

Systematic random Sampling:-

Systematic Sampling is a commonly employed technique of the complete and upto data list of the sampling unit is available. Let us suppose that N Sampling units are serially numbered from 1 to N in some order and a sample of size n is to be drawn such that

$$N = nk \Rightarrow n = N/k$$

Variance of the estimate mean:-

$$\text{var}(\bar{y}_{sy}) = \frac{N-1}{N} S^2 - \frac{(n-1)}{N} k S_{wsy}^2$$

where

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

Proof:-

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

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

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

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

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

$$(N-1)S^2 = k(n-1)S^2_{wsy} + nk \text{ var}(\bar{y}_{sys}) \quad (\because nk=N)$$

$$\text{var}(\bar{y}_{sys}) = \frac{N-1}{N} S^2 - \frac{k(n-1)}{N} S^2_{wsy}$$

Theorem :-

$$\text{var}(\bar{y}_{sys}) = \frac{nk-1}{nk} \frac{S^2}{n} (1+(n-1)p)$$

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

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

Proof :-

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

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

$$nk^2 \text{ var}(\bar{y}_{sys}) = \sum_{i=1}^k \left[\sum_{j=1}^n (y_{ij} - \bar{y}_{i.})^2 \right]$$

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

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

$$= (nk-1) s^2 + (n-1) (nk-1) s^2 P$$

$$= (nk-1) s^2 [1 + (n-1) P]$$

$$\text{var } [y_{\text{sys}}] = \frac{nk-1}{nk} \frac{s^2}{n} [1 + (n-1) P]$$

Theorem:-

$$\text{var } (\bar{y}_{\text{sys}}) = \frac{k-1}{nk} s_{\text{wst}}^2 [1 + (n-1) P_{\text{wst}}]$$

Proof:-

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

$$= \frac{1}{k} \sum_{i=1}^k \left[\frac{1}{n} \sum_{j=1}^n y_{ij} - \frac{1}{n} \sum_{j=1}^n y_{i.} \right]^2$$

$$= \frac{1}{n^2 k} \left[\sum_{i=1}^k (y_{i.} - \bar{y}_{..})^2 + \sum_{i=1}^k \sum_{j \neq j_1}^n [y_{ij} - \bar{y}_{.j}] [y_{ij} - \bar{y}_{.j_1}] \right]$$

$$= \frac{1}{n^2 k} \left[n(k-1) s_{\text{wst}}^2 + n(n-1) (k-1) s_{\text{wst}}^2 + s_{\text{wst}}^2 \right]$$

$$= \frac{k-1}{nk} s^2 [1 + (n-1) P_{\text{wst}}]$$

Systematic sampling vs Stratified Random sampling:

$$\text{Var}(\bar{y}_{st}) = \sum_{j=1}^n \left[\frac{1}{n_j^2} - \frac{1}{N^2} \right] p_j^2 \sigma_j^2$$

But

$NE = K$ and $n_j^0 = 1$ [$j = 1, 2, \dots, n$] and

$$p_j^0 = \frac{N_j^0}{N} = \frac{K}{nK} = \frac{1}{n}$$

$$\text{Var}(\bar{y}_{st}) = \sum_{j=1}^n (1 - 1/K) \frac{1}{n^2} \sigma_j^2 \rightarrow (1)$$

$$= \frac{K}{n^2 K} \sum_{j=1}^n \sigma_j^2$$

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

$$= \frac{K}{n^2 K} \sum_{j=1}^n \sigma_j^2$$

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

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

$$\text{Var}(\bar{y}_{st}) = \frac{K-1}{nK} S_{wse}$$

Merits and demerits of systematic sampling :-

1) Systematic sampling is operationally more convenient than simple random sampling or stratified sampling time and work is involved is also respectively much less.

2) Systematic sampling may be more efficient than simple random sampling provided the frame is arranged wholly at random. The most common approach to randomness is provided to by alphabetical list such as name in telephone directory although even those may have certain non-random characteristics.

Demerits :-

The main advantage of systematic sampling is that systematic samples are not in general random sample since the requirement in merit two is rarely fulfilled.

If N is not multiple of n , then

1) The actual sample size is different from that required

2) Sample mean is not an unbiased estimate of the population mean.

However, these disadvantages can be overcome by adopting a technique known as circular systematic Sampling [C.S.S]

3) In systematic sampling we have,

$$\text{Var}(\bar{y}_{\text{sys}}) = \frac{1}{kn} \sum_{i=1}^n (y_i - \bar{y}_{\cdot})^2$$

