i) \frac{S_i^2}{n_i} = \sum_{i=1}^k \frac{P_i^2 S_i^2}{n_i}$$

Explain the methods used for fixing the number of units to be selected from each stratum

(Q1)

Allocation of sample size of different strata

In stratified sampling, the allocation of the sample size to different strata is done by the consideration of three factors.

- i. stratum size
- ii. The variability within the stratum
- iii. The cost in taking sampling unit in the stratum.

A good allocation method is one which maximises the precision of the estimate with minimum resources.

To draw a sample size from the population, first of all we have to decide the size of the sub samples to be drawn from each strata.

In sample survey there are several types of allocation methods will be adopted, among them most commonly used methods are.

1. Proportional allocation
2. Neymann allocation or optimum allocation.

1. Proportional allocation (Bowley allocation)

This method of allocation is proposed by Bowley (1926)

This procedure of allocation is very common in practice, because of its simplicity, when no other information except stratum sizes is available we use this method.

Allocation of n_i 's to various strata is called proportional if the sample fraction is constant for each stratum.

$$\text{ie } \frac{n_1}{N_1} = \frac{n_2}{N_2} = \dots = \frac{n_k}{N_k} = \frac{\sum_{i=1}^k n_i}{\sum_{i=1}^k N_i} = \frac{n}{N} = c$$

$$\text{ie } \frac{n_i}{N_i} = \frac{n}{N} = c$$

$$n_i = \frac{n}{N} N_i$$

$$n_i = c N_i$$

from the various - - - - -

of optimization.

ie obtaining best results at minimum possible cost.

The rules given by Neyman for an optimum allocation n_i 's ($i=1,2,\dots,k$) are determined so that

- i) Minimize the variance for fixed sample size n
ie minimize $\text{var}(\bar{y}_{st})$ for fixed n .
- ii) minimize the variance (ie, maximizing the precision)
for fixed total cost c
ie $\text{var}(\bar{y}_{st})$ is minimum for fixed total cost c .
- iii) minimize the total cost for fixed value of $V(\bar{y}_{st})$
(desired precision).

under the optimum allocation

$$n_i = n \frac{N_i s_i}{\sum_{i=1}^k N_i s_i}$$

$$n_i \propto N_i s_i \quad \left(C = \sum N_i s_i = \text{constant} \right)$$

Here n is the total sample size

n_i is the i^{th} sub sample size

N_i is the i^{th} stratum size

s_i^2 is the i^{th} stratum mean square deviation

$$\text{and } s_i^2 = \frac{\sum_{j=1}^{N_i} (y_{ij} - \bar{y}_i)^2}{N_i - 1}$$

Theorem 4 :- With the usual notations,

$$\text{var}(\bar{Y}_{st})_{opt} \leq \text{var}(\bar{Y}_{st})_{prop} \leq \text{var}(\bar{Y}_n)_{ran}$$

$$\text{or } v_{opt} \leq v_{prop} \leq v_{ran}$$

Comparison of Stratified Random Sampling with simple Random sampling

$$\text{i.e. } \text{var}(\bar{Y}_n)_{ran} \geq \text{var}(\bar{Y}_{st})_{prop} > \text{var}(\bar{Y}_{st})_{opt}$$

Proof :- we have variance of the estimate of population mean in different methods as

$$v(\bar{Y})_{ran} = \frac{N-n}{N} \frac{S^2}{n} = \left(\frac{1}{n} - \frac{1}{N}\right) S^2$$

$$v(\bar{Y}_{st})_{prop} = \frac{N-n}{N} \frac{1}{n} \sum_{i=1}^k P_i S_i^2$$

$$= \left(\frac{1}{n} - \frac{1}{N}\right) \sum_{i=1}^k P_i S_i^2$$

$$v(\bar{Y}_{st})_{opt} = \frac{1}{n} \left(\sum_{i=1}^k P_i S_i\right)^2 - \frac{1}{N} \sum_{i=1}^k P_i S_i^2$$

$$\text{Consider } v_{prop} - v_{opt} = \left(\frac{1}{n} - \frac{1}{N}\right) \sum_{i=1}^k P_i S_i^2 - \frac{1}{n} \left(\sum_{i=1}^k P_i S_i\right)^2 + \frac{1}{N} \sum_{i=1}^k P_i S_i^2$$

$$= \frac{1}{n} \sum_{i=1}^k P_i S_i^2 - \frac{1}{N} \sum_{i=1}^k P_i S_i^2 - \frac{1}{n} \left(\sum_{i=1}^k P_i S_i\right)^2 + \frac{1}{N} \sum_{i=1}^k P_i S_i^2$$

$$= \frac{1}{n} \left[\sum_{i=1}^k P_i S_i^2 - \left(\sum_{i=1}^k P_i S_i\right)^2 \right]$$

$$= \frac{1}{n} \left[\sum_{i=1}^k P_i S_i^2 - (\bar{S})^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^k P_i (S_i - \bar{S})^2 > 0$$

where $\bar{S} = \sum_{i=1}^k P_i S_i$ is a weighted average

$$\therefore V_{prop} - V_{opt} \geq 0$$

$$V_{prop} \geq V_{opt} \Rightarrow V_{opt} \leq V_{prop} \rightarrow \textcircled{1}$$

$$\text{we know that } V_{ran} = \left(\frac{1}{n} - \frac{1}{N}\right) S^2 \rightarrow \textcircled{2}$$

we shall first express s^2 in terms of s_i^2

$$\text{we have } S^2 = \frac{1}{N-1} \sum_{i=1}^k \sum_{j=1}^{N_i} (Y_{ij} - \bar{Y}_N)^2$$

$$(N-1)S^2 = \sum_{i=1}^k \sum_{j=1}^{N_i} (Y_{ij} - \bar{Y}_{N_i} + \bar{Y}_{N_i} - \bar{Y}_N)^2$$

$$(N-1)S^2 = \sum_{i=1}^k \sum_{j=1}^{N_i} \left[(Y_{ij} - \bar{Y}_{N_i}) + (\bar{Y}_{N_i} - \bar{Y}_N) \right]^2$$

$$(N-1)S^2 = \sum_{i=1}^k \sum_{j=1}^{N_i} (Y_{ij} - \bar{Y}_{N_i})^2 + \sum_{i=1}^k \sum_{j=1}^{N_i} (\bar{Y}_{N_i} - \bar{Y}_N)^2 + 2 \sum_{i=1}^k \sum_{j=1}^{N_i} (Y_{ij} - \bar{Y}_{N_i})(\bar{Y}_{N_i} - \bar{Y}_N)$$

$$= \sum_{i=1}^k (N_i - 1) s_i^2 + \sum_{i=1}^k N_i (\bar{Y}_{N_i} - \bar{Y}_N)^2 + 0$$

$$\text{since } \sum_{j=1}^{N_i} (Y_{ij} - \bar{Y}_{N_i}) = 0$$

$$(N-1)S^2 = \sum_{i=1}^k (N_i - 1) s_i^2 + \sum_{i=1}^k N_i (\bar{Y}_{N_i} - \bar{Y}_N)^2$$

(algebraic sum of deviation take from its mean is always zero)

$$\because s_i^2 = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (Y_{ij} - \bar{Y}_{N_i})^2$$

If we consider N_i and N are sufficiently large,

$$\text{hence } N_i - 1 \approx N_i$$

$$N - 1 \approx N$$

$$\text{we get } NS^2 = \sum_{i=1}^k N_i s_i^2 + \sum_{i=1}^k N_i (\bar{Y}_{N_i} - \bar{Y}_N)^2$$

$$S^2 = \sum_{i=1}^k \frac{N_i}{N} s_i^2 + \sum_{i=1}^k \frac{N_i}{N} (\bar{Y}_{N_i} - \bar{Y}_N)^2$$

$$S^2 = \sum_{i=1}^k P_i s_i^2 + \sum_{i=1}^k P_i (\bar{Y}_{N_i} - \bar{Y}_N)^2 \quad \left(P_i = \frac{N_i}{N} \right)$$

Substitute this value of s^2 in eqn $\textcircled{2}$, we get

$$V(\bar{Y}_n)_{ran} = \left(\frac{1}{n} - \frac{1}{N}\right) \left(\sum_{i=1}^k P_i s_i^2 + \sum_{i=1}^k P_i (\bar{Y}_{N_i} - \bar{Y}_N)^2 \right)$$

$$V(\bar{Y}_n)_{ran} = \left(\frac{1}{n} - \frac{1}{N}\right) \sum_{i=1}^k P_i s_i^2 + \left(\frac{1}{n} - \frac{1}{N}\right) \sum_{i=1}^k P_i (\bar{Y}_{N_i} - \bar{Y}_N)^2$$

$$V(\bar{Y}_n)_{ran} = V_{prop} + \left(\frac{1}{n} - \frac{1}{N}\right) \sum_{i=1}^k P_i (\bar{Y}_{N_i} - \bar{Y}_N)^2$$

$$\text{since } \left(\frac{1}{n} - \frac{1}{N}\right) \sum_{i=1}^k P_i (\bar{Y}_{N_i} - \bar{Y}_N)^2 \geq 0$$

$$V(\bar{Y}_n)_{ran} = V_{prop} \geq 0 \rightarrow \textcircled{2}$$

Cost of obtaining information of a sample in one stratum will be usually different from other strata.

For example, the cost of collecting the data from rural areas will be more because of travelling expenses than from urban areas. Let c_i be the cost per unit in the i th stratum and let 'a' be the overhead cost, then the cost function c in the stratified random sampling is

$$c = a + \sum_{i=1}^k c_i n_i$$

Theorem 5 :- $\text{var}(\bar{y}_{st})$ is minimum for fixed total sample size n if $n_i \propto N_i s_i$

Proof :- We have to minimize

$$V(\bar{y}_{st}) = \frac{1}{N^2} \sum_{i=1}^k N_i (N_i - n_i) \frac{s_i^2}{n_i}$$

subject to the condition that $\sum_{i=1}^k n_i = n$

This is equivalent of minimising the Lagrangian function ϕ for the variations in n_i as given below.

$\phi = \text{var}(\bar{y}_{st}) + \lambda \left(\sum_{i=1}^k n_i - n \right)$ where λ is Lagrange multiple

$$\because \phi = \frac{1}{N^2} \sum_{i=1}^k N_i (N_i - n_i) \frac{s_i^2}{n_i} + \lambda \left(\sum_{i=1}^k n_i - n \right)$$

$$\phi = \frac{1}{N^2} \sum_{i=1}^k N_i \left(\frac{N_i}{n_i} - 1 \right) s_i^2 + \lambda \left(\sum_{i=1}^k n_i - n \right)$$

minimise the function ϕ differentiating ϕ with respect and equate it to zero, we get

$$\frac{\partial \phi}{\partial n_i} = 0 \Rightarrow \frac{1}{N^2} N_i^2 \left(-\frac{1}{n_i^2} \right) s_i^2 + \lambda(1) = 0 \rightarrow (1)$$

$$-\frac{N_i^2 s_i^2}{N^2 n_i^2} + \lambda = 0 \quad \lambda = \frac{N_i^2 s_i^2}{N^2 n_i^2} \quad n_i^2 = \frac{N_i^2 s_i^2}{N^2 \lambda}$$

$$n_i = \frac{N_i s_i}{N \sqrt{\lambda}} \quad \text{--- (2)}$$

$$\frac{\partial^2 \phi}{\partial n_i^2} = \frac{N_i^2 s_i^2}{N^2} \left(\frac{2}{n_i^3} \right) > 0 \text{ from (1)}$$

value of n_i in eqn (2) provides minimum

for this summing over i from $1, 2, \dots, k$ we get

$$\sum_{i=1}^k n_i = \frac{\sum_{i=1}^k N_i S_i}{N\sqrt{1}}$$

$$n = \frac{\sum_{i=1}^k N_i S_i}{N\sqrt{1}}$$

$$\sqrt{1} = \frac{\sum_{i=1}^k N_i S_i}{Nn}$$

Substitute the value of $\sqrt{1}$ in eqn (2)

we get $n_i = \frac{N_i S_i}{\frac{N \sum_{i=1}^k N_i S_i}{Nn}}$

$$n_i = \frac{n \cdot N_i S_i}{\sum_{i=1}^k N_i S_i} \quad \text{--- (3)}$$

this is the value of n_i in the optimum allocation for fixed total sample size n
 $n_i \propto N_i S_i \quad i=1, 2, \dots, k$ $\left(\frac{n}{\sum_{i=1}^k N_i S_i} = \text{Constant} = c \right)$

The value of n_i in eqn (3) is known as Neyman's formula for optimum allocation.

