

Cluster Sampling

In Random Sampling, the popln has been divided into a finite no. of distinct and identifiable units defined as sampling units.

The smallest unit into which the popln can be eliminated is called an element of the popln.

A group of such elements is known as cluster. When the sampling units in the cluster the procedure is called cluster sampling. If the entire area containing popln under study is divided into smaller segments and each element in the popln belongs to one and only one segment, the procedure is called area sampling.

Identification and location of element requires considerable time area an element has been located, the time taken for surveying a few neighboring elements is small. Thus the main fun. in cluster sampling is to specify the clusters or to divided the popln into appropriate clusters. The clusters are generally made up of neighboring elements and the elements within cluster is to have similar characteristic. as the sample rule, the no. of elements in the cluster should be small, and no. of cluster should be large.

After dividing in the popn into specified clusters, the required no. of clusters and can be selected either by equal or unequal Prob. of selection. All the elements in selected clusters are enumerated.

For a given no. of sampling units, cluster sampling is more convenient and less cost the advantages of cluster sampling

Advantages:-

1) Collection of data for neighbouring elements is faster, cheaper, easier and operationally more convenient than observing unit spread over the region.

2) It is less costly than SRS due to the saving of time in journey, identification, etc.

3) When the sampling frame of elements may not be readily available.

Disadvantages:-

1) Even if the sample frame is available it could be expensive to conduct an enquiry on a SRS of elements.

2) From the point of view of statistical efficiency, cluster sampling is generally less efficient than SRS due to the tendency of units in cluster.

[The efficiency of cluster sampling is likely to decrease with increase in cluster size.

For a given sample size, the smaller sampling unit will bring more precise results than the large sampling unit. The last efficiency may be balanced by the reduction in the cost, the efficiency per unit

Cost may be more in cluster sampling than in SRS.

The selection of cluster can be random or by first selecting a unit call a key unit at random and then randomly taking the required no. of neighbouring cluster (to estimate the maize production, the cluster of three villages formed by first selecting a key village in the random and then taking two more villages from a block of some specified area).

Notations:-

Suppose the popln consists of N clusters, each of M elements and that a sample of n clusters is drawn by the method of SRS.

N = no. of clusters in the popln.

n = no. of clusters in the sample.

M = no. of elements in the cluster.

y_{ij} = the value of the characteristic under study for the j^{th} elements ($j=1, 2, \dots, M$) in the i^{th} cluster ($i=1, 2, \dots, N$)

$\bar{y}_i = \frac{\sum_j y_{ij}}{M}$ = the mean per element of the i^{th} cluster.

$\bar{y}_n = \frac{\sum_i \bar{y}_i}{n}$ = the mean of cluster means in a sample of n cluster.

$\bar{Y}_N = \frac{\sum_i \bar{y}_i}{N}$ = the mean of cluster means in the popln.

$\bar{Y} = \frac{\sum_i \sum_j y_{ij}}{NM}$ = the mean per element in the popln.

$S_f^2 = \frac{\sum_j (y_{ij} - \bar{y}_i)^2}{(M-1)}$ = the mean square b/w

elements within the i th cluster ($i = 1, 2, \dots, M$)

$S_w^2 = \frac{1}{N} \sum_{i=1}^M S_i^2 / N =$ the mean square within cluster,
(w for within)

$S_b^2 = \frac{1}{M} \sum_{i=1}^M (\bar{y}_i - \bar{y})^2 / (M-1) =$ the mean square b/w clusters means in the pop'n (b for b/w)

$S^2 = \frac{1}{N} \sum_{i=1}^M \sum_{j=1}^M (y_{ij} - \bar{y})^2 / (NM-1) =$ the mean square b/w elements in the pop'n.

$$\rho = \frac{E(y_{ij} - \bar{y})(y_{ik} - \bar{y})}{E(y_{ij} - \bar{y})^2}$$

$$= \frac{\sum_{i=1}^M \sum_{j=1}^M \sum_{k=1}^M (y_{ij} - \bar{y})(y_{ik} - \bar{y})}{(M-1)(NM-1) S^2}$$

the intra cluster correlation co-efficient b/w elements within clusters.