Systematic Sampling vs simple random sampling:-

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

$$\text{Var}(\bar{y}_{\text{sys}}) = \frac{N-1}{N} s^2 - \frac{k(n-1)}{N} s_{\text{wsy}}^2 \rightarrow (2)$$

$$\text{Var}(\bar{y}_{i \cdot})_n - \text{Var}(\bar{y}_{\text{sys}})$$

$$\Rightarrow \left[\frac{N-n}{N} \right] \frac{s^2}{n} - \left[\frac{N-1}{N} s^2 - \frac{k(n-1)}{N} s_{\text{wsy}}^2 \right]$$

$$\rightarrow \left[\frac{N-n}{Nn} \right] s^2 - \frac{N-1}{N} s^2 + \frac{k(n-1)}{N} s_{\text{wsy}}^2$$

$$\Rightarrow \left[\frac{N-n}{n} - (N-1) \right] \frac{s^2}{N} + \frac{k(n-1)}{N} s_{\text{wsy}}^2$$

$$\Rightarrow \left[\frac{N-n}{n} - N+1 \right] \frac{s^2}{N} + \frac{k(n-1)}{N} s_{\text{wsy}}^2$$

$$= \frac{N(1-n)}{n} \frac{s^2}{N} + \frac{K(n-1)}{N} s_{wsy}^2$$

$$= \frac{K(n-1)}{nK} s_{wsy}^2 - \frac{(n-1)}{n} s^2$$

$$= \frac{(n-1)}{n} s_{sys}^2 - \frac{n-1}{n} s^2$$

$$\text{Var}(y_i) - \text{Var}(y_{sys}) = \left[\frac{n-1}{n} \right] [s_{wsy}^2 - s_{sys}^2]$$

$$\text{Var}(y_i) \cdot n - \text{Var}(y_{sys}) > 0$$

$$s_{wsy}^2 > s^2$$

Theorem :-

$$\text{Var}(\bar{y}_{st}) \leq \text{Var}(\bar{y}_{st}) \leq \text{Var}(\bar{y}_N)$$

Proof :-

$$y_i^0 = i \quad (i = 1, 2, \dots, n)$$

$$\sum_{i=1}^N y_i^0 = \sum_{i=1}^n i = N(N+1)$$

$$\sum_{i=1}^N y_i^0{}^2 = \sum_{i=1}^n i^2 = \frac{N(N+1)(2N+1)}{6}$$

$$\bar{y}_N = \frac{1}{N} \sum_{i=1}^N y_i^0$$

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

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

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

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

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

on simplification :-

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

$$\text{Var}(\bar{y}_s)_R = \left[\frac{N-n}{N} \right] \left[\frac{s^2}{n} \right] \Rightarrow \left[\frac{1}{n} - \frac{1}{N} \right] s^2$$

$$\text{Var}(\bar{y}_n)_R = \left[\frac{1}{nk} - \frac{1}{N} \right] \cdot \frac{N(N+1)}{12} \quad [\because N = nk]$$

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

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

$$\text{Var}(\bar{y}_n)_R = \frac{k-1}{k} \frac{(nk+1)}{12} \rightarrow (7)$$

$$\text{var}(y_{st}) = \frac{k-1}{n^2 k} \sum_{j=1}^n S_{ij}^2$$

$$\Rightarrow \sum_{j=1}^n \frac{k(k+1)}{12} = \frac{n k (k+1)}{12} \quad \left[\because S^2 = \frac{N+1(N)}{12} \right]$$

$$\text{var}(y_{st}) = \frac{k-1}{n^2 k} \frac{n k (k+1)}{12}$$

$$\text{var}(y_{st}) = \frac{k-1}{n} \frac{(k+1)}{12} \Rightarrow \frac{k^2 + (k-1)}{12n} = \frac{k^2 - 1}{12n} \rightarrow \textcircled{2}$$

$$\text{var}(\bar{y}_{sys}) = \frac{1}{k} \sum_{i=1}^k (y_{i0} - \bar{y}_{..})^2$$

$$y_{i0} = \frac{1}{n} \sum_{j=1}^n y_{ij}^0$$

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

$$= \frac{1}{n} [n i (1+2+3+\dots+n-1) + k]$$

$$y_{i0} = \frac{i + (n-1)k}{2}$$

$$\bar{y}_{..} = \bar{y}_{N..} = \frac{N+1}{2} = \frac{n k + 1}{2}$$

$$(y_{i0} - \bar{y}_{..}) = \frac{i + (n-1)k}{2} - \frac{n k + 1}{2}$$

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

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

$$\begin{aligned}
 \text{Var}(\bar{y}_{\text{sys}}) &= \frac{1}{k} \sum_{i=1}^k (y_{i0} - \bar{y}_{\cdot\cdot})^2 \\
 &= \frac{1}{k} \sum_{i=1}^k \left[\frac{i - k + 1}{2} \right]^2 \\
 &= \frac{1}{k} \sum_{i=1}^k \left[i^2 + \left[\frac{k+1}{2} \right]^2 - \frac{2i[k+1]}{2} \right] \\
 &= \left[\frac{1}{k} \sum_{i=1}^k i^2 + \frac{1}{k} \sum_{i=1}^k \frac{(k+1)^2}{4} - \frac{1}{k} \sum_{i=1}^k 2i(k+1) \right] \\
 &= \left[\frac{1}{k} \sum_{i=1}^k i^2 + \frac{(k+1)^2}{4} - \frac{k+1}{k} \sum_{i=1}^k i \right] \\
 &= \frac{1}{k} \left[\frac{k(k+1)(2k+1)}{6} + \frac{k+1)^2}{4} - \frac{k+1}{k} \cdot \frac{k(k+1)}{2} \right] \\
 &= \frac{[(k+1) \cdot [2k+1] + \frac{(k+1)^2}{4} - \frac{k+1)^2}{2}]}{6}
 \end{aligned}$$

on simplification,

$$\text{Var}(\bar{y}_{\text{sys}}) = \frac{k^2 - 1}{12} \rightarrow \textcircled{3}$$