To obtain more precise estimator of the population mean we have to minimise the variance of the sample mean in optimum allocation.

for this, Neyman's optimum allocation suggests that the size of sample selected from the stratum to be more for the value of $N_i S_i$ is large i.e. N_i is large and S_i is large.

Theorem 6 :- In stratified random sampling for a specified cost function, $\text{var}(\bar{y}_{st})$ is minimum if $n_i \propto$

$$\frac{N_i S_i}{\sqrt{C_i}}$$

Proof :- We have to minimise $\text{var}(\bar{y}_{st}) = \frac{1}{N^2} \sum_{i=1}^k N_i (N_i - n_i) \frac{S_i^2}{n_i}$ subject to the condition that $n = \sum_{i=1}^k C_i n_i$ --- (2)

using the Lagrangian function below

where λ being Lagrange's multiplier

$$\therefore \phi = \frac{1}{N^2} \sum_{i=1}^k N_i (N_i - n_i) \frac{s_i^2}{n_i} + \lambda \left(\sum_{i=1}^k c_i n_i - c + a \right)$$

$$\phi = \frac{1}{N^2} \sum_{i=1}^k N_i \left(\frac{N_i}{n_i} - 1 \right) s_i^2 + \lambda \left(\sum_{i=1}^k c_i n_i - c + a \right) \quad \text{--- (3)}$$

To minimise the function ϕ differentiating ϕ with respect to n_i and equate it to zero, we get

$$\frac{d\phi}{dn_i} = 0 \Rightarrow \frac{N_i^2}{N^2} \left(-\frac{1}{n_i^2} \right) s_i^2 + \lambda c_i = 0$$

$$\frac{N_i^2 s_i^2}{N^2 n_i^2} = \lambda c_i$$

$$n_i^2 = \frac{N_i^2 s_i^2}{N^2 \lambda c_i}$$

$$n_i = \frac{N_i s_i}{N \sqrt{\lambda c_i}} \quad \text{--- (4)}$$

$$\text{Now } \frac{d^2\phi}{dn_i^2} = \frac{N_i^2 s_i^2}{N^2} \left(\frac{2}{n_i^3} \right) > 0$$

\therefore The value of n_i in eqn (4) provides minimum for ϕ

$$\therefore n_i = \frac{N_i s_i}{N \sqrt{\lambda c_i}}$$

We have to determine the value of λ for this summing over i from $1, 2, \dots, k$ we get

$$\sum_{i=1}^k n_i = \sum_{i=1}^k \frac{N_i s_i}{N \sqrt{\lambda c_i}}$$

$$n = \frac{\sum_{i=1}^k N_i s_i / \sqrt{c_i}}{N \sqrt{\lambda}}$$

$$\sqrt{\lambda} = \frac{\sum_{i=1}^k N_i s_i / \sqrt{c_i}}{N n}$$

Substitute the value of $\sqrt{\lambda}$ in eqn (4) we get

$$n_i = \frac{N_i s_i / \sqrt{c_i}}{\frac{\sum_{i=1}^k N_i s_i / \sqrt{c_i}}{N n}}$$

$$n_i = n \frac{N_i s_i / \sqrt{c_i}}{\sum_{i=1}^k N_i s_i / \sqrt{c_i}}$$

$$n_i c_i \frac{N_i s_i}{\sqrt{c_i}}$$

Systematic Random Sampling

**
1. Explain systematic sampling, merits and demerits of systematic sampling.

A. Systematic sampling is a very simple technique. It is generally used if the complete and up-to-date sampling frame (units) is available.

This sampling technique has a nice feature of selecting the whole sample with just one random strata.

A sampling technique in which only the first unit is selected with the help of random numbers or lottery method at random and the rest being automatically selecting according to some pre-determined pattern involving regular spacing of units is known as systematic random sampling or systematic sampling.

Suppose, there is a population with N units. If a sample of size n is to be selected, then divide the population into n groups such that each group consists of k units.

i.e. the population is divided into n groups in such a way that $N = nk$ or $k = \frac{N}{n}$ where k is an integer k is called the sampling interval.

Systematic sampling consists in drawing a random number say i ($i \leq k$) and selecting the unit corresponding to this number and every k^{th} unit subsequently.

Thus the systematic sample of size n consists of the units.

$i, i+k, i+2k, \dots, i+(n-1)k$

The random number i is called the random strata and its value determines the whole sample.

Random Sampling.

Example :- If $N=200$ $n=10$ and $k = \frac{N}{n} = \frac{200}{10} = 20$
Suppose unit number 15 is selected at random from the form the first 20 units, then the remaining units at the sample are 35, 55, 75, 95, 115, 135, 155, 175 and 195.

Merits (Advantages)

1. Systematic sampling is a very simple technique operationally it is more convenient than simple random sampling or stratified random sampling.
2. Time, work and money required for selecting a systematic sample is relatively much less.
3. Systematic sampling may be more efficient than simple random sampling provided the sampling frame is arranged completely at random.
4. The most commonly used to achieve randomness is alphabetical order of the population units.
For example, Telephone directory.

Demerits (Disadvantages)

- The main Disadvantages of systematic sampling is that systematic samples are not generally random samples, only the first unit is selected at random.

To achieve randomness, it is very difficult to get randomly arranged sampling frame.

If the population size N is not a multiple of the sample size n , then the actual sample size is different from the required sample size.

In that case the sample mean will not be an unbiased estimate of the population mean.

Such sampling may give highly biased

$$V(\bar{y}_{sys}) = \frac{1}{k} \sum_{i=1}^k (\bar{y}_{i.} - \bar{y})^2$$

It is not possible to obtain an unbiased estimate of this variance (population variance).

NOTATIONS AND TERMINOLOGY

Let y be the population characteristic

N = Total number of population units

n = Total number of sample units

$k = \frac{N}{n}$ = sampling interval

y_{ij} = value of j^{th} unit of the i^{th} sample
 ($i = 1, 2, \dots, k$; $j = 1, 2, \dots, n$)

$\bar{y}_{sys} = \bar{y}_{i.}$ = Mean of i^{th} systematic sample
 (or) mean of the i^{th} random start in the sample.

$$= \frac{1}{n} \sum_{j=1}^n y_{ij}, \quad i = 1, 2, \dots, k.$$

$\bar{y}_{..}$ = population mean

$$= \frac{1}{nk} \sum_{i=1}^k \sum_{j=1}^n y_{ij}$$

$$(or) = \frac{1}{k} \sum_{i=1}^k \bar{y}_{i.}$$

S^2 = population mean square

$$= \frac{1}{N-1} \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{..})^2 = \frac{1}{nk-1} \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{..})^2$$

The k possible systematic samples together with means and probability are given below in the following table.

Random start	Sample Composition	Probability	Mean
1	1, 1+k, ..., 1+jk, ..., 1+(n-1)k	1/k	$\bar{y}_{1.}$
2	2, 2+k, ..., 2+jk, ..., 2+(n-1)k	1/k	$\bar{y}_{2.}$
⋮	⋮	⋮	⋮
i	$i, i+k, \dots, i+jk, \dots, i+(n-1)k$	1/k	$\bar{y}_{i.}$
⋮	⋮	⋮	⋮
k	$k, 2k, \dots, (i+j)k, \dots, nk$	1/k	$\bar{y}_{k.}$

In this table, the units are referred as a strata. It is very clear from the above that each of the N units occur once in only one of the k samples and thus has an equal chance of being included into the sample since the probability of selecting the i th systematic sample is $\frac{1}{k}$, $i=1, 2, \dots, k$.

Theorem 1:— Sample mean is an unbiased estimator of the population mean in systematic sampling.

Proof:— We know,

$$\begin{aligned} \text{Population mean} = \bar{y}_{..} &= \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^n y_{ij} \quad (\text{or}) = \frac{1}{nk} \sum_{i=1}^k \sum_{j=1}^n y_{ij} \\ &= \frac{1}{k} \sum_{i=1}^k \left(\frac{1}{n} \sum_{j=1}^n y_{ij} \right) \end{aligned}$$

$$(\text{or}) = \frac{1}{k} \sum_{i=1}^k \bar{y}_i.$$

Consider

$$\begin{aligned} E(\bar{y}_{\text{sys}}) &= E(\bar{y}_i) = \bar{y}_1 \left(\frac{1}{k} \right) + \bar{y}_2 \left(\frac{1}{k} \right) + \dots + \bar{y}_k \left(\frac{1}{k} \right) \\ &= \frac{1}{k} \sum_{i=1}^k \bar{y}_i = \bar{y}_{..} \end{aligned}$$

\therefore The sample mean is an unbiased estimator of the population mean.

$$\text{i.e., } E(\bar{y}_i) = E(\bar{y}_{\text{sys}}) = \bar{y}_{..}$$

Theorem 2:— Variance of the systematic sample mean is given by $\text{var}(\bar{y}_{\text{sys}}) = \frac{N-1}{N} \cdot S^2 - \frac{k(n-1)}{N} \cdot S_{\text{wsy}}^2$

Where S_{wsy}^2 = population mean square among the units which i.e., within the same systematic sample.

$$= \frac{1}{k(n-1)} \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2$$

Proof:— Variance of the systematic sample mean is $\text{var}(\bar{y}_{\text{sys}}) = \text{var}(\bar{y}_i)$

$$= E \left[\bar{y}_i - E(\bar{y}_i) \right]^2$$

$$= E \left[\bar{y}_i - \bar{y}_{..} \right]^2 \quad (\text{from theorem})$$

$$= \frac{1}{k} \sum_{i=1}^k (\bar{y}_i - \bar{y}_{..})^2$$

As we know, the population mean square in the systematic sampling, $S^2 = \frac{1}{N-1} \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{..})^2$

$$= (N-1) S^2 = \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_i + \bar{y}_i - \bar{y}_{..})^2$$

$$= \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{i.})^2 + \sum_{i=1}^k \sum_{j=1}^n (\bar{y}_{i.} - \bar{y}_{..})^2 + 2 \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{i.}) (\bar{y}_{i.} - \bar{y}_{..})$$

Covariance term vanish

$$\begin{aligned} \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{i.}) (\bar{y}_{i.} - \bar{y}_{..}) &= \sum_{i=1}^k (\bar{y}_{i.} - \bar{y}_{..}) \sum_{j=1}^n (y_{ij} - \bar{y}_{i.}) \\ &= \sum_{i=1}^k (\bar{y}_{i.} - \bar{y}_{..}) \left(\sum_{j=1}^n y_{ij} - n \bar{y}_{i.} \right) \\ &= 0 \quad \left(\because \sum_{j=1}^n y_{ij} = n \bar{y}_{i.} \right) \end{aligned}$$

$$\therefore (N-1)S^2 = \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{i.})^2 + \sum_{i=1}^k \sum_{j=1}^n (\bar{y}_{i.} - \bar{y}_{..})^2$$

$$(N-1)S^2 = \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{i.})^2 + n \sum_{i=1}^k (\bar{y}_{i.} - \bar{y}_{..})^2$$

$$(N-1)S^2 = k(n-1)S_{\text{swsy}}^2 + nk v(\bar{y}_{\text{sys}}) \quad (\text{from } \textcircled{1})$$

Divide with $N = nk$, we get

$$\frac{N-1}{N} S^2 = \frac{k(n-1)}{N} S_{\text{swsy}}^2 + \frac{nk v(\bar{y}_{\text{sys}})}{N}$$

$$v(\bar{y}_{\text{sys}}) = \frac{N-1}{N} S^2 - \frac{k(n-1)}{nk} S_{\text{swsy}}^2 \quad (\because N = nk)$$

$$v(\bar{y}_{\text{sys}}) = \frac{N-1}{N} S^2 - \frac{k(n-1)}{nk} S_{\text{swsy}}^2$$

Theorem $\textcircled{3}$:- variance of the systematic sample mean

$$v(\bar{y}_{\text{sys}}) = \frac{nk-1}{nk} \frac{S^2}{n} [1 + (n-1)\rho]$$

where ρ is the intraclass correlation coefficient between the units of the same systematic sample and is given by

$$\rho = \frac{\sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_{..}) (y'_{ij} - \bar{y}_{..})}{(n-1)(nk-1)S^2}$$

