

## CLUSTER SAMPLING

Sampling Scheme :

Its advantages and disadvantages

Cluster Sampling consists in forming suitable clusters of units and surveying all the units in a sample of clusters selected according to an appropriate sampling scheme. The advantages of cluster sampling from the point of view of cost arise mainly due to the fact that collection of data for nearby units is easier, faster, cheaper and more convenient than observing units scattered over a region. For instance, in ~~application~~ a population survey it may be cheaper to collect data from all the persons in a sample of households than from a sample of the same no. of persons selected directly from all the persons. Similarly, it would be operationally more convenient to survey all households situated in a sample of villages than to survey a sample of the same no. of households selected at random from a list of all households.

Another example of utility of cluster sampling is provided by crop survey where locating a randomly selected plot requires a considerable part of total time taken for the survey. But once the plot is located, the time taken for identifying and surveying a few neighbouring plots will be only marginal.

Because of its operational convenience and the possible reduction of cost, cluster sampling is resorted to in many surveys using mutually exclusive non-overlapping clusters formed by grouping nearby units which can be conveniently observed together. For a given total number of sampling units, cluster sampling is less efficient than sampling of individual units from the view point of sampling variance as the latter expected to provide a better cross sections of the population than the former due to the usual tendency of units in a cluster, to be similar. In fact the sampling efficiency of cluster sampling is likely to decrease with increasing cluster size. (However, cluster sampling is operationally more convenient and less costly than sampling of units directly due to <sup>the</sup> possible saving in time for journey, identification, contact etc. and hence in many principal situations the loss in sampling efficiency is likely to be offset by the reduction

- Why cluster sampling is used instead of its disadvantages?

Total: NM units

N clusters (mutually exclusive)

M units within each cluster

Case-I: One cluster selected (SRS)

$Y_{ij}$ : jth unit in the ith cluster,  $i=1(1)N, j=1(1)M$

$\bar{Y}_i$ :  $\frac{1}{M} \sum_{j=1}^M Y_{ij}$  = ith cluster mean

$\bar{Y}$ :  $\frac{1}{MN} \sum_{i=1}^N \sum_{j=1}^M Y_{ij} = \frac{1}{N} \sum_{i=1}^N \bar{Y}_i$  = Population Mean

$\hat{Y}_c$  = Mean of Selected Cluster

Claim:  $\hat{Y}_c$  is unbiased for  $\bar{Y}$ .

Proof:  $\hat{Y}_c$  can take values  $\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_N$  each with probability  $\frac{1}{N}$ .

$$\therefore E(\hat{Y}_c) = \frac{1}{N} \sum_{i=1}^N \bar{Y}_i = \bar{Y}$$

Variance:  $\text{Var}(\hat{Y}_c) = \frac{1}{N} \sum_{i=1}^N (\bar{Y}_i - \bar{Y})^2 = \sigma_b^2$

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

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

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

Total Variance of all observations:  $\sigma^2 = \frac{1}{NM} \sum_{i=1}^N \sum_{j=1}^M (Y_{ij} - \bar{Y})^2$

Between Cluster Variance:  $\sigma_b^2 = \frac{1}{N} \sum_{i=1}^N (\bar{Y}_i - \bar{Y})^2$

Intra-cluster correlation coefficient:  $\rho_c = \frac{\sum_{i=1}^N \sum_{j \neq j'=1}^M (Y_{ij} - \bar{Y})(Y_{ij'} - \bar{Y})}{(M-1) \sum_{i=1}^N \sum_{j=1}^M (Y_{ij} - \bar{Y})^2}$

$$\therefore \text{Var}(\hat{\bar{y}}_c) = \frac{1}{NM^2} [NM\sigma^2 + (M-1)\rho_c NM\sigma^2]$$

$$= \frac{1}{M} \sigma^2 [1 + (M-1)\rho_c]$$

If ~~just~~ instead of cluster sampling, a SRSWR of size  $M$  is drawn from the population of size  $NM$ , then the unbiased estimator of  $\bar{y}$  will be the sample mean  $\bar{y}_r$  and its variance is