From ① ② a ③, we get:

$$\begin{aligned}
 \text{Var}(\bar{y}_{\text{st}}) : \text{Var}(\bar{y}_{\text{sys}}) : \text{Var}(y_n)R \\
 \frac{k^2 - 1}{12n} : \frac{k^2 - 1}{12} : \frac{k - 1(nk + 1)}{12}
 \end{aligned}$$

$$\frac{1}{n} : 1 : n$$

$$\text{var}(y_{st}) \leq \text{var}(\bar{y}_{sys}) \leq \text{var}(\bar{y}_n) \cdot k$$

Ratio And Regression Estimators :-

In most of the sample survey information auxiliary variable is available along with the information on the variable under study. This information can be suitably used in estimating the mean of the character under study mean efficiently. The auxiliary information can be utilized either is the section of the sample or at the estimation of the stage. The common estimation which at the information on the auxiliary variable are ratio and the regression estimates.

Ratio Estimates :-

In the ratio method an auxiliary variable x_i correlated with y_i is obtained

for such unit of the population total x of x_i must be known. In practice x_i is often the value of y_i at some previous time when a complete census was taken. The aim in this method is obtained increased precision by taking advantages of the correlation coefficient between x_i and y_i .

Sampling Theory

UNIT - IV

Cluster Sampling :-

It is one of the basic assumptions in any sampling procedure that the population can be divided into a finite number of distinct and identifiable units called sampling units. The smallest units into which the population can be divided are called elements of the population. The group of such elements are called clusters.

Equal cluster sampling :-

i) Suppose the population is divided into N clusters and each cluster is of size n .

ii) select a sample of n clusters from N clusters by the method of SRS, generally

WOR, so

$$\text{total population size} = NM$$

$$\text{Total sample size} = nM.$$

Estimate of mean and its variance :-

(2)

First select n clusters from N clusters by SRSWOR. Based on n clusters find the mean of each cluster separately based on all the units in every cluster. So, we have the cluster means as $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_n$. Consider the mean of all such cluster means.

$$\bar{y}_{et} = \frac{1}{n} \sum_{i=1}^n \bar{y}_i$$

Bias:

$$\begin{aligned} E(\bar{y}_{et}) &= \frac{1}{n} \sum_{i=1}^n E(\bar{y}_i) \\ &= \frac{1}{n} \sum_{i=1}^n \bar{y} \\ &= \bar{y} \end{aligned}$$

Thus, \bar{y}_{et} is an unbiased estimator of \bar{y}

Variance:

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

$$\text{var}(\bar{y}_{et}) = E(\bar{y}_{et} - \bar{y})^2 = \frac{N-n}{Nn} S_b^2$$

$$S_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{y}_i - \bar{y})^2$$

In case SRSWOR,

$$\widehat{\text{Var}}(\bar{Y}_{ct}) = \frac{N-n}{Nn} S_b^2$$

$$S_b^2 = \frac{1}{n-1} \sum_{t=1}^n (\bar{Y}_t - \bar{Y}_{ct})^2$$

Multistage Sampling :-

Multistage sampling is defined as a sampling method that divides the population into groups for conducting research. It is a complex form of cluster sampling, sometimes, also known as multistage sampling. During this sampling method, significant clusters of the selected people are split into subgroups at various stages to make it simpler for primary data collection.

Estimation of relative Efficiency :-

$$E = \frac{s^2}{MS_b^2}$$

An estimator of E can be obtained by the

Substituting the estimates of s^2 and S_b^2

$$E(S_b^2) = E\left[\frac{1}{n} \sum_{i=1}^n (\bar{y}_i - \bar{y}_e)^2\right]$$

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

$$= S_b^2$$

Y_b^2 is an unbiased estimator of S_b^2

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

$s_i^2, i = 1, 2, \dots, N$

$$E(S_w^2) = E\left[\frac{1}{N} \sum_{i=1}^N s_i^2\right] = \frac{1}{N} \sum_{i=1}^N E(s_i^2)$$

$$= \frac{1}{N} \sum_{i=1}^N \left[\frac{1}{N} \sum_{i=1}^N s_i^2\right]$$

$$= \frac{1}{N} \sum_{i=1}^N s_i^2$$

$$= S_w^2$$

consider,

$$S^2 = \frac{1}{MN-1} \sum_{i=1}^N \sum_{j=1}^M (Y_{ij} - \bar{Y})^2$$

$$\text{OR } (MN-1)S^2 = \sum_{i=1}^N \sum_{j=1}^M \left[(Y_{ij} - \bar{Y}_i) + (\bar{Y}_i - \bar{Y}) \right]^2$$

$$= \sum_{i=1}^N \sum_{j=1}^M \left[(Y_{ij} - \bar{Y}_i)^2 + (\bar{Y}_i - \bar{Y})^2 \right]$$