Proof :- We know

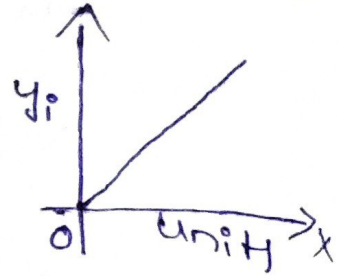
$$v(\bar{y}_{\text{sys}}) = \frac{1}{k} \sum (\bar{y}_{i.} - \bar{y}_{..})^2$$

... at the population ... trend, then prove that

$$\text{var}(\bar{Y}_{st}) \leq \text{var}(\bar{Y}_{sys}) \leq \text{var}(\bar{Y}_n)_{ran}$$

Proof :- Let us consider the population has a linear trend with the population units

Y_1, Y_2, \dots, Y_N takes the values
 $1, 2, \dots, N$ i.e. $Y_i = i$ $i = 1, 2, \dots, N$



Then population Total

$$\begin{aligned} \sum_{i=1}^N Y_i &= \sum_{i=1}^N i = 1 + 2 + \dots + N \\ &= \frac{N(N+1)}{2} \end{aligned}$$

$$\begin{aligned} Y_i &= a_i + b \\ a &= 1 \quad b = 0 \\ Y_i &= i \end{aligned}$$

$$\text{Population mean} = \bar{Y}_N = \frac{1}{N} \sum_{i=1}^N Y_i \Rightarrow \frac{1}{N} (1 + 2 + \dots + N)$$

$$\bar{Y}_N = \frac{N(N+1)}{2N} = \frac{N+1}{2}$$

$$\sum_{i=1}^N Y_i^2 = \sum_{i=1}^N i^2 = 1^2 + 2^2 + \dots + N^2 = \frac{N(N+1)(2N+1)}{6}$$

$$\text{Population mean square} = S^2 = \frac{1}{N-1} \left[\sum_{i=1}^N Y_i^2 - N \bar{Y}_N^2 \right]$$

$$S^2 = \frac{1}{N-1} \left[\frac{N(N+1)(2N+1)}{6} - N \left(\frac{N+1}{2} \right)^2 \right]$$

$$= \frac{1}{N-1} \left[\frac{N(N+1)(2N+1)}{6} - \frac{N(N+1)^2}{4} \right]$$

$$= \frac{N(N+1)}{2(N-1)} \left[\frac{2N+1}{3} - \frac{N+1}{2} \right]$$

$$s^2 = \frac{N(N+1)}{2(N+1)} \left[\frac{4N+2-3N-3}{6} \right]$$

$$s^2 = \frac{N(N+1)}{2(N+1)} \left[\frac{N-1}{6} \right]$$

$$s^2 = \frac{N(N+1)}{12}$$

∴ variance of the Sample mean in SRSWOR

$$V(\bar{y}_n)_{\text{Ran}} = \frac{N-n}{Nn} s^2$$

$$= \frac{N-n}{Nn} \frac{N(N+1)}{12} = \frac{N-k-n}{n} \frac{(Nk+1)}{12}$$

$$V(\bar{y}_n)_{\text{Ran}} = \frac{N-k-n}{n} \frac{(Nk+1)}{12} = \frac{(k-1)(Nk+1)}{12} \quad \text{--- (1)}$$

In Stratified Random Sampling variance of the Estimate of Population mean

$$V(\bar{y}_{\text{st}}) = \frac{1}{N^2} \sum_{i=1}^K N_i (N_i - n_i) \frac{S_i^2}{n_i}$$

Let us assume that the variance in each stratum of the population is equal and it is S_i^2 and $n_i = \frac{n}{N} N_i$

$$V(\bar{y}_{\text{st}}) = \frac{1}{N^2} \sum_{i=1}^K N_i \left(N_i - \frac{n}{N} N_i \right) \frac{S_i^2}{\frac{n}{N} N_i}$$

$$\text{Let us } = \frac{1}{N^2} \sum_{i=1}^K N_i \left(1 - \frac{n}{N} \right) \frac{S_i^2}{\frac{n}{N}}$$

$$= \frac{S_i^2}{N^2} \sum_{i=1}^K N_i \left(1 - \frac{n}{N} \right) \frac{N}{n}$$

$$= \frac{S_i^2}{N^2} \frac{N-n}{N} \sum_{i=1}^K N_i$$

$$= \frac{S_i^2}{N^2} \left(\frac{N-n}{N} \right) \frac{N}{n} N$$

$$= \frac{S_i^2}{N^2} \left(\frac{N-n}{N} \right) \frac{N}{n} N$$

$$V(\bar{y}_{\text{st}}) = \frac{S_i^2 (N-n)}{Nn}$$

we have $s^2 = \frac{N(N+1)}{12}$ for population as N units.

But in each k units in each stratum

$$\begin{aligned} & \frac{S_i^2}{N^2} \sum_{i=1}^K N_i \left(1 - \frac{n}{N} \right) \frac{N}{n} \\ &= \frac{S_i^2}{N^2} \left[\frac{N-n}{N} \right] \frac{N}{n} \\ &= \frac{S_i^2 (N-n)}{Nn} \end{aligned}$$

$$= \frac{k(k+1)(nk-n)}{12nk} \quad (N=nk)$$

$$= \frac{k(k+1)n(k-1)}{12nk} = \frac{k^2-1}{12n} \quad (2)$$

Population units which have the linear trend can be divided as systematic samples

$$\bar{y}_{i.} = \frac{1}{n} \sum_{j=1}^n y_{ij}$$

$$= \frac{1}{n} [i + (i+k) + (i+2k) + \dots + i + (n-1)k]$$

$$= \frac{1}{n} [ni + [1+2+\dots+(n-1)]k]$$

$$= \frac{1}{n} (ni + \frac{(n-1)n}{2}k) = i + \frac{(n-1)}{2}k$$

$$\text{and } \bar{y}_{..} = \bar{y}_{N..} = \frac{N+1}{2} = \frac{nk+1}{2}$$

$$\therefore \bar{y}_{i.} - \bar{y}_{..} = i + \frac{(n-1)k}{2} - \frac{nk+1}{2} = i - \frac{k+1}{2}$$

$$\text{Now } v(\bar{y}_{sys}) = \frac{1}{k} \sum_{i=1}^k (\bar{y}_{i.} - \bar{y}_{..})^2$$

$$= \frac{1}{k} \sum_{i=1}^k \left(i - \frac{k+1}{2}\right)^2$$

$$= \frac{1}{k} \left[\sum_{i=1}^k i^2 + \sum_{i=1}^k \left(\frac{k+1}{2}\right)^2 - 2 \frac{(k+1)}{2} \sum_{i=1}^k i \right]$$

$$= \frac{1}{k} \left[\frac{k(k+1)(2k+1)}{6} + \frac{k(k+1)^2}{4} - \frac{(k+1)k(k+1)}{2} \right]$$

$$= \frac{k+1}{2} \left[\frac{2k+1}{3} + \frac{k+1}{2} - (k+1) \right]$$

$$= \frac{k+1}{2} [4k+2+3k+3-6k-6]$$

$$= \frac{(k+1)(k-1)}{2 \cdot 6}$$

$$v(\bar{y}_{sys}) = \frac{k^2-1}{12} \quad (3)$$

From equations (1), (2) and (3) we get

$$v(\bar{y}_{st}) : v(\bar{y}_{sys}) : v(\bar{y}_{n})_{ran} ::$$

$$\frac{k^2-1}{12n} : \frac{k^2-1}{12} : (k-1)(nk+1)$$

$$\frac{k+1}{n} : k+1 : nk+1$$

which is approximately equal to $\frac{1}{n} : 1 : n$.

from (1), (2) & (3) and putting them in a ratio we get

$$V(\bar{Y}_{st}) : V(\bar{Y}_{sys}) : V(\bar{Y}_{ran}) :: \frac{k^2-1}{12n} : \frac{k^2-1}{12} : \frac{(k-1)(nk-1)}{12}$$

Dividing by $\frac{k-1}{12}$ we get the ratio as $\frac{1}{n} : 1 : \frac{nk+1}{k+1}$

$$\text{But } \frac{nk+1}{k+1} = \frac{(nk+1)/k}{(k+1)/k} \quad (\text{dividing } N_1 \text{ and } D_1 \text{ by } k)$$

$$= \frac{n + \frac{1}{k}}{1 + \frac{1}{k}} = n \quad \because \text{for large } k, \frac{1}{k} \rightarrow 0$$

The variance are in the ratio

$$\frac{1}{n} : 1 : n \quad (\text{approx})$$

$$\Rightarrow \frac{1}{n} \leq 1 \leq n$$

$$\text{Hence } V(\bar{Y}_{st}) \leq V(\bar{Y}_{sys}) \leq V(\bar{Y}_{ran})$$

The equality sign holds good if $n=1$

In that case each variance = $(k+1)$

Theorem 5 :- With usual notations prove that mean of a systematic sample is more precise (more efficient) than the mean of a simple random sample if and only if $S_{wsy}^2 > S^2$

In other words $V(\bar{Y}_{sys}) < V(\bar{Y}_{ran})$ if $S_{wsy}^2 > S^2$

Proof :- In simple random sampling without replacement variance of the sample mean is given by

$$V(\bar{Y}) = \frac{N-n}{N} \frac{S^2}{n}$$

Similarly in systematic sampling

$$V(\bar{Y}_{sys}) = \frac{N-1}{N} S^2 - \frac{k(n-1)}{N} S_{wsy}^2$$

$V(\bar{Y}_{sys})$ should be less than $V(\bar{Y})$

$$\frac{N-1}{N} S^2 - \frac{k(n-1)}{N} S_{wsy}^2 < \frac{N-n}{N} \frac{S^2}{n}$$

$$\left[(N-1) - \frac{(N-n)}{n} \right] S^2 < k(n-1) S_{wsy}^2$$

$$(N-n-N+n) S^2 < k(n-1) S_{wsy}^2$$

$$\frac{N(n-1)}{N} s^2 < k(n-1) S_{wsy}^2$$

$$\frac{nk(n-1)}{n} s^2 < k(n-1) S_{wsy}^2$$

$$k(n-1) s^2 < k(n-1) S_{wsy}^2$$

$$s^2 < S_{wsy}^2 \Rightarrow S_{wsy}^2 > s^2$$

$$v(\bar{y}_{sys}) < v(\bar{y}_{ran})$$

Thus $v(\bar{y}_{sys}) < v(\bar{y}_{ran})$ if $S_{wsy}^2 > s^2$

Hence mean of a systematic sample.

UNIT III - Analysis of Variance (ANOVA)

1. Explain the concept of analysis of variance (ANOVA) and also assumption, uses.

Sol :- The analysis of Variance is a powerful statistical tool for the test of significance (hypothesis)

1. The test of significance based on t-distribution.
2. The significant of difference between two sample means using t-distribution.
3. But if we want test the significance difference of three or more sample means.
4. Hence t-test can't be used.
5. We need an alternative procedure is required for testing the homogeneity of several means.

Eg :- five fertilizers are applied to four parts we may be interested to find out if the effect of these fertilizers. The yield is significantly different.

- * The answer to these problem is provided by the technique of analysis of variance to test the homogeneity of several means.
- * The alternative procedure is ANOVA.
- * The term analysis of variance was introduced by Prof. R.A. Fisher in 1920. He said variance is inherent in nature.
- * ANOVA is based on F-distribution.
- * The analysis of variance is a powerful statistical tool which support the controlled and uncontrolled variations in the data.

Defination :- In the words of R.A. Fisher ANOVA is the separation of variance ascribable to one group of causes from the variance ascribable to the other group.

Causes of variation :-

1. The variation in any experiment is inherent in nature.
2. The total variation in any numerical data is due to number of factors or causes. These are classified as

3. The variations due to assignable causes are called controlled and measured. Hence these are called controlled variations.

4. The variation due to chance causes is beyond the control of human hand and they can't be controlled. This variation is also known as error variation.

5. This variation is also known as error variation.

6. It plays an important role in ANOVA.

Uses:-

The technique of ANOVA is useful in many fields like agriculture, economics, biology, education, psychology, business and so on.

Assumptions:-

The technique of ANOVA is based on f-test. For the validity of f-test, the following assumptions are made in analysis of variance.

1. The sample observations are independent.

2. The parent population from which the sample are taken is normal.

3. Various effects are additive in nature.

4. E_{ij} is a random error.

$$E_{ij} \sim N(0, \sigma^2)$$

5. The analysis of variance may be classified as

(i) one-way classification

(ii) Two-way classification

Eg:-

Suppose four types of fertilizers are used to test their effects on the yield of paddy is significantly different or they have the same effect.