Equal Cluster Sampling :-

Estimator of mean and its variance :-

No new principles are involved in making estimator when a prob. sample of 'n' equal sized clusters has been taken and each cluster is enumerated. completely since the clusters are of equal size, it is clear that $\bar{y}_N = \bar{y}$. For the sampling variance of the estimator \bar{y}_n in the form developed by Hansen and Hurvitz (1942), we shall begin by proving the following theorem.

Theorem:

In simple random sampling, with of n clusters each containing M elements from a population of M clusters, the sample mean \bar{y}_n is an unbiased estimator of \bar{Y} and its variance is given by

$$V(\bar{y}_n) = \frac{(1-f)}{n} S_b^2 \\ \approx \frac{(1-f)}{n} S_M^2 [1 + (M-1)\rho]$$

where ρ is the intra cluster correlation coefficient

Proof

W.k.T, $\bar{y}_n = \sum y_i / n$

$$E(\bar{y}_n) = \frac{1}{n} \sum E(y_i) \quad \because E(y_i) = \sum P_i y_i$$

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

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

Thus, \bar{y}_n is an unbiased estimator of \bar{Y} .

Also we know that SRSWOR $V(\bar{y}) = \frac{1-f}{n} \cdot s^2$,

with this result we can write

$$V(\bar{y}_n) = \frac{(1-f)}{n} S_b^2 \\ = \frac{(1-f)}{n} \sum_{i=1}^N \frac{(y_i - \bar{Y}_N)^2}{N-1} \quad \leftarrow \textcircled{1}$$

After substituting the value of

$$\rho = \frac{\sum \sum (y_{ij} - \bar{Y})(y_{ik} - \bar{Y})}{(M-1)(NM-1) \cdot s^2}$$

in the above relation, we get

$$\sum_{i=1}^N (\bar{y}_i - \bar{Y}_N)^2 = \frac{1}{M^2} \sum_i \sum_j (\bar{y}_{ij} - \bar{y})^2 + \frac{1}{M^2} \sum_i \sum_j \sum_{j \neq k} (\bar{y}_{ij} - \bar{y})(\bar{y}_{ik} - \bar{y})$$

Sub in ①, we get

$$v(\bar{Y}_n) = \left(\frac{1-f}{n}\right) \cdot \frac{NM-1}{M^2(N-1)} s^2 [1 + (M-1)\rho]$$

$$\approx \left(\frac{1-f}{nM}\right) s^2 [1 + (M-1)\rho] \quad \text{for large } N$$

It is shown that the variance in cluster sampling depends on the no. of cluster in the ~~cluster~~ sample, the size of the cluster, the inter cluster correlation coefficient ρ and the variance s^2 . If $M=1$ it gives the sampling variance of a simple random sample of nM elements taken individually.

Corollary 1:

In simple random sampling work of n clusters each containing M elements from a population of N clusters. the population total is estimated unbiasedly of

$$\bar{Y}_c = NM \bar{Y}_n$$

with its sampling variance

$$v(\hat{Y}_c) = (1-f) \frac{N^2 M^2 s_b^2}{n}$$

$$\approx MN^2 \frac{(1-f)}{n} s^2 [1 + (M-1)\rho]$$

Corollary 2:

Instead of sampling in clusters, if a simple random sample of nM elements be directly taken from the population of NM elements

the variance of the mean per element would be

$$v(\bar{y}_n) = \frac{(1-f)}{nM} s^2$$

Corollary 3:

In simple random sampling, w.o.r. an estimator of variance of \bar{y}_n is given by

$$v(\bar{y}_n) = \frac{(1-f)}{n} s_b^2$$

$$\text{where, } s_b^2 = \frac{\sum_i (\bar{y}_i - \bar{y}_n)^2}{(n-1)}$$

Corollary 4:

In simple random sampling w.o.r. the cluster mean \hat{y}_n is an unbiased estimator of \bar{Y} and its variance is given by

$$v(\hat{y}_n) = \frac{s_b^2}{n}$$

Relative Efficiency of Cluster Sampling:

In sampling nM elements from the population by the simple random sampling, the variance of the sample mean \hat{y} is given by

$$v(\hat{y}) = \frac{(1-f) s^2}{nM}$$

The relative efficiency of cluster sampling compared with simple random sample is

$$\text{Rel. Efficiency} = \frac{V_{SR}(\bar{y})}{V_e(\hat{y}_n)} = \frac{s^2}{M s_b^2}$$

This shows that the efficiency of cluster sampling increases as the mean square between clusters decreases.

$$\text{Also } (N-1)MS_b^2 = (NM-1)s^2 - N(M-1)s_w^2$$

Therefore, relative Efficiency will increase with increase in the mean square within clusters.

Another way to express relative efficiency is to make use of the concept of intra cluster correlation coefficient.

For large N , the relative efficiency of cluster sampling in terms of intracluster coefficient ρ is given by,

$$\text{Rel Efficiency (E)} = [1 + (M-1)\rho]^{-1}$$

The Efficiency can easily be calculated by estimating the value of ρ from the sample.

An estimator of ρ can be written as

$$\hat{\rho} = \frac{(n-1)MS_b^2 - n s_w^2}{(n-1)MS_b^2 - n(M-1)s_w^2}$$

$$\text{where, } s_w^2 = \sum_i \sum_j \frac{(y_{ij} - \bar{y}_i)^2}{n(M-1)}$$

Thus for large N , an estimate of the relative efficiency of the cluster sampling can be written as

$$\text{Est. Rel Efficiency (e)} = \frac{1}{M} + \frac{(M-1)s_w^2}{M^2 s_b^2}$$

Accordingly ρ can be estimated by

$$\hat{\rho} = \frac{(1-e)}{(M-1)e}$$

For a sample of n clusters, the results in the form of analysis of variance can be written in the form

Analysis of variance of sample of n clusters

Source of variation	Degrees of freedom	Mean Square
Between clusters	$n-1$	$MS_b^2 = \frac{\sum_i (\bar{y}_i - \bar{y}_n)^2}{(n-1)}$
Within clusters	$n(M-1)$	$s_w^2 = \frac{\sum_i \sum_j (\bar{y}_{ij} - \bar{y}_i)^2}{n(M-1)}$
Total sample	$n(M-1)$	$s^2 = \frac{\sum_i \sum_j (\bar{y}_{ij} - \bar{y}_n)^2}{(nM-1)}$ $= \frac{(n-1)MS_b^2 + n(M-1)s_w^2}{(nM-1)}$

It should be noted clearly that in a random sample of n clusters, s_b^2 and s_w^2 will provide unbiased estimates of S_b^2 and S_w^2 respectively, while s^2 will not be an unbiased estimate of S^2 . The reason is that a sample of nM elements is not taken randomly from the population of NM elements.

A unbiased estimate may be obtained easily by substituting the values in the relation,

When

$$\hat{\sigma}^2 = \frac{(N-1) M S_b^2 + N(M-1) S_w^2}{(NM-1)}$$

Optimum cluster size

- For a given sample size the sampling variance increase with cluster size and decreases with increasing number of clusters.
- On otherhand, the cost decrease with the cluster size and increases with the number of clusters.
- It is necessary to determine a balancing point by finding out the optimum cluster size and the number of cluster in the sample which can minimise the sampling variance for a given cost or alternatively, minimise the cost for a fixed variance.
- we shall discuss the problem of optimum size of cluster for which maximum precision is attained with a given cost.

The cost of a survey, apart from overhead cost, will be made up of two components.

- i) Cost due to expended in enumerating the elements in the sample and in travelling within cluster, which is proportional to the number of elements in the sample
- ii) Cost due to expended on the travelling between clusters, it has been shown empirically that the expected value of minimum distance between n points located at random is proportional to $n^{1/2}$.

The Cost of the survey can be expressed as

$$C = c_1 n M + c_2 n^{1/2}$$

where,

c_1 = Cost of enumerating an element including the cost of travel between unit within the cluster.

c_2 = Cost per unit distance travelled between clusters.