$$= \sum_{i=1}^N (M-1) S_i^2 + M(N-1) S_b^2$$

$$= N(M-1) \bar{S}_w^2 + M(N-1) S_b^2$$

$$S^2 = \frac{1}{MN-1} \left[N(M-1) \bar{S}_w^2 + M(N-1) S_b^2 \right]$$

So,

$$\text{var}(Y_{ij}) = \frac{N-n}{Nn} S_b^2$$

$$\text{var}(\bar{Y}_{nm}) = \frac{N-n}{Nn} \frac{\bar{S}_w^2}{M}$$

where $S_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{Y}_i - \bar{Y}_{at})^2$

$$E = \frac{S^2}{MS_b^2} \text{ or}$$

$$E = \frac{N(M-1) \bar{S}_w^2 + M(N-1) S_b^2}{M(NM-1) S_b^2}$$

$$M(N-1) = MN - 1 ; MN - 1 = MN$$

$$E = \frac{1}{M} + \left(\frac{M-1}{M} \right) \frac{\bar{S}_w^2}{MS_b^2}$$

$$\hat{E} = \frac{1}{M} + \left(\frac{M-1}{M} \right) \frac{\bar{S}_w^2}{MS_b^2}$$

Two Stage Sampling With Equal First Stage

Units :-

- 1) The population consist of NM elements.
- 2) NM elements are grouped into N first stage units of M second stage units each.
- 3) sample of n_1 first stage units is selected.
- 4) units at each stage are selected with SRSWOR.

Two stage sampling,

$$E(\hat{\theta}) = E_1[E_2(\hat{\theta})].$$

$$E(\hat{\theta}) = E_1 [E_2 \{E_3(\hat{\theta})\}] \dots$$

$$\begin{aligned} \text{var}(\hat{\theta}) &= E[(\hat{\theta} - \theta)^2] \\ &= E_1 E_2 (\hat{\theta} - \theta)^2 \end{aligned}$$

Consider,

$$\begin{aligned} E_2 (\hat{\theta} - \theta)^2 &= E_2 (\hat{\theta}^2) - 2\theta E_2(\hat{\theta}) + \theta^2 \\ &= [E_2(\hat{\theta}^2) + V_2(\hat{\theta})] - 2\theta E_2(\hat{\theta}) + \theta^2 \end{aligned}$$

Average over first stage selection,

$$\begin{aligned} E_1 E_2 (\hat{\theta} - \theta)^2 &= E_1 [E_2(\hat{\theta}^2)] + E_1 [V_2(\hat{\theta}) - 2\theta E_1 E_2(\hat{\theta}) + E_1(\theta^2)] \\ &= E_1 [E_2(\hat{\theta}^2) - \theta^2] + E_1 [V_2(\hat{\theta})] \end{aligned}$$

$$\text{var}(\hat{\theta}) = V_1 [E_2(\hat{\theta}^2)] + E_1 [V_2(\hat{\theta})]$$

Three stage sampling:-

$$\text{var}(\hat{\theta}) = V_1 [E_2 \{E_3(\hat{\theta})\}] + E_1 [V_2 \{E_3(\hat{\theta})\}] + E_1 [E_2 \{V_3(\hat{\theta})\}]$$

Estimation of population mean:-

$$\bar{Y} = \bar{Y}_{mn}$$

Bias :-

$$E(\bar{Y}) = E_1 [E_2(\bar{Y}_{mn})]$$

$$= E_1 [E_2(\bar{Y}_{im}/i)]$$

$$= E_1 [E_2(\bar{y}_{im}/i)]$$

$$= E_1 \left[\frac{1}{n} \sum_{i=1}^n \bar{y}_e \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \bar{y}_i$$

$$= \bar{Y}$$

Variance :-

$$\text{var}(\bar{Y}) = E_1 [V_2(\bar{Y}/i)] + V_1 [E_2(\bar{Y}/i)]$$

$$= E_1 \left[V_2 \left\{ \frac{1}{n} \sum_{i=1}^n \bar{y}_e / i \right\} \right] + V_1 \left[E_2 \left\{ \frac{1}{n} \sum_{i=1}^n \bar{y}_e / i \right\} \right]$$

$$= E_1 \left[\frac{1}{n^2} \sum_{i=1}^n V(\bar{Y}_i / i) \right] + V_1 \left[\frac{1}{n} \sum_{i=1}^n E_2(\bar{Y}_i / i) \right]$$

$$\begin{aligned}
 &= E \left[\frac{1}{n^2} \sum_{i=1}^n \left(\frac{1}{m} - \frac{1}{M} \right) S_i^2 \right] + V \left[\frac{1}{n} \sum_{i=1}^n \bar{y}_i \right] \\
 &= \frac{1}{n^2} \sum_{i=1}^n \left[\frac{1}{m} - \frac{1}{M} \right] E_i(S_i^2) + V_1(\bar{y}_c) \\
 &= \frac{1}{n^2} n \left(\frac{1}{m} - \frac{1}{M} \right) \bar{S}_w^2 + \frac{N-n}{Nn} S_b^2 \\
 &= \frac{1}{n} \left[\frac{1}{m} - \frac{1}{M} \right] \bar{S}_w^2 + \left(\frac{1}{n} - \frac{1}{N} \right) S_b^2
 \end{aligned}$$

where $\bar{S}_w^2 = \frac{1}{N} \sum_{i=1}^N S_i^2 = \frac{1}{N(M-1)} \sum_{i=1}^N \sum_{j=1}^M (y_{ij} - \bar{y}_i)^2$

$$S_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{y}_i - \bar{Y})^2$$

Two stage sampling with unequal first stage units :-

Let y_{ij} : value of j^{th} second stage unit of i^{th} first stage unit

M_i : number of second stage units in i^{th} first stage unit

$M_0 = \sum_{i=1}^N M_i$: total no. of second stage units in the population

m_i : num. of second stage units to be selected from i^{th} first stage unit if it is in the sample