② Discuss one-way classification of data and write down the analysis of variance table.

⇒ ANOVA is the simplest technique.

⇒ In one way ANOVA observations are classified into groups (or) classes (or) samples on the basis of single criterion i.e; one assignable cause

for example 'n' students of a class are arranged

According to their marks

⇒ Let us suppose that N observations y_{ij} $\left[\begin{matrix} i=1,2,\dots,k \\ j=1,2,\dots,n_i \end{matrix} \right]$ of sizes n_1, n_2, \dots, n_k

⇒ The sample data can be arranged into k classes as shown below.

classes (treatment)	1	2	...	j	...	n_i	total	Mean
1	y_{11}	y_{12}	...	y_{1j}	...	y_{1n_1}	$T_{1.}$	$\bar{y}_{1.}$
2	y_{21}	y_{22}	...	y_{2j}	...	y_{2n_2}	$T_{2.}$	$\bar{y}_{2.}$
...
i	y_{i1}	y_{i2}	...	y_{ij}	...	y_{in_i}	$T_{i.}$	$\bar{y}_{i.}$
...
k	y_{k1}	y_{k2}	...	y_{kj}	...	y_{kn_k}	$T_{k.}$	$\bar{y}_{k.}$
							G	$\bar{y}_{..}$

$$\bar{y} = \frac{\sum y_{ij}}{k}$$

G = Grand total

\bar{y} = over all mean

If $n_1 = n_2 = \dots = n_k = n$ (say) then the data is called one way classified data with equal no. of observation. otherwise the data is known as one-way classified data with unequal no. of observations.

The total variations in the observations y_{ij} can be divided into the following two components.

i) The variation between the classes i.e., Treatment Effect; this is due to assignable causes which can be detected and Control by human hand.

ii) The variation within the classes i.e., inherent variation is known as Error effect. This is due to chance causes which cannot be control of human hand.

⇒ The sources of variation in the data are

i) Effect of the treatment of d_i $i=1,2,\dots,k$.

ii) Error due to chance causes [Random causes]

$$E_{ij} \sim N(0, \sigma^2)$$

The main objective of ANOVA technique is to examine if the variation due to different classes is significant.

Assumptions :-

is normal.

⇒ various effects are additive in nature

⇒ E_{ij} is a random error.

$$E_{ij} \sim N(0, \sigma^2)$$

Statistical Analysis of one-way classification:

T_i = Total yield of i^{th} treatment.

$$T_i = \sum_{j=1}^{n_i} y_{ij}$$

$$G = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}$$

G = Grand total
 N = Total No. of observation

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

$$\bar{y}_{..} = \frac{G}{N} \text{ [over all mean]}$$

$$\bar{y}_{i.} = \frac{\sum_{j=1}^{n_i} y_{ij}}{n_i} = \text{Mean of } i^{\text{th}} \text{ class.}$$

NULL Hypothesis:

$NH: H_0$: All the treatment are homogeneous

$NH: H_0$: population means are equal

i.e., $H_0: \mu_1 = \mu_2 = \dots = \mu_k$

Mathematical Model of one-way classification:

(linear additive model)

$$y_{ij} = \mu + \alpha_i + E_{ij}$$

$$y_{ij} = \mu + \alpha_i + \beta_j + E_{ij}$$

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_k + E_{ij}$$

In case of one-way classification, the linear mathematical model is given by $y_{ij} = \mu + \alpha_i + E_{ij}$

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

$j = 1, 2, \dots, n_i$

where y_{ij} = The yield from the j^{th} unit by applying i^{th} treatment.

α_i = The i^{th} treatment effect.

E_{ij} = Error effect due to random.

μ = General mean effect

the equation is the basis for calculation

Estimation of Parameters:-

$$y_{ij} = \mu + \alpha_i + \beta_j$$

$$E_{ij} = y_{ij} - \mu - \alpha_i$$

The parameters μ, α_i are estimated by using the principle of least squares an minimizing the error sum of squares

$$F = \sum_{i=1}^k \sum_{j=1}^{n_i} E_{ij}^2$$

$$F = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \mu - \alpha_i)^2$$

If minimize F w.r. to μ

$$\frac{dF}{d\mu} = 0$$

$$\text{we get } \mu = \bar{y}_{..}$$

If minimize F w.r. to α_i

$$\frac{dF}{d\alpha_i} = 0$$

$$\text{we get } \alpha_i = \bar{y}_{i.} - \bar{y}_{..}$$

Various sum of squares:-

In one way classification the total variation can be split into two parts.

i) frequent sum of squares

ii) error sum of squares

The total sum of squares for corrected mean is

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

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

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

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

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

Since cross products vanish.

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

Total sum of squares = error of squares + treatment sum of squares.

$$S_T^2 = S_E^2 + S_t^2$$

In numerical data for simple arithmetic calculations we make the use of following formulae.

i) Correction factor (CF) = $\frac{G^2}{N}$

Where G = Grand total

N = Total No. of observations

ii) Total sum of squares $S_T^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}^2 - CF$

iii) Treatment sum of squares $S_t^2 = \frac{\sum T_i^2}{n} - CF$ [equal size]

$S_t^2 = \sum_{i=1}^k \left(\frac{T_i^2}{n_i} \right) - CF$ [unequal size]

iv) Error sum of squares $S_E^2 = S_T^2 - S_t^2$

Degrees of freedom :-

i) The dof for total sum of squares = $N-1$

ii) The dof for treatment sum of squares = $k-1$

iii) The dof for error sum of squares = $N-1-(k-1) = N-k$

Mean sum of squares :- (MSS)

The sum of squares divided by its dof gives the variance, it is also called Mean sum of squares.

i) Mean sum of squares due to Treatment = $S_t^2 = \frac{S_t^2}{k-1}$

$a_1 = \frac{S_t^2}{k-1}$

ii) Mean sum of squares due to Error = $S_E^2 = \frac{S_E^2}{N-k}$

$a_2 = \frac{S_E^2}{N-k}$

Test statistic :-

Hence the test statistic for H_0 is

$f_{cal} = \frac{MSS \text{ for treatment}}{MSS \text{ for error}}$
 $= \frac{S_t^2}{S_E^2} = \frac{a_1}{a_2} \text{ in } F \left(\begin{matrix} k-1, N-k \\ \text{at } 5\% \text{ \& } 1\% \text{ LOS} \end{matrix} \right)$

ANOVA Table :-

f-ratio
 $f_{cal} = \frac{S_t^2}{S_E^2}$
 f_{table} for $(k-1, N-k)$ at 5% & 1%.

Conclusion :-

If $f_{cal} \leq f_{table}$ for $(k-1, N-k)$ dof at required Probability level (5% or 1%) then we accept H_0 otherwise we reject H_0 .

③ Explain Estimation of parameters in one way classification.

Ans: * The parameters in the model μ, α_i are estimated by using the principle of least squares. on minimizing the Error sum of squares.

* In one-way classification the linear mathematical model is

$$Y_{ij} = \mu + \alpha_i + E_{ij} \quad \begin{array}{l} i=1, 2, \dots, k \\ j=1, 2, \dots, n_i \end{array}$$

* Random Error = $E_{ij} = Y_{ij} - \mu - \alpha_i$

$$\begin{aligned} * \text{Error sum of squares} = E &= \sum_{i=1}^k \sum_{j=1}^{n_i} E_{ij}^2 \\ &= \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \mu - \alpha_i)^2 \end{aligned}$$

We minimize E w.r. to μ

$$\frac{dE}{d\mu} = 0$$

$$\frac{d}{d\mu} \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \mu - \alpha_i)^2 = 0$$

$$2 \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \mu - \alpha_i) (-1) = 0$$

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \mu - \alpha_i) = 0$$

$$\sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij} - \sum_{i=1}^k \sum_{j=1}^{n_i} \mu - \sum_{i=1}^k \sum_{j=1}^{n_i} \alpha_i = 0$$

$$\sum \sum Y_{ij} - kn\mu - n_i \sum_{i=1}^k \alpha_i = 0$$

$$\text{since } \sum_{i=1}^k \alpha_i = 0 \quad \begin{array}{l} [nk=N] \\ [nik=N] \end{array}$$

$$\sum \sum Y_{ij} - N\mu = 0$$

$$N\mu = \sum \sum Y_{ij}$$

$$\mu = \frac{\sum \sum Y_{ij}}{N} = \bar{Y}_{..}$$

$$\boxed{\mu = \bar{Y}_{..}}$$

We minimize E w.r. to α_i

$$\frac{dE}{d\alpha_i} = 0$$

$$\frac{d}{d\alpha_i} \sum_{j=1}^{n_i} (Y_{ij} - \mu - \alpha_i) (-1) = 0$$

$$\sum_{j=1}^{n_i} (y_{ij} - \mu - \alpha_i) (-1) = 0$$

$$\sum_{j=1}^{n_i} (y_{ij} - \mu - \alpha_i) = 0$$

$$\sum_{j=1}^{n_i} y_{ij} - \sum_{j=1}^{n_i} \mu - \sum_{j=1}^{n_i} \alpha_i = 0$$

$$\sum_{j=1}^{n_i} y_{ij} - n_i \mu - n_i \alpha_i = 0$$

$$n_i \alpha_i = \sum_{j=1}^{n_i} y_{ij} - n_i \mu$$

$$\alpha_i = \frac{\sum_{j=1}^{n_i} y_{ij}}{n_i} - \frac{n_i \mu}{n_i}$$

$$\alpha_i = \bar{y}_{i.} - \mu$$

$$\boxed{\alpha_i = \bar{y}_{i.} - \bar{y}_{..}}$$

DISCUSS TWO-WAY CLASSIFICATION OF DATA AND ALSO WRITE DOWN ANOVA TABLE.

4. Explain statistical analysis of two-way classification

A. Two-way Classification :-

In some cases the resulting observations are affected by two factors and we need to examine the effect of these factors on the observations. Then statistical data is divided into two on the basis of two factors. This type of classification is known as Two-way classification.

Ex :- 1) A sample of n individuals can be classified according to their height and weight & according to age and weight.
 > Let us take the yield of Paddy may be affected by differences in seeds and differences in different levels of fertilizers.

Let us know N observations are divided with respect to two characteristics say A, B.

First we divide the observations into k classes w.r.t to characteristic A as rows and divide the observations into classes w.r.t to characteristic B as columns.

$$\boxed{\text{i.e., } N = hk}$$

Here y_{ij} denote the yield from experimental unit of i^{th} row and j^{th} column. In this case of observations are

$$u_{i.} = \frac{\sum_{j=1}^h y_{ij}}{h} \quad \forall i=1, \dots, k$$

$$\bar{y}_{.j} = \frac{\sum_{i=1}^k y_{ij}}{k} \quad \forall j=1, 2, \dots, h$$

$$\bar{y}_{..} = \frac{\sum_{i=1}^k \sum_{j=1}^h y_{ij}}{N} = \frac{\sum_{i=1}^k \sum_{j=1}^h y_{ij}}{hk} = \frac{\sum_{i=1}^k \bar{y}_{i.}}{k} = \frac{\sum_{j=1}^h \bar{y}_{.j}}{h}$$

$$N = hk$$

k = No. of Rows

h = No. of columns.

Null Hypothesis :-

H_{01} : All the rows are homogeneous.

i.e., $H_{01} : R_1 = R_2 = \dots = R_k \quad \left(\sum_{i=1}^k \alpha_i = 0 \right)$

H_{02} : All the columns are homogeneous

i.e., $H_{02} : C_1 = C_2 = \dots = C_h \quad \left(\sum_{j=1}^h \beta_j = 0 \right)$

Mathematical model of two-way classification :-

The linear mathematical model of the observations for two-way classification is

$$y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij} \quad \begin{matrix} i=1, 2, \dots, k \\ j=1, 2, \dots, h \end{matrix}$$

y_{ij} = yield from the i^{th} row and j^{th} column.

μ = General mean effect.

α_i = i^{th} row effect

β_j = j^{th} column effect.

ϵ_{ij} = Random effect. (Error effect due to random)

This equation is the basic equation for the calculation of VSS.

Estimation of μ .