Variance of estimator \bar{y}_n is based on a sample of n clusters of size m each,

$$V(\bar{y}_n) = \frac{(1-f) S_b^2}{n}$$

S_b^2 can be obtained if we know

- i) Variance ' s^2 ' between all elements in the population
- ii) The variance ' S_w^2 ' within cluster.

Hence, an approach has always been made to predict S_w^2 as it is affected by the cluster size while s^2 remains unchanged by it.

On the basis of several agriculture surveys, they observed that S_w^2 appears to bear a relation with m , which can be written by the empirical relation

$$S_w^2 = a m^b \quad (b > 0)$$

where, a and b are positive constants and find by survey from ANOVA

$$S_b^2 = \frac{(MN-1)S^2 - N(M-1)S_w^2}{M(N-1)}$$

Substituting value of S_w^2

$$S_b^2 = S^2 - (M-1) a m^{b-1}$$

$$V(\bar{y}_n) = \frac{1}{n} [S^2 - (M-1) a M^{b-1}]$$

To calculate the values of n and M by minimizing V for given C . To minimize $\phi = C + \lambda V$

Differentiate ϕ w.r.t n and M and equate to zero

$$n = \left[\frac{-C_2 + (C_2^2 + 4C_1 C M)^{1/2}}{2C_1 M} \right]^2$$

' M ' can be obtained from the equation.

$$\frac{a M^{b-1} [bM - (b-1)]}{S^2 - (b-1) a M^{b-1}} = 1 - \left(1 + \frac{4C_1 C M}{C_2^2} \right)^{-1/2}$$

By iterative method. On substituting the value of M we can obtain the optimum value of ' n '.

M will change according to change in

C_1, C_2 and C such that $C_1 C M / C_2^2$ is nearly constant.

Conclusions:

The optimum size of the unit will be smaller when

- i) The cost of enumeration of an element increase.
- ii) The cost of travel between units decrease;
- iii) The cost of the survey is sufficiently large

CLUSTER SAMPLING FOR PROPORTIONS:

Suppose it is required to estimate the proportion of elements belonging to a specified class when the population consists of N clusters, each of size M , and a random sample, w.o.r, of n clusters is selected. Suppose the M elements in any cluster can be classified into two classes, to the class, and 0 otherwise. It can be easily seen that

$$y_{ij} = 1, \text{ otherwise } 0$$

$$P_i = \frac{a_i}{M}$$

a_i being the number of elements in the i^{th} cluster belonging to the specified class.

An unbiased estimator of the population proportion

$$P = \frac{\sum_{i=1}^N a_i}{NM}$$

$$\therefore P_i = \frac{a_i}{M}$$

$$= \sum_{i=1}^N \frac{P_i}{N}$$

$$\hat{P}_c = \sum_{i=1}^n \frac{P_i}{n} = \bar{P}$$

where P_i is the proportion of elements belonging to the specified class in the i^{th} cluster of the sample

The sampling variance of \hat{P}_c is given by

$$V(\hat{P}_c) = \frac{(1-f)}{n} S_b^2$$

where S_b^2 is the variance between cluster proportions and is given by

$$M S_b^2 = \sum_{i=1}^N \frac{(P_i - P)^2}{N} = PQ - \sum_{i=1}^N \frac{P_i Q_i}{N}$$

For large N , we have $S^2 \cong S_b^2 + S_w^2 = PQ$

and the within - variance S_w^2 is given by $\sum_{i=1}^N P_i Q_i / N$. Hence, the intracluster Correlation Coefficient ρ can be written as

$$\rho = 1 - \frac{\sum_{i=1}^N M P_i Q_i}{CM - DP}$$

\therefore The Sampling Variance, in terms of the intracluster correlation coefficient, can be expressed as

$$V(\hat{P}_c) = \frac{(1-f)NPQ}{(N-1)NM} [1 + CM - DP] \rightarrow (*)$$

An estimator of the total number of units belonging to the specified category can be obtained by multiplying \hat{P}_c by NM and the expression for its sampling variance is $N^2 M^2$ times that given by (*).