$m_0 = \sum_{i=1}^n m_i$: total no. of second stage units in sample

$$\bar{y}_i(m_i) = \frac{1}{m_i} \sum_{j=1}^{m_i} y_{ij}$$

$$\bar{y}_i = \frac{1}{M_i} \sum_{j=1}^{M_i} y_{ij}$$

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N \bar{y}_i = \bar{y}_N$$

$$\bar{y} = \frac{\sum_{i=1}^N \sum_{j=1}^{M_i} y_{ij}}{\sum_{i=1}^N M_i} = \frac{\sum_{i=1}^N M_i \bar{y}_i}{\bar{M} N} = \frac{1}{N} \sum_{i=1}^N u_i \bar{y}_i$$

$$u_i = \frac{M_i}{\bar{M}}$$

$$\bar{M} = \frac{1}{N} \sum_{i=1}^N M_i$$

Estimation of mean and variance :-

$$\hat{Y} = \bar{y}_{s_2} = \frac{1}{n} \sum_{i=1}^n \bar{y}_i(m_i)$$

Bias :-

$$E(\bar{y}_{s_2}) = E\left[\frac{1}{n} \sum_{i=1}^n \bar{y}_i(m_i)\right]$$

$$= E_1\left[\frac{1}{n} \sum_{i=1}^n E_2 \bar{y}_i(m_i)\right]$$

$$\bar{y}_{s2} + \frac{N-1}{Nn} \frac{1}{N-1} \sum_{i=1}^N (M_i - \bar{m}) (y_{i(m_i)} - \bar{y}_{s2})$$

Variance :-

$$\text{var}(\bar{y}_{s2}) = \text{var} [E(\bar{y}_{s2}/n)] + E [\text{var}(\bar{y}_{s2}/n)]$$

$$= \text{var} \left[\frac{1}{n} \sum_{i=1}^n \bar{y}_i \right] + E \left[\frac{1}{n^2} \sum_{i=1}^n \text{var}(y_{i(m_i)}/i) \right]$$

$$= \left(\frac{1}{n} - \frac{1}{N} \right) S_b^2 + E \left[\frac{1}{n^2} \sum_{i=1}^n \left(\frac{1}{m_i} - \frac{1}{M_i} \right) S_i^2 \right]$$

$$= \left(\frac{1}{n} - \frac{1}{N} \right) S_b^2 + \frac{1}{Nn} \sum_{i=1}^N \left(\frac{1}{m_i} - \frac{1}{M_i} \right) S_i^2$$

$$S_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{y}_i - \bar{y}_N)^2$$

$$S_i^2 = \frac{1}{M_i - 1} \sum_{j=1}^{M_i} (y_{ij} - \bar{y}_i)^2$$

$$MSE(\bar{y}_{s2}) = \text{var}(\bar{y}_{s2}) + [\text{Bias}(\bar{y}_{s2})]^2$$

$$= E_i \left[\frac{1}{n} \sum_{i=1}^n \bar{y}_i \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \bar{y}_i = \bar{y}_N \neq \bar{y}$$

$$\text{Bias}(\bar{y}_{SE}) = E(\bar{y}_{SE}) - \bar{y}$$

$$= \frac{1}{N} \sum_{i=1}^N \bar{y}_i - \frac{1}{NM} \sum_{i=1}^N M_i \bar{y}_i$$

$$= \frac{1}{NM} \left[\sum_{i=1}^N M_i \bar{y}_i - \frac{1}{N} \left(\sum_{i=1}^N \bar{y}_i \right) \left(\sum_{i=1}^N M_i \right) \right]$$

$$= \frac{1}{NM} \sum_{i=1}^N (M_i - \bar{M}) (\bar{y}_i - \bar{y}_N)$$

$$\text{Bias}(\bar{y}_{SE}) = \frac{-N-1}{NM(N-1)} \sum_{i=1}^N (M_i - \bar{M}) (\bar{y}_{i(m)} - \bar{y}_{SE})$$

$$E[\text{Bias}(\bar{y}_{SE})] = \frac{-N-1}{NM} E_i \frac{1}{n-1} \sum_{i=1}^N E_2 (M_i - \bar{M}) (\bar{y}_{i(m)} - \bar{y}_{SE})$$

$$= \frac{-N-1}{NM} E \left[\frac{1}{n-1} \sum_{i=1}^N (M_i - \bar{M}) (\bar{y}_i - \bar{y}_N) \right]$$

$$= \frac{-1}{NM} \sum_{i=1}^N (M_i - \bar{M}) (\bar{y}_i - \bar{y}_N) = \bar{y}_N - \bar{y}$$

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

Sampling Theory

①

UNIT - V

Double Sampling :-

The ratio and regression methods of estimation require the knowledge of population mean of an auxiliary variable (\bar{X}) to estimate the population mean of study variable (\bar{Y}). If the information on the auxiliary variable is not available, then there are two options - one option is to collect a sample only on study variable and use sample mean as an estimator of the population mean.