$$\text{Var}(\bar{y}_r) = \frac{\sigma^2}{M}$$

$\therefore$  Cluster sampling will be more efficient than SRSWR, if

$$\text{Var}(\hat{\bar{y}}_c) < \text{Var}(\bar{y}_r)$$

$$\Rightarrow \frac{\sigma^2}{M} [1 + (M-1)\rho_c] < \frac{\sigma^2}{M}$$

$$\Rightarrow 1 + (M-1)\rho_c < 1 \quad [\because \sigma^2 > 0, M \in \mathbb{N}]$$

$$\Rightarrow (M-1)\rho_c < 0$$

$$\Rightarrow \rho_c < 0 \quad [M > 1]$$

Therefore, cluster sampling will be more efficient than SRSWR only if  $\rho_c$  is negative. But in practice  $\rho_c$  is usually positive when nearby units are grouped to form clusters and hence cluster sampling is usually less efficient than SRSWR.

Case-II : Sampling of n-clusters :

$\bar{y}_i$  : Mean of the  $i$ th sample cluster,  $i=1(1)n$

Claim: U.E. of  $\bar{Y}$  is  $\frac{1}{n} \sum_{i=1}^n \bar{y}_i = \hat{\bar{Y}}_c$

where,  $\bar{y}_i = \frac{1}{M} \sum_{j=1}^M y_{ij}$

Proof:  $E(\bar{y}_i) = \frac{1}{N} \sum_{i=1}^N \bar{Y}_i = \bar{Y}$ , since  $\bar{y}_i$  can take values  $\bar{Y}_1, \dots, \bar{Y}_N$  with equal probability  $\frac{1}{N}$ .

$\therefore E(\hat{\bar{Y}}_c) = \frac{1}{n} \sum_{i=1}^n E(\bar{y}_i) = \bar{Y}$

Variance

$$\text{Var}\left(\frac{1}{n} \sum_{i=1}^n \bar{y}_i\right) = \frac{1}{n^2} \left[ \sum_{i=1}^n \text{Var}(\bar{y}_i) + 2 \sum_{i \neq i'} \text{Cov}(\bar{y}_i, \bar{y}_{i'}) \right]$$

$$\begin{aligned} \text{Var}(\hat{\bar{Y}}_c) &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(\bar{y}_i) \quad [\because \text{Cov}(\bar{y}_i, \bar{y}_{i'}) = 0 \quad \forall i \neq i' = 1(1)n \\ & \quad \text{the clusters are mutually} \\ & \quad \text{exclusive \& independent}] \\ &= \frac{S_b^2}{n} \end{aligned}$$

Carefully note that in case of sampling of n clusters, we are actually doing a simple random sampling from a population whose elements are cluster means as  $\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_N$ .

Let, the sample cluster means are  $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_n$ .

Claim:  $\hat{\bar{Y}}_c = \frac{1}{n} \sum_{i=1}^n \bar{y}_i$  is an unbiased estimator of  $\bar{Y}$ , the population mean.

This is proved previously.

Now,  $\text{Var}(\hat{\bar{Y}}_c) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(\bar{y}_i) + \frac{1}{n^2} \sum_{i \neq i'} \text{Cov}(\bar{y}_i, \bar{y}_{i'})$

$\text{Cov}(\bar{y}_i, \bar{y}_{i'}) = 0 \quad \forall i \neq i' = 1(1)n$ , since each clusters are mutually independent.

We can directly do this as

$$\begin{aligned} \text{Var}(\hat{\bar{Y}}_c) &= \frac{N-n}{nN} S_b^2, \text{ where } S_b^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{Y}_i - \bar{Y})^2 \\ &= \text{Between Cluster SS for population} \\ \text{and } \text{Var}(\hat{\bar{Y}}_c) &= \frac{N-n}{nN} s_b^2, \text{ where } s_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{y}_i - \bar{y})^2 \\ &= \text{Between Cluster SS for sample} \end{aligned}$$

$$\bar{y}_i = \frac{1}{M} \sum_{j=1}^M y_{ij} \quad \text{and} \quad \sigma_b^2 = \frac{1}{N} \sum_{i=1}^N (\bar{y}_i - \bar{y})^2$$

The u.e. of  $\sigma_b^2$  is:  $s_b^2 = \frac{1}{(n-1)} \sum_{i=1}^n (\bar{y}_i - \hat{\bar{y}}_c)^2$

= Variance of sample cluster means

$$\therefore \widehat{\text{Var}}(\hat{\bar{y}}_c) = \frac{s_b^2}{n} \quad \text{if clusters are selected by SRSWR}$$