If a simple random sample of nM elements could be taken, the variance of the sample proportion \hat{P} would be given by

$$V(\hat{P}) = (1-f) \frac{NPQ}{nM(N-1)} \rightarrow (**)$$

The efficiency of Cluster Sampling as compared to simple random sampling, var. can be obtained as

$$RE = \frac{(N-1)}{(NM-1)} \frac{NPQ}{NPQ - \sum_{i=1}^N P_i Q_i}$$

An estimator of $V(\hat{P}_c)$ is given by

$$\begin{aligned}v(\hat{P}_c) &= (1-f) \frac{S_b^2}{n} \\ &= (1-f) \frac{\sum_{i=1}^n (P_i - \bar{P})^2}{n(n-1)} \rightarrow (***)\end{aligned}$$

An estimator for the Sampling Variance of the total number of units belonging to a specified class can be obtained by multiplying $N^2 M^2$ to the value in relation (***)

Unequal Cluster Sampling: Estimators of Mean & their

Variance:

Consider the case when the size of all the cluster is the same. But in many practical situation cluster, sizes vary. Now we shall discuss the case of ~~unequal~~ unequal cluster.

Suppose N cluster. let i th cluster of M_i element ($i=1, 2, \dots, N$) and $\sum_i^N M_i = M_0$. The population mean \bar{y} defined by

$$\bar{y} = \frac{\sum_i^N \sum_j^{M_i} y_{ij}}{\sum_i^N M_i} = \frac{\sum_i^N M_i \bar{y}_i}{\sum_i^N M_i} = \frac{\sum_i^N M_i \bar{y}_i}{M_0}$$

\bar{y}_i = mean per element of the i th cluster.

The Pooled mean of the cluster means in the popn,

$$\bar{y}_N = \sum_i^N \bar{y}_i / N.$$

Let a random sample, n or, of n cluster be drawn and all elements of the cluster surveyed.

$$i) \bar{y}_n = \sum_i^n \bar{y}_i / n$$

$$ii) \bar{y}'_n = \sum_i^n \frac{M_i \bar{y}_i}{\sum_i^n M_i}$$

$$iii) \bar{y}^*_n = \frac{N}{M_0} \frac{\sum_i^n M_i \bar{y}_i}{n} = \frac{N}{n M_0} \sum_i^n M_i \bar{y}_i$$

$$[\therefore \bar{M} = \frac{\sum_i^N M_i}{N} = \frac{1}{N} \sum_i^N M_i]$$

$$\bar{M} = \frac{M_0}{N}$$

$$\frac{1}{\bar{M}} = \frac{N}{M_0}]$$

Theorem.

S.T the simple arithmetic mean given by relation (i) is not an unbiased estimator. The bias and sampling variance of the estimator are given by.

$$B(\bar{y}_n) = - \frac{\text{COV}(\bar{y}_i, m_i)}{\bar{M}} \quad \text{and}$$

$$V(\bar{y}_n) = \frac{1-f}{n} s^2_b$$

$$\text{Where } s^2_b = \frac{\sum_i^N (\bar{y}_i - \bar{Y}_N)^2}{N-1}$$

pf

To Prove \bar{y}_n is not unbiased

$$\begin{aligned} E(\bar{y}_n) &= E\left[\frac{1}{n} \sum_i^h \bar{y}_i\right] \\ &= \frac{1}{n} E\left[\sum_i^h \bar{y}_i\right] \\ &= \frac{1}{n} \sum_{i=1}^n E(\bar{y}_i) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{1}{N} \sum_{i=1}^N \bar{y}_i \\ &= \frac{\bar{y}_N}{N} \neq \bar{Y} \end{aligned}$$

Thus, \bar{y}_n is a biased estimator of the Population mean \bar{Y} .