Optimal allocation :-

optimal allocation is a procedure for dividing the sample survey. A stratified sample selects separate samples from subgroups called strata of the population and can often increase the accuracy of survey results.

Double sampling in ratio method of estimation :-

If the population mean \bar{Y} is not known, then the double sampling technique is applied.

$$\bar{Y} = \bar{X}' = \frac{1}{n'} \sum_{i=1}^{n'} x_i$$

$$\frac{\bar{Y}}{\bar{X}} = \frac{\bar{Y}'}{\bar{X}'}$$

Let

$$\varepsilon_0 = \frac{\bar{y} - \bar{Y}}{\bar{Y}}, \quad \varepsilon_1 = \frac{\bar{x} - \bar{X}}{\bar{X}}, \quad \varepsilon_2 = \frac{\bar{X}' - \bar{X}}{\bar{X}}$$

$$E(\varepsilon_0) = E(\varepsilon_1) = E(\varepsilon_2) = 0$$

$$E(\varepsilon_1^2) = \left(\frac{1}{n} - \frac{1}{N} \right) C_x^2$$

$$E(\varepsilon_1 \varepsilon_2) = \frac{1}{\bar{X}^2} E(\bar{x} - \bar{X})(\bar{X}' - \bar{X})$$

$$= \frac{1}{\bar{X}^2} E_1 \left[E_2(\bar{x} - \bar{X})(\bar{X}' - \bar{X}) \mid n \right]$$

$$= \frac{1}{\bar{X}^2} E_1 \left[(\bar{X}' - \bar{X})^2 \right]$$

$$= \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{S_x^2}{\bar{X}^2}$$

$$= \left(\frac{1}{n'} - \frac{1}{N} \right) C_x^2 = E(\varepsilon_2^2)$$

$$E(\varepsilon_0 \varepsilon_2) = \frac{1}{\bar{X} \bar{Y}} \text{cov}(\bar{y}, \bar{x})$$

$$= \frac{1}{\bar{x}\bar{y}} \text{cov}[E(\bar{y}/n'), E(\bar{x}/n')] + \frac{1}{\bar{x}\bar{y}} E[\text{cov}(\bar{y}, \bar{y})/n']$$

$$= \frac{1}{\bar{x}\bar{y}} \text{cov}[\bar{y}, \bar{x}] + \frac{1}{\bar{x}\bar{y}} E[\text{cov}(\bar{y}, \bar{y})]$$

$$= \frac{1}{\bar{x}\bar{y}} \text{cov}[\bar{y}, \bar{x}]$$

$$= \left(\frac{1}{n'} - \frac{1}{N}\right) \frac{S_{xy}}{\bar{x}\bar{y}}$$

$$= \left(\frac{1}{n'} - \frac{1}{N}\right) \rho \frac{S_x}{\bar{x}} \frac{S_y}{\bar{y}}$$

$$= \left(\frac{1}{n'} - \frac{1}{N}\right) \rho C_x C_y$$

where \bar{y} is the sample mean of y 's based on the sample size n' .

$$E(\varepsilon_0 \varepsilon_1) = \frac{1}{\bar{x}\bar{y}} \text{cov}(\bar{y}, \bar{x})$$

$$= \left(\frac{1}{n} - \frac{1}{N}\right) \frac{S_{xy}}{\bar{x}\bar{y}}$$

$$= \left(\frac{1}{n} - \frac{1}{N}\right) \rho \frac{S_x}{\bar{x}} \frac{S_y}{\bar{y}}$$

$$= \left(\frac{1}{n} - \frac{1}{N}\right) \rho C_x C_y$$

$$E(\varepsilon_0^2) = \frac{1}{\bar{y}^2} \text{var}(\bar{y})$$

$$= \frac{1}{\bar{y}^2} [V_1 \{E_2(\bar{y}/n)\}^2 + E_1 \{V_2(\bar{y}_n/n)\}^2]$$

$$= \frac{1}{\bar{y}^2} [V_1(\bar{y}/n) + E_1 \left\{ \left(\frac{1}{n} - \frac{1}{n_1} \right) S_y^2 \right\}]$$

$$= \frac{1}{\bar{y}^2} \left[\left(\frac{1}{n_1} - \frac{1}{N} \right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n} \right) S_y^2 \right]$$

$$= \left(\frac{1}{n} - \frac{1}{N} \right) \frac{S_y^2}{\bar{y}^2}$$

$$= \left(\frac{1}{n} - \frac{1}{N} \right) C_y^2$$

$$E(\varepsilon_1 \varepsilon_2) = \frac{1}{\bar{x}^2} \text{cov}(\bar{x}, \bar{x}')_1$$

$$= \frac{1}{\bar{x}^2} [\text{cov} \{ E(\bar{x}/n), E(\bar{x}'/n_1) \} + 0]$$

$$= \frac{1}{\bar{x}^2} \text{var}(\bar{x}')$$

where $\text{var}(\bar{x}')$ is the variance of mean of SC based on an initial sample of size n_1 .