Sum of squares (EES)

$$y_{ij} = \mu + \alpha_i + \beta_j + E_{ij}$$

$$E_{ij} = y_{ij} - \mu - \alpha_i - \beta_j$$

$$\text{Error sum of squares} = E = \sum_{i=1}^k \sum_{j=1}^h E_{ij}^2$$

we minimize E w.r.to μ

$$\frac{dE}{d\mu} = 0$$

$$\text{we get, } \mu = \bar{y}_{..}$$

we minimize E w.r.to α_i

$$\frac{dE}{d\alpha_i} = 0$$

$$\text{we get } \alpha_i = \bar{y}_{i.} - \bar{y}_{..}$$

we minimize E w.r.to β_j

$$\frac{dE}{d\beta_j} = 0$$

$$\text{we get } \beta_j = \bar{y}_{.j} - \bar{y}_{..}$$

Various sum of squares :- (VSS)

In two-way classification the total sum of squares can be split into three components

- i) variation due to Rows
- ii) variation due to Columns
- iii) variation due to Error

$$\left[(a+b+c)^2 = a^2 + b^2 + c^2 + \text{cross product} \right]$$

The total sum of squares for corrected mean.

$$\sum_{i=1}^k \sum_{j=1}^h (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^k \sum_{j=1}^h \left[(y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..}) + (\bar{y}_{i.} - \bar{y}_{..}) + (\bar{y}_{.j} - \bar{y}_{..}) \right]^2$$

$$\sum_{i=1}^k \sum_{j=1}^h (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^k \sum_{j=1}^h (y_{ij} - \bar{y}_{i.} + \bar{y}_{..})^2 + \sum_{i=1}^k \sum_{j=1}^h (\bar{y}_{i.} - \bar{y}_{..})^2 + \sum_{i=1}^k \sum_{j=1}^h (\bar{y}_{.j} - \bar{y}_{..})^2 + \text{cross product.}$$

Since the cross product vanishes.

The algebraic sum of the deviation of observations

1. Describe randomised block design (RBD & CBD). Give its layout, advantages and disadvantages.

RBD:— A method dividing the heterogeneous experimental material into relatively homogeneous subgroups or blocks and the treatments are applied randomly to relatively homogeneous experimental units within each block and replicated over all the blocks is known as Randomised Block Design (R.B.D)

For example, In agricultural experimentation, the experimental area

i.e. field is not homogeneous and the fertility gradient is only one direction, then a simple method of controlling the variability of the experimental material consists of dividing the heterogeneous field into homogeneous subgroups or blocks & replicates perpendicular to the direction of the fertility gradient.

Now if the treatments applied at random to homogeneous units within each block and replicated over all the blocks, the design is called Randomised block design. In a CRD we do not resort to the grouping of the experimental site and allocate the treatments at random to the experimental units.

But in RBD treatments are allocated at random within the units of each block. i.e. Randomisation restricted.

Also variation among blocks is removed from variation due to error.

Layout of RBD:—

In agricultural experiment, if we consider 5 treatments A, B, C, D, E each replicated four times. Then we divide the whole experimental area into four homogeneous blocks and each block into five units.

Treatments (fertilizers) are then allocated at random the units of a block. fresh randomisation being done for each block. The

layout of RBD as follows. Blocks Treatments

To allocate treatments randomly, we use any one of the methods, lottery method or random number tables method.

For randomization, we may use Tippett's random number tables.

Advantages (merits) :-

1. In this design three principles are used so it may be considered as a good design of Experiment.
2. The statistical analysis of RBD data is simple and easy to understand.
3. RBD provides more accurate results than CRD since the experimental material is divided into blocks this results in decreasing error variance, therefore experimental error is considerably reduced.
4. In RBD, no restrictions are placed on the number of treatments and on the number of replications. So the design is flexible. But at least two replications are required to test the significance to treatments of blocks.
5. It is more accurate (efficient) design than C.R.D
6. RBD is a very popular experiment and extensively used design in almost all the scientific experiments.

Disadvantages :-

1. RBD can not be used with unequal replications.
2. The randomization in RBD is restricted within each block.
3. It is suitable only if there is one factor of heterogeneity. If there is more than one factor of heterogeneity, it is not suitable.
4. RBD is not suitable for a large no. of treatments.
5. The analysis of RBD should be difficult in case of missing observation.

Statistical analysis of RBD :-

The statistical analysis of RBD is similar to the ANOVA for two way classified data.

Suppose a RBD experiment is carried out with k

Treatments	Blocks				Total	mean		
	1	2	...	j			...	h
1	y_{11}	y_{12}	...	y_{1j}	...	y_{1h}	T_1	\bar{y}_1
2	y_{21}	y_{22}	...	y_{2j}	...	y_{2h}	T_2	\bar{y}_2
...
i	y_{i1}	y_{i2}	...	y_{ij}	...	y_{ih}	T_i	\bar{y}_i
...
k	y_{k1}	y_{k2}	...	y_{kj}	...	y_{kh}	T_k	\bar{y}_k
Total	T_1	T_2	...	T_j	...	T_h	G	-
Mean	\bar{y}_1	\bar{y}_2	...	\bar{y}_j	...	\bar{y}_h	-	-

Total variation is splitted into three parts
 \therefore Total variation = variation in between Rows (Treatments effect) + variation in between Columns (Block effect) + Error effect.

$N = nk =$ is the total number of experimental units
 $k =$ no. of treatments.

$h =$ no. of blocks

$G = \sum_{i=1}^k \sum_{j=1}^h y_{ij}$ is grand Total

$T_i = \sum_{j=1}^h y_{ij}$ $T_j = \sum_{i=1}^k y_{ij}$

$\bar{y}_i = \frac{1}{h} \sum_{j=1}^h y_{ij}$ $\bar{y}_j = \frac{1}{k} \sum_{i=1}^k y_{ij}$ $\bar{y}_{..} = \frac{1}{hk} \sum_{i=1}^k \sum_{j=1}^h y_{ij}$

Null hypothesis :-

H_{01} : All the treatments are homogeneous
 i.e. $H_{01} : \alpha_1 = \alpha_2 = \dots = \alpha_k = 0$

H_{02} : All the blocks are homogeneous
 i.e. $H_{02} : \beta_1 = \beta_2 = \dots = \beta_h = 0$

Mathematical Model :- The linear mathematical model

becomes $y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$

- where y_{ij} is the yield from j^{th} block by receiving treatment. μ is general mean effect.
- α_i is effect due to i^{th} treatment
- β_j is effect due to j^{th} block
- ϵ_{ij} is error effect due to random and $\epsilon_{ij} \sim \text{iid } N(0, \sigma_e^2)$

Estimator of parameters :-

The parameters in the model μ, α_i, β_j are estimated by using the principle of least squares on minimizing the error sum of squares.

$\mu = \bar{y}_{..}$ $\alpha_i = \bar{y}_{i.} - \bar{y}_{..}$ $\beta_j = \bar{y}_{.j} - \bar{y}_{..}$
 Various sum of squares :- The total sum of squares (Total variation) for the corrected mean.

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

$$\sum_{i=1}^k \sum_{j=1}^h (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^k \sum_{j=1}^h (y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2 + \sum_{i=1}^k \sum_{j=1}^h (\bar{y}_{i.} - \bar{y}_{..})^2 + \sum_{i=1}^k \sum_{j=1}^h (\bar{y}_{.j} - \bar{y}_{..})^2 + \text{Cross Product}$$

The cross product in the above eqn vanish

$$\sum_{i=1}^k \sum_{j=1}^h (y_{ij} - \bar{y}_{..})^2 = \sum_{i=1}^k \sum_{j=1}^h (y_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2 + h \sum_{i=1}^k (\bar{y}_{i.} - \bar{y}_{..})^2 + k \sum_{j=1}^h (\bar{y}_{.j} - \bar{y}_{..})^2 + 0$$

$$TSS = ESS + \text{Treatment SS} + \text{Block SS}$$

$$S_T^2 = S_E^2 + S_t^2 + S_B^2$$

In numerical data for simple arithmetic calculations, we make the use of following formulae.

$$TSS = S_T^2 = \sum \sum y_{ij}^2 - CF \quad CF = G^2/N$$

$$\text{Treatment SS} = S_t^2 = \frac{\sum T_i^2}{h} - CF$$

$$\text{Block SS} = S_B^2 = \frac{\sum T_{.j}^2}{k} - CF$$

$$\text{Error SS} = S_E^2 = S_T^2 - S_t^2 - S_B^2$$

A NOVA TABLE FOR RBD

SV	Dof	SS	mss	F-ratio variance ratio
Treatments	k-1	S_t^2	$s_t^2 = \frac{S_t^2}{k-1}$	$f_1 = \frac{S_t^2}{S_E^2} \sim F(k, (k-1)(h-1))$
Blocks	h-1	S_B^2	$s_B^2 = \frac{S_B^2}{h-1}$	$f_2 = \frac{S_B^2}{S_E^2} \sim F(h-1, (k-1)(h-1))$
Error	(k-1)(h-1)	S_E^2	$s_E^2 = \frac{S_E^2}{(k-1)(h-1)}$	at 5% & 1% LOS
Total	hk-1 = N-1	S_T^2		

Conclusion :- for treatments, Blocks.

If $f_{cal} \leq f_{table}$ value we accept our hypothesis otherwise we reject our hypothesis at required probability level (5% & 1%)

Problem: Analyse the RBD layout

Sol: - H_0 : All the treatments (A, B, C, D, E) are homogeneous
 H_1 : All the blocks (I, II, III, IV) are homogeneous
 LOS: $\alpha = 5\%$

Treatments	Blocks				$T_{i.}$	$T_{i.}^2$	$\sum y_{ij}^2$
	I	II	III	IV			
A	12.5	12.4	14.2	11.6	50.7	2570.49	646.21
B	10.5	13.6	12.6	15.2	51.9	2693.61	685.01
C	11.2	12.5	13.7	12.3	49.7	2470.09	620.67
D	13.3	10.3	11.6	14.3	49.5	2450.25	622.03
E	11.1	14.1	10.4	12.3	47.9	2294.41	581.47
$T_{.j}$	58.6	62.9	62.5	65.7	249.7	12478.85	3155.39
$T_{.j}^2$	3433.96	3956.41	3906.25	4316.49	→ 15613.11		-

Here $k = \text{no. of treatments} = 5$ $h = \text{no. of blocks} = 4$
 $N = hk = 5 \times 4 = 20$ $G = 249.7$ $\sum \sum y_{ij}^2 = 3155.39$
 $\sum T_{i.}^2 = 12478.85$ $\sum T_{.j}^2 = 15613.11$

$C.F = \frac{G^2}{N} = \left(\frac{249.7^2}{20} \right) = 3117.5045$
 $TSS = S_T^2 = \sum \sum y_{ij}^2 - C.F = 3155.39 - 3117.5045 = 37.8855$

Treatment SS = $S_T^2 = \frac{\sum T_{i.}^2}{h} - C.F = \frac{12478.85}{4} - 3117.5045 = 2.208$
 Block SS = $S_B^2 = \frac{\sum T_{.j}^2}{k} - C.F = \frac{15613.11}{5} - 3117.5045 = 5.1175$
 Error SS = $S_E^2 = S_T^2 - S_T^2 - S_B^2 = 37.8855 - 2.208 - 5.1175 = 30.5$

ANOVA Table :-

SV	do f	S.S	M.S.S	f-ratio f _{cal}	f _{table}
----	------	-----	-------	-----------------------------	--------------------

Conclusion :-

Treatments: $f_{cal} = 0.2167$ $f_{table} = 3.26$
 $f_{cal} < f_{table}$ value so, we accept our H_0 at 5% LOS
 All the treatments (A, B, C, D, E) are homogeneous.
 Blocks: $f_{2cal} = 0.6698$ $f_{2table} \text{ value} = 3.49$
 $f_{2cal} < f_{2table}$ value so, we accept our H_0 at 5% LOS.
 we conclude that all the blocks are homogeneous.

I B(120) F(135) A(111) C(121) D(127)
 II C(100) -A(121) D(128) E(136) B(123)
 III C(121) -A(137) E(101) B(99) D(138)

$\bar{y} = 1818$ $\sum \sum y_{ij}^2 = 222822$ $\sum T_i^2 = 662722$ $\sum T_j^2 = 1101876$
 $k=5$ $h=3$ $N=15$ $CF = 220341.6$ $S_T^2 = 2480.4$
 $S_t^2 = 632.4$ $S_B^2 = 33.6$ $S_E^2 = 1814.4$ $f_{ical} = 0.6990 (3.84)$
 $f_{2cal} = 0.0740 (4.46)$

③ In an agricultural field experimentation, three fertilizers are applied in four randomized blocks and the yield of wheat in kilos are given below.