The bias of the Estimator is given by

$$B(\bar{y}_n) = E(\bar{y}_n) - \bar{Y} = \frac{\sum_i^N \bar{y}_i}{N} - \frac{\sum_i^N m_i \bar{y}_i}{N\bar{M}}$$

$$= \frac{\sum_i^N (\bar{y}_i) (\bar{m}) - \sum_i^N m_i \bar{y}_i}{N \bar{m}}$$

$$= \frac{\sum_i^N (\bar{m} - m_i) \bar{y}_i}{N \bar{m}}$$

$$= \frac{-\text{cov}(y_i, m_i)}{\bar{m}}$$

For a Popn in which m_i 's do not appreciably vary from one cluster to another, the bias may be materially significant. If m_i and y_i are uncorrelated the bias is zero and \bar{y}_n is an unbiased estimate in the case

This shows bias is expected to be small when m_i and \bar{y}_i are not highly correlated. In such a case, it is advisable to use this estimator. Its sampling variance is given by

$$V(\bar{y}_n) = E[\bar{y}_n - E(\bar{y}_n)]^2 = E(\bar{y}_n - \bar{y}_n)^2$$

$$= \frac{(1-f)}{n} \frac{\sum_i^N (\bar{y}_i - \bar{y}_n)^2}{(N-1)} = \frac{(1-f)}{n} s^2_b$$

Corollary An unbiased estimator of $V(\bar{y}_n)$ is given by

$$v(\bar{y}_n) = \frac{(1-f)}{n} s^2_b$$

where $s^2_b = \frac{\sum_i^N (\bar{y}_i - \bar{y}_n)^2}{(n-1)}$

Theorem

Show that the estimator of the mean given by relation (iii) is unbiased and its variance is given by.

$$V(\bar{y}_n^*) = \frac{(1-f)}{n} S_b^{*2}$$

where

$$S_b^{*2} = \frac{1}{(N-1)} \sum_i^N \left(\frac{M_i}{M} \bar{y}_i - \bar{y} \right)^2$$

Pr

Let us consider the estimator $\bar{y}_n^* = \frac{\sum_i^n M_i \bar{y}_i}{nM}$

$$E(\bar{y}_n^*) = E\left(\frac{\sum_i^n M_i \bar{y}_i}{nM}\right) = \frac{\sum_i^n E(M_i \bar{y}_i)}{nM} = \bar{y}$$

Hence \bar{y}_n^* is an unbiased estimator. The sampling

variance of the estimator is given by.

$$\begin{aligned} V(\bar{y}_n^*) &= \frac{(1-f)}{n} \sum_i^n \left(\frac{M_i}{M} \bar{y}_i - \bar{y} \right)^2 \\ &= \frac{(1-f)}{n} S_b^{*2} \end{aligned}$$

If may be noticed that the estimator \bar{y}_n^* will often be less precise this occurs because the variance depends upon the variation of the product $M_i \bar{y}_i$ and is likely to be larger than \bar{y}_n unless \bar{y}_i and M_i vary in such a way that

their product is almost constant

Corollary An unbiased estimator of $V(\bar{y}^*_n)$ is given by,

$$v(\bar{y}^*_n) = \frac{(1-f)}{n} s_b^{*2}$$

$$s_b^{*2} = \frac{\sum_i^h \left(\frac{M_i \bar{y}_i}{\bar{M}} - \bar{y}^*_n \right)^2}{(h-1)}$$

Where

Relative efficiency of unequal cluster sampling^o

In a number of situations it is easier to take some naturally formed groups of elements usually in such cases cluster size would be unequal for example villages which are groups of households or households which are groups of persons are usually taken as clusters for the purpose of sampling on account of operational convenience thus unequal cluster sampling^o is the most practical situation and its relative efficiency with respect to simple random sampling^o should be worked out.

In unequal cluster sampling^o the total number of elements $\sum_i^h M_i$ in the sample is a random variable with expected value nM .

If an equivalent simple random sample of size $n\bar{M}$ had been selected directly from the population of $N\bar{M}$ elements the variance of the mean per element would be given by

$$V_{SR}(\bar{y}) = \frac{(N\bar{M} - n\bar{M})}{N\bar{M} n\bar{M}} S^2 = \frac{(1-f)}{n\bar{M}} S^2$$

Comparing this with the value given by,

$$R.E = \frac{S^2}{M s_b^2}$$

Hence it is observed that the efficiency increases as the variation between clusters decreases. In general cluster sampling will be efficient only when the variation between clusters is as small as possible.