Double sampling in regression method of

Estimation :-

When the population's mean of the auxiliary variable \bar{X} , is not known, then double sampling is used as follows :-

i) A large sample of size n' is taken from the population of SRSWOR from which the population mean \bar{X} is estimated as \bar{x}' , (i.e) $\bar{X} = \bar{x}'$

ii) Then a subsample of size n is chosen from the larger sample and both the variables X and Y are measured from it by taking \bar{x}' in place of \bar{X} and treat it as known.

Then $E(\bar{x}') = \bar{X}$, $E(\bar{x}) = \bar{X}$, $E(\bar{y}) = \bar{Y}$.

$$\bar{Y}_{regl} = \bar{Y} + \beta (\bar{x}' - \bar{x})$$

where,
$$\beta = \frac{S_{xy}}{S_x} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$\beta = \frac{S_{xy}}{S_x^2}$ based on the sample of size n .

Let,

$$\epsilon_0 = \frac{\bar{y} - \bar{Y}}{\bar{Y}} \Rightarrow \bar{y} = (1 + \epsilon_0) \bar{Y}$$

$$\epsilon_1 = \frac{\bar{x} - \bar{X}}{\bar{X}} \Rightarrow \bar{x} = (1 + \epsilon_1) \bar{X}$$

$$\epsilon_2 = \frac{\bar{x}' - \bar{x}}{\bar{x}} \Rightarrow \bar{x}' = (1 + \epsilon_2) \bar{x}$$

$$\epsilon_3 = \frac{S_{oxy} - S_{oxy}}{S_{oxy}} \Rightarrow S_{oxy} = (1 + \epsilon_3) S_{oxy}$$

$$\epsilon_4 = \frac{S_x^2 - S_x^2}{S_x^2} \Rightarrow S_x^2 = (1 + \epsilon_4) S_x^2$$

$$E(\epsilon_1) = 0, E(\epsilon_2) = 0, E(\epsilon_3) = 0, E(\epsilon_4) = 0$$

Define

$$\mu_{21} = E[(\bar{x} - \bar{X})^2 (y - \bar{Y})]$$

$$\mu_{30} = E[(\bar{x} - \bar{X})^3]$$

Estimation error of \bar{Y}_{rd}

$$\frac{1}{Y_{rd}} = \frac{(1+\varepsilon_0)}{(1+\varepsilon_1)} (1+\varepsilon_2) \frac{\bar{Y}}{\bar{X}}$$

$$= \bar{Y} (1+\varepsilon_0) (1+\varepsilon_2) (1+\varepsilon_1)^{-1}$$

$$= \bar{Y} (1+\varepsilon_0) (1+\varepsilon_2) (1-\varepsilon_1+\varepsilon_1^2-\dots)$$

$$= \bar{Y} (1+\varepsilon_0 + \varepsilon_2 + \varepsilon_0 \varepsilon_2 - \varepsilon_1 - \varepsilon_0 \varepsilon_1 - \varepsilon_1 \varepsilon_2 + \varepsilon_1^2)$$

Bias of \bar{Y}_{rd}

$$E(\bar{Y}_{rd}) = \bar{Y} [1 + 0 + 0 + E(\varepsilon_0 \varepsilon_2) - 0 - E(\varepsilon_0 \varepsilon_1) - E(\varepsilon_1 \varepsilon_2) + E(\varepsilon_1^2)]$$

$$\text{Bias}^0(\bar{Y}_{rd}) = E(\bar{Y}_{rd}) - \bar{Y}$$

$$= \bar{Y} [E(\varepsilon_0 \varepsilon_2) - E(\varepsilon_0 \varepsilon_1) - E(\varepsilon_1 \varepsilon_2) + E(\varepsilon_1^2)]$$

$$= \bar{Y} \left[\left(\frac{1}{n} - \frac{1}{N} \right) \rho_{xy} c_y - \left(\frac{1}{n} - \frac{1}{N} \right) \rho_{xy} c_y - \left(\frac{1}{n} - \frac{1}{N} \right) c_y^2 + \left(\frac{1}{n} - \frac{1}{N} \right) c_x^2 \right]$$

$$= \bar{Y} \left(\frac{1}{n} - \frac{1}{n_1} \right) (C_x^2 - \rho C_x C_y)$$

$$= \bar{Y} \left(\frac{1}{n} - \frac{1}{n_1} \right) C_x (C_x - \rho C_y)$$

MSE of \hat{Y}_{Rd} :

$$MSE(\hat{Y}_{Rd}) = E(\hat{Y}_{Rd} - \bar{Y})^2$$

$$= \bar{Y}^2 E(\epsilon_0 + \epsilon_2 - \epsilon_1)^2$$

$$= \bar{Y}^2 E[\epsilon_0^2 + \epsilon_1^2 + \epsilon_2^2 + 2\epsilon_0\epsilon_2 - 2\epsilon_0\epsilon_1 - 2\epsilon_1\epsilon_2]$$