Blocks			
I	II	III	IV
A 8	C 10	A 6	B 10
C 12	B 8	B 9	A 8
B 10	A 8	C 10	C 9

1. Analyse the data using RBD and state conclusion

$f_1 = 7.8 (5.14)$ $f_2 = 1.61 (4.76)$

2. Find efficiency of RBD over CRD

$E = \frac{h(k-1)S_E^2 + (h-1)S_E^2}{(hk-1)S_E^2} = 1.17$

Latin Square Design (LSD)

③ Describe LSD, layout, mention its advantages and disadvantages.

Ans:- Latin square design (LSD)
 The RBD attempts to control the variability in experimental material in one direction only. But there may be situations that heterogeneity may be in two perpendicular directions, i.e. horizontally as well as vertically (11). For example in agricultural experiments soil variation may be in two perpendicular directions. In such situations RBD may be of little use.

A layout known as LSD is of great use to control the variation due to two factors. The entire heterogeneous experimental material is relatively homogeneous block with respect to rows and columns, treatment in

rows and columns in such a way that every treatment occurs once and only once in each row and in each column, such a design is called Latin square design. This helps in elimination of row and column effects from the experimental error and tries to make the experiment more sensitive. LSD is extensively used in agricultural experiments. Also it is used in industrial and animal husband experiments.

Layout of LSD:

In LSD, the number of treatments and number of rows is equal to number of columns. If we consider m treatments, then there will be $m \times m$ experimental units. The whole experimental material is divided into m^2 experimental units arranged in a square so that each row and each column consists m experimental units. Then m treatments are allocated at random to these rows and columns in such a way that each and every treatment occurs once and only once in each row and in each column. This layout is called $m \times m$ LSD.

For example, if there are five treatments, then 5×5 LSD layout can be explained in the following table.

A	B	C	D	E
B	C	D	E	A
C	D	E	A	B
D	E	A	B	C
E	A	B	C	D

Advantages of LSD:

- 1) In field experiments, if the fertility gradient is in two directions then LSD is more efficient than RBD.
- 2) It is based on three principles of experimentation.
- 3) In LSD more number of factors are compared using less number of experimental units.
- 4) The statistical analysis of LSD is slightly complicated than that of RBD, but it is more efficient than RBD.
- 5) Local control is done with two types of grouping. Therefore experimental error will be much less in LSD compared to CRD and RBD.

DISADVANTAGES:

① number of replication is equal to the number of treatments
 ② therefore, if the number is too large

number of replication.

2) LSD is suitable for the number of treatments between 5 and 10. The design is not suitable and impracticable for more than 10 treatments.

3) If several units are missing in LSD, the statistical analysis is more difficult.

4) Randomization is also restricted within each row and each column in this design.

5) The fundamental assumption that there is no interaction between different factors may not be true in general.

A) Explain the statistical analysis of LSD data.

Suppose we have to compare m treatments for their effects, using a L.S.D. The m^2 experimental units needed for the LSD experiment are divided into m rows and m columns. The m treatments are allocated to the units of rows and columns in such a way that every treatment occurs only once in row & column.

From m^2 units we get m^2 yield.

Null hypothesis :-

H_{01} : All the Rows are homogeneous i.e.

$$H_{01}: \mu_1 = \mu_2 = \dots = \mu_m = 0$$

H_{02} : All the Columns are homogeneous

$$\text{i.e. } H_{02}: \beta_1 = \beta_2 = \dots = \beta_m = 0$$

H_{03} : All the treatments are homogeneous

$$\text{i.e. } H_{03}: \nu_1 = \nu_2 = \dots = \nu_m = 0.$$

Mathematical Model :-

Let us suppose that y_{ijk} ($i, j, k = 1, 2, \dots, m$) be the yield from the experimental unit in the i^{th} row, j^{th} column and by receiving the k^{th} treatment.

The triplet (i, j, k) assumes m^2 experimental units. In LSD of the possible m^3 values.

Let us denote the set of m^2 values by usual, we write $(i, j, k) \in S$

The linear mathematical model of LSD becomes

$$y_{ijk} = \mu + d_i + \beta_j + \nu_k + \epsilon_{ijk} \quad (i, j, k) \in S$$

where μ = general mean effect

ν_k - due to the k^{th} treatment
 ϵ_{ijk} = Error effect due to chance and
 ϵ_{ijk} iid $N(0, \sigma_e^2)$

consider $G = y_{...}$ = Total of all the m^2 observations
 $R_i = y_{i..}$ = Total of m observations in the i^{th} row
 $C_j = y_{.j.}$ = Total of m observations in the j^{th} column
 $T_k = y_{..k}$ = Total of m observations from the k^{th} treatment.

Estimation of parameters:-

The parameters $\mu, \alpha_i, \beta_j, \nu_k$ are estimated by the principle of least squares. The least squares estimates of the parameters are $\mu = \frac{G}{N} = \frac{y_{...}}{m^2} = \bar{y}_{...}$

$$\alpha_i = \bar{y}_{i..} - \bar{y}_{...} \quad i=1, 2, \dots, m$$

$$\beta_j = \bar{y}_{.j.} - \bar{y}_{...} \quad j=1, 2, \dots, m$$

$$\nu_k = \bar{y}_{..k} - \bar{y}_{...} \quad k=1, 2, \dots, m$$

$$\bar{y}_{i..} = \frac{1}{m} \sum_{(j,k)} y_{ijk} \quad \bar{y}_{.j.} = \frac{1}{m} \sum_{(i,k)} y_{ijk} \quad \bar{y}_{..k} = \frac{1}{m} \sum_{(i,j)} y_{ijk}$$

$$\bar{y}_{...} = \frac{1}{m^2} \sum_{(i,j,k) \in S} y_{ijk}$$

Various sum of squares:-

The total variation can be split into 4 parts total variation = variation b/w the rows + variation b/w the columns + variation b/w the treatments + Error variation.

Total sum of squares for corrected mean

$$\sum_{(i,j,k) \in S} (y_{ijk} - \bar{y}_{...})^2 = \sum_{(i,j,k) \in S} [(y_{ijk} - \bar{y}_{i..} - \bar{y}_{.j.} - \bar{y}_{..k} + 2\bar{y}_{...}) + (\bar{y}_{i..} - \bar{y}_{...}) + (\bar{y}_{.j.} - \bar{y}_{...}) + (\bar{y}_{..k} - \bar{y}_{...})]^2$$

$$\sum_{(i,j,k) \in S} (y_{ijk} - \bar{y}_{...})^2 = \sum_{(i,j,k) \in S} (y_{ijk} - \bar{y}_{i..} - \bar{y}_{.j.} - \bar{y}_{..k} + 2\bar{y}_{...})^2 + m \sum (\bar{y}_{i..} - \bar{y}_{...})^2 + m \sum (\bar{y}_{.j.} - \bar{y}_{...})^2 + m \sum (\bar{y}_{..k} - \bar{y}_{...})^2 + \text{Cross Product}$$

The product terms vanish since the algebraic sum of deviations from the mean is zero.

$$TSS = ESS + RSS + CSS + \text{Treatments}$$

$$S_T^2 = S_E^2 + S_R^2 + S_C^2 + S_t^2$$

For numerical calculations of various sum of squares the following formulae can be used

$$mv_k = \sum_{(i,j)} y_{ijk} - m\mu$$

$$v_k = \sum_{(i,j)} \frac{y_{ijk}}{m} - \frac{m\mu}{m}$$

$$v_k = \bar{y}_{..k} - \mu$$

$$\hat{v}_k = \bar{y}_{..k} - \bar{y}_{...}$$

Problems ①: - Analyse the LSD layout

A(12) B(14) C(11) D(10)
 B(12) C(14) D(13) A(11)
 C(9) D(15) A(13) B(10)
 D(11) A(16) B(14) C(10)

Sol

H_0 : All the rows are homogeneous.

H_{02} : All the columns are homogeneous

H_{03} : All the treatments are homogeneous

LOS: $\alpha = 5\%$

Row	1	2	3	4	R_i	R_i^2	$\sum y_{ij}^2$
1	12	14	11	10	47	2209	561
2	12	14	13	11	50	2500	630
3	9	15	13	10	47	2209	575
4	11	16	14	10	51	2601	673
C_j	44	59	51	41	195	9519	2439
C_j^2	1936	3481	2601	1681	-	9699	-

$$m=4 \quad N=m^2=4 \times 4=16 \quad G=195 \quad \sum y_{ij}^2 = 2439$$

$$\sum R_i^2 = 9519 \quad \sum C_j^2 = 9699$$

Treatments	$\sum A$	$\sum B$	$\sum C$	$\sum D$	Total
Total T_k	52	50	44	49	9541
T_k^2	2704	2500	1936	2401	9541

$$\sum T_k^2 = 9541$$

$$CF = \frac{G^2}{N} = \frac{(195)^2}{16} = 2376.5625$$

$$S_T^2 = \sum y_{ij}^2 - CF = 2439 - 2376.5625 = 62.4375$$

$$S_R^2 = \frac{\sum R_i^2}{m} - CF = \frac{9519}{4} - 2376.5625 = 3.1875$$

$$S_C^2 = \frac{\sum C_j^2}{m} - CF = \frac{9699}{4} - 2376.5625 = 48.1875$$

$$S_t^2 = \frac{\sum T_k^2}{m} - CF = \frac{9541}{4} - 2376.5625 = 8.6875$$

$$S_F^2 = S_T^2 - S_R^2 - S_C^2 - S_t^2 = 62.4375 - 3.1875 - 48.1875 - 8.6875$$

ANOVA TABLE :-

SV	df	SS	MSS	f-ratio	f _{table}
Rows	4-1=3	$S_R^2 = 3.1875$	$\bar{x}_R^2 = \frac{3.1875}{3} = 1.0625$	$f_1 = \frac{S_R^2}{S_E^2} = 2.6844$	$f_1(3,6)$ at 5% = 4.76
Columns	4-1=3	$S_C^2 = 48.1875$	$\bar{x}_C^2 = 16.0625$	$f_2 = \frac{S_C^2}{S_E^2} = 40.5823$	"
Treat	4-1=3	$S_T^2 = 8.6875$	$\bar{x}_T^2 = 2.8958$	$f_3 = \frac{S_T^2}{S_E^2} = 7.3163$	"
Error	$(4-1)(4-2) = 6$	$S_E^2 = 2.375$	$\bar{x}_E^2 = 0.3958$	—	—
Total	16-1=15	$S_T = 62.4375$	—	—	—

Inference :- Rows $f_{1cal} = 2.6844$ $f_{table} = 4.76$
 $f_{1cal} < f_{table}$ value so we accept our H_0
 All the rows are homogeneous
 Columns $f_{2cal} = 40.5823$ $f_{table} = 4.76$
 $f_{2cal} > f_{table}$ value we reject H_0 at 5% LOS
 All the columns are not homogeneous.

Treatments :-

$f_3 cal = 7.3163$ $f_3 table value = 4.76$
 $f_3 cal > f_3 table value$ we reject H_0 at 5% LOS
 All the treatments are not homogeneous.

Efficiency of LSD over RBD

(i) when Rows are taken as blocks $e = \frac{S_C^2 + (m-1)\bar{x}_E^2}{m\bar{x}_E^2}$

(ii) when Columns are taken as blocks $e = \frac{16.0625 + (4-1)0.3958}{4 \times 0.3958} = 1.7063$

Efficiency of LSD over CRD

$$e = \frac{\bar{x}_R^2 + S_C^2 + (m-1)\bar{x}_E^2}{(m+1)\bar{x}_E^2}$$

Problem no (2)

The table below gives the yield of rice in kilos observed in a field experiment carried out in 4x4

Factorial Experiments.

Q Explain factorial Experiments, merits and Demerits.

In the Experiment of CRD & RBD & LSD we were mainly concerned with the comparison of single set of treatments. Such experiments dealing with one factor only are called simple Experiments.

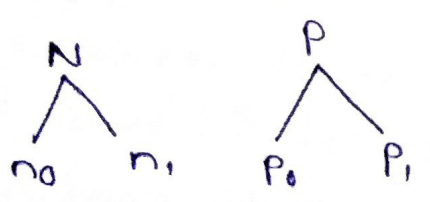
In factorial Experiments, the effects of several factors of variation are investigated. In these Experiments the treatments be all the combinations of different factors under study.

In factorial Experiments the interaction effects shall be studied and Estimated.

Here we study Experiments involving several factors & Experiments involving combinations of different levels of different factors. These factorial Experiments were considered as complex Experiments in olden days.

Def:- An Experiment in which the treatments are different combinations of different levels of several factors is called a factorial Experiments.