If clusters are selected by SRSWOR:

$$\widehat{\text{Var}}(\hat{\bar{y}}_c) = \frac{N-n}{n(N-1)} \hat{\sigma}_b^2 = \frac{s_b^2}{n} \cdot \frac{N-n}{N-1}$$

If a SRSWOR of size  $nM$  is selected out of  $NM$  units instead of cluster sampling, the variance of the u.e. of  $\bar{y}_r$  is

$$\text{Var}(\bar{y}_r) = \frac{\sigma^2}{nM} \cdot \frac{NM-nM}{NM-1} = \frac{\sigma^2}{n} \cdot \frac{N-n}{NM-1}$$

In this case, the efficiency of cluster sampling w.r. to SRSWOR

is  $E = \frac{\text{Var}(\bar{y}_r)}{\widehat{\text{Var}}(\hat{\bar{y}}_c)} = \frac{\frac{\sigma^2}{n} \cdot \frac{N-n}{NM-1}}{\frac{s_b^2}{n} \cdot \frac{N-n}{N-1}}$

$$= \frac{\sigma^2}{s_b^2} \cdot \frac{N-1}{NM-1} = \frac{\sigma^2}{\frac{\sigma^2}{M} [1+(M-1)\rho_c]} \cdot \frac{N-1}{NM-1}$$

$$= \frac{M}{[1+(M-1)\rho_c]} \cdot \frac{(N-1)}{(NM-1)} \quad \left[ \because \sigma_b^2 = \frac{\sigma^2}{M} (1+(M-1)\rho_c) \right]$$

Remark:  $E > 1$ , if cluster sampling is more efficient than SRSWOR.

$$\Rightarrow \left( \frac{M}{1+(M-1)\rho_c} \right) \cdot \left( \frac{N-1}{NM-1} \right) > 1$$

$$\Rightarrow NM - M > NM - 1 + \rho_c (M-1) (NM-1)$$

$$\Rightarrow -(M-1) > \rho_c (M-1) (NM-1)$$

$$\Rightarrow \rho_c < -\frac{1}{NM-1} \quad [M > 1]$$

But the intraclass cluster correlation coefficient is positive when nearby units are grouped into form clusters and hence cluster sampling is

usually less efficient than SRSWOR.

Case-III: n-clusters selected (each cluster consists of  $M_i$  units):

$Y_{ij}$ :  $j^{\text{th}}$  unit in the  $i^{\text{th}}$  cluster ( $j=1(1)M_i$ ) in population  
 $i=1(1)N$

$$\bar{Y}_i = \frac{1}{M_i} \sum_{j=1}^{M_i} Y_{ij} = \text{ith cluster mean } (= \frac{Y_i}{M_i}) \text{ in population}$$

$$Y_i = \text{ith cluster total in } N \text{ population} = \sum_{j=1}^{M_i} Y_{ij}$$

$$\bar{Y} = \text{Mean of the total } \sum_{i=1}^N M_i \text{ Population} = \frac{\sum_{i=1}^N \sum_{j=1}^{M_i} Y_{ij}}{\sum_{i=1}^N M_i} = \frac{\sum_{i=1}^N M_i \bar{Y}_i}{\sum_{i=1}^N M_i}$$

We take  $n$ -clusters from  $N$ -clusters and survey all the units in the selected clusters.

Consider,  $y_{ij}$  =  $j^{\text{th}}$  unit in the  $i^{\text{th}}$  cluster in sample ( $j=1(1)M_i$ )  
 $i=1(1)n$

$$\bar{y}_i = \frac{1}{M_i} \sum_{j=1}^{M_i} y_{ij} = \text{ith cluster mean in sample } (i=1(1)n)$$

$$y_i = \sum_{j=1}^{M_i} y_{ij} = \text{ith cluster total in sample } (i=1(1)n)$$

We know,

$$\frac{y_1 + y_2 + \dots + y_n}{n} \text{ estimates}$$

$$\frac{Y_1 + Y_2 + \dots + Y_N}{N}$$

ie.  $\frac{1}{n} \sum_{i=1}^n y_i$  "