$$= \bar{Y}^2 E[\epsilon_0^2 + \epsilon_1^2 + \epsilon_2^2 + 2\epsilon_0\epsilon_2 - 2\epsilon_0\epsilon_1 - 2\epsilon_1\epsilon_2]$$

$$= \bar{Y}^2 \left[\left(\frac{1}{n} - \frac{1}{N} \right) C_y^2 + \left(\frac{1}{n} - \frac{1}{N} \right) C_x^2 - \left(\frac{1}{n_1} - \frac{1}{N} \right) C_x^2 + 2 \left(\frac{1}{n} - \frac{1}{N} \right) \rho C_x C_y \right]$$

$$- 2 \left(\frac{1}{n} - \frac{1}{N} \right) \rho C_x C_y$$

$$= \bar{Y}^2 \left(\frac{1}{n} - \frac{1}{N} \right) (C_x^2 + C_y^2 - 2\rho C_x C_y) + \bar{Y}^2 \left(\frac{1}{n_1} - \frac{1}{n} \right) C_x (C_x - \rho C_y)$$

$$(2\rho C_x C_y - C_x^2)$$

$$= MSE(\text{ratio estimator}) + \bar{Y}^2 \left(\frac{1}{n_1} - \frac{1}{n} \right) (2\rho C_x C_y - C_x^2)$$

(9)

$$2P. C_x C_y - C_x^2 > 0$$

or

$$P > \frac{1}{2} \frac{C_x}{C_y}$$

Estimation Error :-

Then,

$$\hat{Y}_{regd} = \bar{y} + \hat{\beta} (\bar{x}' - \bar{x})$$

$$= \bar{y} + \frac{S_{oxy} (1 + \varepsilon_3)}{S_x^2 (1 + \varepsilon_4)} (\varepsilon_2 - \varepsilon_1) \bar{x}$$

$$= \bar{y} + \bar{x} \frac{S_{oxy}}{S_x^2} (1 + \varepsilon_3) (\varepsilon_2 - \varepsilon_1) (1 + \varepsilon_4)^{-1}$$

$$= \bar{y} + \bar{x} \beta (1 + \varepsilon_3) (\varepsilon_2 - \varepsilon_1) (1 - \varepsilon_4 + \varepsilon_4^2 - \dots)$$

$$\hat{Y}_{regd} = \bar{y} + \bar{x} \beta (\varepsilon_2 + \varepsilon_2 \varepsilon_3 - \varepsilon_2 \varepsilon_4 - \varepsilon_1 - \varepsilon_1 \varepsilon_3 + \varepsilon_1 \varepsilon_4)$$

Bias :-

$$E(\hat{Y}_{regd}) = \bar{y} + \bar{x} \beta [E(\varepsilon_2 \varepsilon_3) - E(\varepsilon_2 \varepsilon_4) - E(\varepsilon_1 \varepsilon_3) + E(\varepsilon_1 \varepsilon_4)]$$

$$\text{Bias}(\hat{Y}_{regd}) = E(\hat{Y}_{regd}) - \bar{y}$$

$$= \bar{y} \beta \left[\left(\frac{1}{n'} - \frac{1}{N} \right) \frac{1}{N} \sum \frac{(\bar{x}' - \bar{x})(S_{xy} - S_y)}{\bar{x} S_{xy}} \right]$$

$$= \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{1}{N} \left(\sum \frac{(\bar{x}' - \bar{x})(S_x^2 - S_y^2)}{\bar{x} S_x^2} \right)$$

$$= \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{1}{N} \sum \left[\frac{(\bar{x} - \bar{x})(S_{xy} - S_{xy})}{\bar{x} S_{xy}} \right]$$

$$+ \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{1}{N} \left[\sum \frac{(\bar{x} - \bar{x})(S_x^2 - S_x^2)}{\bar{x} S_x^2} \right]$$

$$= \bar{y} \beta \left[\left(\frac{1}{n'} - \frac{1}{N} \right) \frac{M_{21}}{\bar{x} S_{xy}} - \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{M_{30}}{\bar{x} S_x^2} \right]$$

$$= \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{M_{21}}{\bar{x} S_{xy}} + \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{M_{30}}{\bar{x} S_x^2}$$

$$= -\beta \left(\frac{1}{n'} - \frac{1}{n} \right) \left(\frac{M_{21}}{S_{xy}} - \frac{M_{30}}{S_x^2} \right)$$

Double Sampling for Stratification :-

The method of post-stratification is useful only if the relative proportion of each stratum in the population $w_h = \frac{N_h}{N}$ is known for each stratum h . If these proportions are not known, double sampling may be used, with an initial (large) sample used to classify the units into strata and then a stratified sample selected from the initial sample.

Step 1: n' initial simple random sample are selected from a population of N units. These units are classified into strata. With n'_h observed to be in stratum h . The population proportion $w_h = \frac{N_h}{N}$ is estimated by the sample proportion:

$$w_h = \frac{n'_h}{n'} \quad h = 1, \dots, L.$$

Step 2: A second sample is then selected by stratified random sampling from the first sample. These units are classified into strata with n_h units selected from the n'_h sample units in stratum h . Measurement of Y_{hi} is recorded for each unit in the second sample.

We denote the sample mean in stratum h in the second sample as :

$$\bar{Y}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} Y_{hi}$$