Example:- we have two factors N and P each with two levels (0, 1) say. Then we have factorial combinations as follows.



$N_i P_j \rightarrow$ factorial combination
i.e., treatments in factorial Experiment.

Symmetrical factorial Experiment :-

In General f factors, l-levels for each factor the factorial Experiments is l^f and this l^f factorial Experiments is called Symmetrical factorial

$f=n$ $l=2 \Rightarrow 2^n$ factorial experiments

$f=2$ $l=3 \Rightarrow 3^2$ " " "

$f=3$ $l=3 \Rightarrow 3^3$ " " "

$f=n$ $l=3 \Rightarrow 3^n$ " " "

$2^1, 2^2, \dots, 2^n$ and $3^1, 3^2, \dots, 3^n$ are symmetric factorial experiments.

Merits and demerits of factorial experiments.

Merits (advantages):—

① When compared with single factor experiments factorial experiments have more flexibility regarding the study of different levels of several factors. In the factorial experiments each factor is to be studied over a wide set of experiments.

② Interaction between the effects of various factors can be studied only in factorial experiments.

③ Factorial experiments have the advantage of economising on experimental resources.

When experiments are conducted factor by factor much more resources are required for the same precision than when they are tried in factorial experiments.

④ In factorial experiments on allocating all the combinations once we will get one complete replication. Besides such a replication there is hidden replication.

⑤ Factorial experiments are most servisable than standard designs in biological & clinical experiments.

A single factorial is easy to conduct than elaborate (extensive) many single factorial experiments at a time.

Demerits:—

(1) If n is large the experiment becomes too complex. For example if $n=10$ in 2^n factorial design will have 1024 treatments.

different levels of factors.

Explain 2^2 design.

Here we have two factors each at two levels
so that there are $2 \times 2 = 4$ treatment combinations

ii. the capital letters A and B indicates the names
two factors under study. Also let the small letters
a and b denote the levels of corresponding factors.
The first level of A and B is generally expressed
in the absence of corresponding level in the treat-

combinations.
The four treatment combinations can be enumerated
as follows.

- aob₁ | factors A and B both at first level
- a₁bo | A at second level and B at first level
- aob₂ | A at first level and B at second level.
- a₁b₁ | A and B both at second level.

These four treatment combinations can be compared
by laying out (putting)

- i) CRD with 4 treatments
- ii) RBD with 8 replicates (say) each replicate
(block) containing 4 units
- iii) LSD as a 4×4 layout.

ANOVA can be carried out accordingly.
In the above cases there are 3 d.f associated with
treatment effects.

In factorial experiments our main objective is
to carry out separate tests for the main effects
and the interaction AB.

We split the 3 d.f of treatment sum of
squares into 3 orthogonal components each with

1 d.f for main effects and 1 d.f for interaction effects of 2^2 design

Factorial experiment with $2^2 = 4$

Let $[1]$ $[a]$ $[b]$ and $[ab]$ denote the sum of n units.

$[1]$ $[a]$ $[b]$ $[ab]$ receiving the treatments $1, a, b$ and ab respectively.

Let the corresponding mean values obtained on dividing these totals by n be denoted by (1) (a) (b) and (ab) respectively. The letters A, B and AB when they refer to numbers will represent the main effects due to the factors A and B and their interaction AB respectively.

The effect of A can be represented by the difference between mean yields obtained at each level.

Thus the effect of the factor 'A' at the first level b_0 of $B = (a, b_0) - (a_0, b_0) = (a) - (1) \rightarrow \textcircled{1}$

Similarly the effect of A at the second level b_1 of $B = (a, b_1) - (a_0, b_1) = (ab) - (b) \rightarrow \textcircled{2}$

These two effects $\textcircled{1}$ and $\textcircled{2}$ are termed as simple effects of the factor A .

The average effect of A over the two levels of B is called the main effect of A over the two levels of B and is defined by

The average observed effect of A over the two levels of B is called the main effect due to A and is defined by $A = \frac{1}{2} [(ab) - (b) + (a) - (1)]$

$$A = \frac{1}{2} (a-1)(b+1) \rightarrow \textcircled{I}$$

where the right hand side to be expanded algebraically and then the treatment combination are to be replaced by treatment mean.

Similarly the effect of factor B at first level a_0 of $A = (a_0, b_1) - (a_0, b_0) = (b) - (1) \rightarrow \textcircled{I}$

the effect of B at second level a_1 of $A = (a_1, b_1) - (a_1, b_0) = (ab) - (a) \rightarrow \textcircled{2}$

$$= (a, b) - (a) \rightarrow \textcircled{2}$$

the average of $\textcircled{1}$ and $\textcircled{2}$ is known as the

$$\text{Main Effect of B} = \frac{1}{2} ((b) \cdot (1) + (ab) - (a))$$

$$= \frac{1}{2} (a+1)(b-1) \quad \text{--- II}$$

The right hand side to be expanded algebraically and the treatment combinations are to be replaced by their means.

Interaction Effect AB

If the two factors act independently of one another if means ① and ② are the estimates of the same thing.

If the two factors are not act independently the two expressions ① and ② will not be same.

Then the difference of these two numbers is the generation between the factors A and B.

The generation effect between the factors A and B defined as

$$\text{Interaction effect AB} = \frac{1}{2} [(ab) - (b) - (a) + (1)]$$

$$= \frac{1}{2} (a-1)(b-1) \quad \text{--- III}$$

Where the RHS is to be expanded algebraically and then the treatment combination are to be replaced by the corresponding treatments.

From III we notice that the interaction between B and A is as

$$\text{Interaction effect BA} = \frac{1}{2} (b-1)(a-1) \quad \text{--- III}$$

which are same as the expressions III and III a

It means the interaction does not depend on the order of the factors.

④ Explain Statistical Analysis of 2^2 factorial design.

Factorial experiments are conducted either in CRD or RBD or LSD and thus they can be analysed in the usual manner except that in this case the treatment sum of squares is split into three orthogonal components each with 1 d.f.

Suppose a 2^2 factorial experiment is conducted

by a Randomised block design (RBD) layout with 2 blocks each containing 4 plots.
 The data from the layout may be tabulated as follows

Blocks	Yields Treatments			
	a_0b_0	a_1b_0	a_0b_1	a_1b_1
1	Y_{11}	Y_{21}	Y_{31}	Y_{41}
2	Y_{12}	Y_{22}	Y_{32}	Y_{42}
...
j	Y_{1j}	Y_{2j}	Y_{3j}	Y_{4j}
...
q	Y_{1q}	Y_{2q}	Y_{3q}	Y_{4q}

The yield totals under the 4 treatments can be written as $[a_0b_0]$ $[a_1b_0]$ $[a_0b_1]$ $[a_1b_1]$
 where $[]$ denotes the sum of yields from q plots.

Null hypothesis :-

- H_{01} : All the blocks (replicates) are homogeneous
- H_{02} : The main effect A is not significant
(81)
- H_{02} : There is no significant effect of main effect A on increasing the yield.
(81)
- H_{02} : There is no significant difference due to MEA.
- H_{03} : There is no significant difference due to MeB
- H_{04} : The Interaction Effect AB is not significant

In testing these hypothesis the analysis is just like as in RBD except in the case of treatment sum of squares.

total = $4q-1$

sum of squares due to blocks = $S_{Block}^2 = \frac{\sum B_j^2}{4} - CF$
 Here $B_j = j^{th}$ block total

DOF for blocks = $q-1$

sum of squares due to treatments is divided into three components sum of squares due to main effect B, sum of squares due to interaction effect AB.

The dof for treatments are $4-1=3$ is distributed to the three components equally. To find sum of squares due to different effects we may use traditional method.

sum of squares due to main effect A = $S_A^2 = \frac{[A]^2}{4q}$

$[A] = [ab] + [a] - [b] - [1]$

sum of square due to main effect B = $S_B^2 = \frac{[B]^2}{4q}$

$[B] = [ab] - [a] + [b] - [1]$

sum of squares due to interaction effect AB

= $S_{AB}^2 = \frac{[AB]^2}{4q}$ $[AB] = [ab] - [a] - [b] + [1]$

Here $[A]$ $[B]$ $[AB]$ are the different effects totals

\therefore sum of squares due to treatments = $S_A^2 + S_B^2 + S_{AB}^2$

\therefore sum of squares due to error = $TSS - SSB - S_A^2 - S_B^2 - S_{AB}^2$

Error dof = $4q-1 - q+1 - 3 = 3(q-1)$

To test the above null hypothesis we construct the following ANOVA Table using the above calculations.

ANOVA table

Conclusion :-

for blocks, Main effect μ and Interaction effect AB .

If $F_{cal} \leq F_{table}$ value we accept our hypothesis otherwise we reject our hypothesis at 5% & 1% LOS.

5) Yates method of computing factorial effect Totals.

For the calculations of various effect Totals for 2^n factorial experiments Yates developed a special rule which enables us to avoid specific algebraic formulae.

This method is very useful if the number of factors is more. Yates method for 2^n factorial experiment consists in the following steps.

1) In the first column we write the treatment combinations in the standard order. For example 2^2 factorial experiment with 2 factor A and B. The order of treatment combination will be 1 a b ab

For 2^3 experiment with 3 factors A, B, C the order is 1 a b ab c ac bc abc.

2) In the second column we write the corresponding treatment totals
[1] [a] [b] [ab]

3) The third column can be splitted into two halves. The first half is obtained by writing the pairwise sum of the values in column 2 in the given order and second half is obtained by writing in the same order the pairwise differences of the values in the column 2. The difference should be obtained by subtracting the first value from the second value in the pair.

4) The column is constructed by adopting the

gives Total interaction effect of AB.

Table for a 2^2 factorial experiments

① Treatment combination	② Total yield from all replicates	③	④ Effect totals
1	[1]	[ab] + [1]	[ab] + [b] + [a] + [1] = G
a	[a]	[ab] + [b]	[ab] - [b] + [a] - [1] = [A]
b	[b]	[a] - [1]	[ab] + [b] - [a] - [1] = [B]
ab	[ab]	[ab] - [b]	[ab] - [b] - [a] + [1] = [AB]

Main effect of A = $\frac{[A]}{2d}$ $[A] = [ab] + [b] - [a] - [1]$

Main effect of B = $\frac{[B]}{2d}$ $[B] = [ab] + [b] - [a] - [1]$

Interaction effect of AB = $\frac{[AB]}{2d}$ $[AB] = [ab] - [a] - [b] + [1]$

Problems: In a 2^2 factorial experiment, the factors are A & B. The yield from 3 replications of the four factors are given below.

	1	b	ab	a
Replication I	25	27	32	26
Replication II	25	31	24	27
Replication III	28	23	26	32

Find the various effects of treatments using traditional method and Yates method. Analyse the design.

Traditional method.

Replication	Treatments.			
	1	a	b	ab
I	25	26	27	32
II	24	25	27	31
III	28	26	28	32
	72	77	82	95

$[1] = 72$ $[a] = 77$ $[b] = 82$ $[ab] = 95$

Main effect due to A = $\frac{(a-1)(b+1)}{2d} = \frac{[ab] + [a] - [b] - [1]}{2d}$

$$= \frac{95 + 77 - 82 - 72}{2 \times 3} = 3$$

$$\text{Main effect due to B} = \frac{(a+1)(b-1)}{2 \times 2} = \frac{[ab] - [a] - [b] - [1]}{2 \times 2}$$

$$= 4.6667$$

$$\text{Interaction effect AB} = \frac{(a-1)(b-1)}{2 \times 2} = \frac{[ab] - [a] - [b] + [1]}{2 \times 2}$$

$$= 1.3333$$

Yates method

	①	②	③	④	
1		72	149	326	= G
a		77	177	18	= [A]
b		82	5	28	= [B]
ab		95	13	8	= [AB]

Main effect due to A

$$= \frac{[A]}{2 \times 2} = \frac{18}{6} = 3$$

$$ME B = \frac{[B]}{2 \times 2} = 4.6667$$

$$I \in AB = \frac{[AB]}{2 \times 2} = \frac{8}{6} = 1.3333$$

ab, ac, bc & abc.

Here we have 7 different treatment effects namely main effect $\left. \begin{matrix} A \\ B \\ C \end{matrix} \right\} 3c_1$,

Two factor interaction effects $\left. \begin{matrix} AB \\ AC \\ BC \end{matrix} \right\} 3c_2$

Three factor interaction effects $ABC \left\} 3c_3$

$$\text{Total} = 2^3 - 1 = 7$$

2^3 factorial experiment can be performed as a CRD with 8 treatments (or) RBD with 2 replicates (say) each replicate containing 8 treatments or LSD with $m=8$ and data can be analysed accordingly.