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

ie.  $\frac{1}{n} \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n M_i} = \hat{Y}_c$  "

$$\frac{\sum_{i=1}^n Y_i}{\sum_{i=1}^n M_i} = \frac{\sum_{i=1}^n \sum_{j=1}^{M_i} Y_{ij}}{\sum_{i=1}^n M_i} = \bar{Y}$$

Claim:  $\hat{Y}_c$  is an u.e. of  $\bar{Y}$ , where  $\hat{Y}_c = \frac{1}{n} \sum_{i=1}^n \frac{1}{M_i} \sum_{j=1}^{M_i} y_{ij}$

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

Proof: Now,  $E(\bar{y}_i) = \frac{1}{N} \sum_{i=1}^N \bar{y}_i$ , since,  $\bar{y}_i$  takes any value  
 $\bar{y}_i =$  Sample cluster ~~that~~ mean can take any value among  $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_N$

$$E(\bar{y}_i) = \frac{1}{N} \sum_{i=1}^N \bar{y}_i = \frac{1}{M_i} \sum_{i=1}^N \sum_{j=1}^{M_i} y_{ij}$$

Proof: Since,  $E\left(\frac{y_1 + y_2 + \dots + y_n}{n}\right) = \frac{1}{n} \sum_{i=1}^n y_i$

$$\Rightarrow E\left(\frac{N}{n} \sum_{i=1}^n \sum_{j=1}^{M_i} y_{ij}\right) = \sum_{i=1}^n \sum_{j=1}^{M_i} y_{ij}$$

$$\Rightarrow E\left[\frac{N}{n \sum_{i=1}^n M_i} \sum_{i=1}^n M_i \bar{y}_i\right] = \frac{1}{\sum_{i=1}^n M_i} \sum_{i=1}^n \sum_{j=1}^{M_i} y_{ij} = \bar{y}$$

$$\Rightarrow E\left[\hat{Y}_c\right] = \bar{y}$$

$\therefore \hat{Y}_c$  is an u.e. of  $\bar{y}$ .

Variance:

$$\text{Var}\left(\hat{Y}_c\right) = \left(\frac{N}{n \sum_{i=1}^n M_i}\right)^2 \text{Var}\left(\sum_{i=1}^n M_i \bar{y}_i\right)$$

$$= \left(\frac{N}{n \sum_{i=1}^n M_i}\right)^2 \sum_{i=1}^n M_i^2 \text{Var}(\bar{y}_i) \quad \left[\because \text{Cov}(\bar{y}_i, \bar{y}_j) = 0 \text{ for } i \neq j\right]$$

$$= \left(\frac{N}{n \sum_{i=1}^n M_i}\right)^2 \left(\sum_{i=1}^n M_i^2\right) \frac{\sigma_b^2}{n} \left(\frac{N-n}{N-1}\right)$$

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

$$\text{Var}\left(\hat{Y}_c\right) = \left(\frac{N}{\sum_{i=1}^n M_i}\right)^2 \text{Var}\left(\frac{1}{n} \sum_{i=1}^n M_i \bar{y}_i\right) = \left(\frac{N}{\sum_{i=1}^n M_i}\right)^2 \frac{1}{n^2} \sum_{i=1}^n M_i^2 \text{Var}(\bar{y}_i) \left(\frac{N-n}{N-1}\right)$$

$$= \left(\frac{N}{\sum_{i=1}^n M_i}\right)^2 \frac{1}{n^2} \sum_{i=1}^n M_i^2 \sigma_b^2 \left(\frac{N-n}{N-1}\right) \quad \left[\because \text{Cov}(\bar{y}_i, \bar{y}_j) = 0 \text{ for } i \neq j\right]$$

where,  $\sigma_b^2 = \frac{1}{\sum_{i=1}^N M_i} \sum_{i=1}^N (\bar{y}_i - \bar{y})^2 M_i = \text{Var}(\bar{y}_i)$ ,  $i=1(1)n$

$$\begin{aligned} \text{Var}(\bar{y}_i) &= E [M_i(\bar{y}_i - E(\bar{y}_i))]^2 = E [\bar{y}_i - \bar{y}]^2 \\ &= E \left[ \frac{1}{M_i} \sum_{j=1}^{M_i} (y_{ij} - \bar{y}) \right]^2 \end{aligned}$$

$$\left( \frac{N-n}{N-1} \right)$$

$$[n(i=1(1)n)]$$