In 2^3 experiment can be performed as a CRD with 8 treatments (or) RBD with 2 replicates (say) each replicate containing 8 treatments we split up the treatment S.S with d.f into 7 orthogonal components corresponding to the three main effects A, B and C three first order (or) two factor interaction (or) three factor) ABC each carrying d.f

Here A, B, C, AB, AC, BC, ABC etc when they refer to numbers will represent the corresponding factorial effects.

7) Explain main effects and Interaction effects in 2^3 factorial experiments

Suppose the factorial experiment with $2^3=8$ treatment is conducted in 2 blocks (replicates) Let $[1][a][b][c][ab][ac][bc][abc]$ denotes the total yields of the units (plots) receiving the treatments 1, a, b, c, ab, ac, bc, abc respectively.

Let the number of observations in each treatment be equal. They refer to the number of observations due to factors A, B, C and ABC respectively. The simple effects of A are calculated in the following manner.

level of B	level of C
b ₀	c ₀
b ₀	c ₁
b ₁	c ₀
b ₁	c ₁

Simple effects of A

$$(a_1 b_0 c_0) - (a_0 b_0 c_0) = (a_1) - (a_0)$$

$$(a_1 b_0 c_1) - (a_0 b_0 c_1) = (a_1) - (a_0)$$

$$(a_1 b_1 c_0) - (a_0 b_1 c_0) = (a_1) - (a_0)$$

$$(a_1 b_1 c_1) - (a_0 b_1 c_1) = (a_1) - (a_0)$$

The main effect of A is the average of the above 4 simple effects.

Main effects of A = $\frac{1}{4} [(a_1) - (a_0) + (a_1) - (a_0) + (a_1) - (a_0) + (a_1) - (a_0)]$
 $= \frac{1}{4} [4(a_1 - a_0)] \rightarrow \text{①}$

ME of B = $\frac{1}{4} [(a_1) - (a_0) + (a_1) - (a_0) + (a_1) - (a_0) + (a_1) - (a_0)] \rightarrow \text{②}$

ME of C = $\frac{1}{4} [(a_1) - (a_0) + (a_1) - (a_0) + (a_1) - (a_0) + (a_1) - (a_0)] \rightarrow \text{③}$

By expanding the expressions ①, ② and ③ algebraically and the treatment combinations are replaced by the corresponding average yield then we get the main effect of A, main effect of B & Main effect of C.

First order interaction effects:-

The average effect of A (at one level of C) at the level b₀ of B is $\frac{1}{2} [(a_1) - (a_0) + (a_1) - (a_0)]$

The average effect of A (at one level of C) at the level b₁ of B is $\frac{1}{2} [(a_1) - (a_0) + (a_1) - (a_0)]$

∴ the interaction effect of AB is given by half the difference b/w the average effect of A at the second & first level of B is given by.

Interaction effect of AB = $\frac{1}{4} [(a_1) - (a_0) + (a_1) - (a_0) - (a_1) - (a_0) + (a_1) - (a_0)]$

Interaction effect of AB = $\frac{(a_1 - a_0)(b_1 - b_0)(c_1 - c_0)}{4} \rightarrow \text{④}$

Similarly we can obtain expression for the interaction Bc and Ac.

Interaction effect of Bc = $\frac{(a_1 - a_0)(b_1 - b_0)(c_1 - c_0)}{4} \rightarrow \text{⑤}$

Interaction effect of Ac = $\frac{(a_1 - a_0)(b_1 - b_0)(c_1 - c_0)}{4} \rightarrow \text{⑥}$

Second order interaction effects:-

We obtained the expression for the second order interaction effects:-

at the level c_0 at c is
 $\frac{1}{2} \{ (ab) - (b) - (a) - 1 \}$

of AB at the level c_1 of c is

$$\frac{1}{2} \{ (a, b, c) - (bc) - (ac) + 1 \}$$

Effect of AB with c is

$$1) \text{ Effect of } ABC = \frac{1}{4} \{ (abc) - (bc) - (ac) + (c) - (ab) + (b) + (a) - 1 \}$$

$$c = \frac{1}{4} \{ (abc) - (ab) - (ac) - (bc) + (a) + (b) + (c) \}$$

$$= \frac{(a-1)(b-1)(c-1)}{4} \rightarrow \text{④}$$

Adding ④ ⑤ ⑥ & ⑦ algebraically and substituting treatment combinations by the corresponding mean yields then we get the interaction = BC & ABC.

Analysis of 2^3 design
 (8)

3-factorial experiment.

Factorial experiments are conducted either in R.D. They can be analysed in the usual way. The treatment sum of squares is decomposed into components each with 1 d.f.

A 3-factorial experiment can be performed in RBD. It has 8 treatments i.e. 1, a, b, ab, c, ac, bc, abc replicated 8 times.

The yields are tabulated as follows

Block	j	1	total
1	y_{11}	y_{12}	[1]
2	y_{21}	y_{22}	[a]
3	y_{31}	y_{32}	[b]
4	y_{41}	y_{42}	[ab]
5	y_{51}	y_{52}	[c]
6	y_{61}	y_{62}	[ac]
7	y_{71}	y_{72}	[bc]
8	y_{81}	y_{82}	[abc]
B ₂	B_j	B_2	G

... units from each block receiving all

$$y_{ij} \quad i=1,2,\dots,8$$

$$= \frac{G_i^2}{8q} \quad i=1,4,\dots,8$$

$$S y_{ij} \quad j=1,4,\dots,8$$

$$= S_T^2 = S y_{ij}^2 - CF = \frac{EBq^2}{8} - CF$$

th block total $j=1,2,\dots,8$ and is divided into 7 parts and as follows.

we to main effect $A = \frac{|A|^2}{8q}$

$$+ (bc) + (a) - (b) - (c) - (1)$$

we to main effect $B = \frac{|B|^2}{8q}$

$$+ (bc) - (a) + (b) + (c) - (1)$$

we to main effect $C = \frac{|C|^2}{8q}$

$$+ (bc) - (a) - (b) + (c) - (1)$$

we to interaction effect $AB = \frac{|AB|^2}{8q}$

$$+ (bc) - (a) - (b) + (c) + (1)$$

we to interaction effect $AC = \frac{|AC|^2}{8q}$

$$+ (bc) + (a) - (b) - (c) + (1)$$

we to interaction effect of $ABC = \frac{|ABC|^2}{8q}$

$$- (bc) + (a) + (b) + (c) - (1)$$

of squares are also calculated easily
ent sum of squares is the sum of

$$S_{AB}^2 + S_{AC}^2 + S_{BC}^2 + S_{ABC}^2$$

we to error

$$S_t$$

$$S_t$$

$$\text{total} = 8q = 1$$

Source of variation	df	SS	MS	F-cal	F-table
Blocks	2-1	S^2_{Block}	$\lambda^2_{Block} = \frac{S^2_{Block}}{2-1}$	$f_1 = \frac{\lambda^2_{Block}}{\lambda^2_E}$	$f_1(2-1) \{ (2-1) \}$
MEA	1	S^2_A	$\lambda^2_A = \frac{S^2_A}{1}$	$f_2 = \frac{\lambda^2_A}{\lambda^2_E}$	$F_2(1, 7(2-1))$
MEB	1	S^2_B	$\lambda^2_B = \frac{S^2_B}{1}$	$f_3 = \frac{\lambda^2_B}{\lambda^2_E}$	"
MEC	1	S^2_C	$\lambda^2_C = \frac{S^2_C}{1}$	$f_4 = \frac{\lambda^2_C}{\lambda^2_E}$	"
IFAB	1	S^2_{AB}	$\lambda^2_{AB} = \frac{S^2_{AB}}{1}$	$f_5 = \frac{\lambda^2_{AB}}{\lambda^2_E}$	"
IFAC	1	S^2_{AC}	$\lambda^2_{AC} = \frac{S^2_{AC}}{1}$	$f_6 = \frac{\lambda^2_{AC}}{\lambda^2_E}$	"
IFBC	1	S^2_{BC}	$\lambda^2_{BC} = \frac{S^2_{BC}}{1}$	$f_7 = \frac{\lambda^2_{BC}}{\lambda^2_E}$	"
IFABC	1	S^2_{ABC}	$\lambda^2_{ABC} = \frac{S^2_{ABC}}{1}$	$f_8 = \frac{\lambda^2_{ABC}}{\lambda^2_E}$	"
Error	$7(2-1)$	S^2_E	$\lambda^2_E = \frac{S^2_E}{7(2-1)}$	—	—
Total	$8(2-1)$	S^2_T	—	—	—

Conclusion :- We compare the f-calculated values with the corresponding f-table values at specified LOS and draw the conclusion accordingly.

2) Yates method for a 2^3 factorial experiment. Yates developed a specific algebraic formula. This is very useful if the no. of factors is more we construct the Yates as described below.

- In the first column we write the treatment combination in the standard order i.e., 1, a, b, ab, c, ac, bc, abc.
- In the second column we write the corresponding treatment totals from all the replications. i.e. [1][a][b][ab], [c][ac][bc][abc]
- The entries in the third column is splitted into two halves, the first half is obtained by writing the pairwise sums of column 2 in the given order. The second half is obtained by subtracting the first value from the second value of the pair.
- The fourth column is constructed by adopting the same procedure on column 3
- The fifth column is also constructed by adopting the same procedure as explained in step 3 on column 4. The column gives the Grand total & Total effect of A i.e

total. Effect of B i.e [B]. Total effect interaction effect of effect of C i.e (c) total interaction effect of BC and

	(1)	$(a+b)+c = x_1$	$(abc)+bc+(ac)+(c)$	$x_2+x_1 = y_1$	$[a][1] = G$	$y_2+y_2 = \dots G$
a	[a]	$[ab]+[b] = x_2$	$(abc)+bc+(ac)+(c)$	$x_4+x_3 = y_2$	$(abc)+(bc)+(ac)+(c)+(ab)$	$y_4+y_3 = (A)$
b	[b]	$[ac]+[c] = x_3$	$(abc)+bc+(a)+(1)$	$x_6+x_5 = y_3$	(abc)	$y_6+y_5 = (B)$
ab	[ab]	$(abc)+(bc) = x_4$	$(abc)-(bc)+(ac)-(c)$	$x_8+x_7 = y_4$		$y_8+y_7 = (AB)$
c	[c]	$(a)-(1) = x_5$	$(ab)+(b)-(a)-(1)$	$x_2-x_1 = y_5$		$y_2-y_1 = (C)$
ac	[ac]	$(ab)-(b) = x_6$	$(abc)+(bc)-(ac)-(c)$	$x_4-x_3 = y_6$		$y_4-y_3 = (Ac)$
bc	[bc]	$(ac)-(c) = x_7$	$(ab)-(b)-(a)+(1)$	$x_6-x_5 = y_7$		$y_6-y_5 = (Bc)$
(abc)	(abc)	$(abc)-(bc) = x_8$	$(abc)-(bc)-(ac)+(c)$	$x_8-x_7 = y_8$		$y_8-y_7 = (ABC)$

Main effects and interaction effect can be obtained by dividing the effect totals by 4g.

Critical difference :-

If there is a significant difference in between the treatments then we would be interested to find out which pair of treatments differ significantly. For the instead of calculating students 't' for different pairs treatment means, we calculate the test significant difference of the given level of significance. The least difference is known as the critical difference (CD).

CD at α level of significance is given by

$$CD = SE(\bar{x}_1 - \bar{x}_2) \times t_{\alpha\%} \text{ for error d.f}$$

For example :- $H_0: \mu_1 = \mu_2$ i.e., two-treatment means do not differ significantly. If H_0 is rejected, then H_0 is accepted i.e. $H_1: \mu_1 \neq \mu_2$ may be accepted. Now we have to determine which pair of treatment differ significantly. To obtain this we have to calculate a student 't' statistic for every pair. So instead of this we compute C.D as follows

$$CD(\bar{x}_1 - \bar{x}_2) = CD = SE(\bar{x}_1 - \bar{x}_2) \times t_{\alpha\%} \text{ for error d.f}$$

$$= \sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \times t_{\alpha\%} \text{ for error d.f}$$

$E(\sigma^2) = \frac{\sigma^2}{n}$ and if each treatment replicated n

times.

... n = 1, 2, ... k then $CD = t_{\alpha} \sqrt{\frac{\sigma^2}{n}}$ for error d.f

