

VIVEKANANDA COLLEGE  
THAKURPUKUR  
KOLKATA-700063

NAAC ACCREDITED 'A' GRADE



Topic: Theory of Systematic Sampling

Course Title: Sample Survey

Paper: CC-8

Unit: 2

Semester: 4

Name of the Teacher: Riddhi Das Majumder

Name of the Department: Statistics

## Systematic Sampling

### Description: Linear Systematic Sampling:

Suppose that the  $N$  units in the population are numbered 1 to  $N$  in some order. Suppose  $N = nk$ , where  $n$  is the sample size desired and  $k$  is an integer. A number is taken at random from the numbers 1 to  $k$  (using a table of random numbers). Suppose the random number is  $i$ . Then starting from the  $i^{\text{th}}$  unit in the popl., every  $k^{\text{th}}$  unit is selected till a sample of size  $n$  is obtained. Then the sample contain  $n$  units with serial numbers  $i, i+k, i+2k, \dots, i+(n-1)k$ .

Thus the sample consists of the first unit selected at random and every  $k^{\text{th}}$  unit thereafter. It is therefore called a systematic sample (with  $k$  as the sampling interval), and the procedure of selection is known as Systematic Sampling or Linear Systematic Sampling.

For example, when  $N = 24, n = 6$  and  $k = 4$ , the four possible linear systematic samples are:

Sample Number	Random start	Sampled units
1	1	1, 5, 9, 13, 17, 21
2	2	2, 6, 10, 14, 18, 22
3	3	3, 7, 11, 15, 19, 23
4	4	4, 8, 12, 16, 20, 24

The Linear Systematic Sampling scheme described above can be regarded as dividing the popl. of  $N$  units into  $k$  mutually exclusive and exhaustive groups (or clusters)  $\{S_1, S_2, \dots, S_k\}$  of  $n$  units each and choosing one of them at random.

### Composition of the $k$ -Systematic Samples:

Sample number			
1	2	...	$k$
$y_1$	$y_2$	$\dots$	$y_k$
$y_{k+1}$	$y_{k+2}$	$\dots$	$y_{2k}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$y_{n-k+1}$	$y_{n-k+2}$	$\dots$	$y_n$

Means  $\bar{y}_1 \quad \bar{y}_2 \quad \dots \quad \bar{y}_i \quad \dots \quad \bar{y}_k$

Thus the systematic sampling amounts to the selection of a single group or cluster with probability  $k^{-1}$ , from the  $k$  groups forming the entire popl. .

Each unit  $U_i$  in the popl. belongs to one and only one cluster. The probability of selecting a cluster is  $1/k$ , which is therefore

The probability with which any member of the cluster is selected in the sample i.e.  $\pi_i = 1/k$ . This shows that a systematic sampling procedure is a probability sampling procedure. A linear systematic sample is a simple random sample of one cluster unit from a popl. of  $k$  cluster units.

A linear systematic sampling is a mixed sampling, which is partly probabilistic and partly non-probabilistic. This is probabilistic since the 1<sup>st</sup> member of the sample is selected at random (with equal probabilities) from the 1<sup>st</sup>  $k$  units and non-probabilistic since the other members in the sample are fixed by the choice of the first member.

Unbiased Estimator for the Population Total and its variance:

Theorem: An unbiased estimator for the population total  $Y$  under linear systematic sampling corresponding to the random start  $r$  is given by  $\hat{Y}_{st} = \frac{N}{n} \sum_{j=1}^n Y_{r+(j-1)k}$  and its variance is  $V(\hat{Y}_{st}) = \frac{1}{k} \sum_{r=1}^k (\hat{Y}_r - \bar{Y})^2$  where  $\hat{Y}_r$  is the value of  $\hat{Y}_{st}$  corresponding to the random start  $r$ .

Proof: Note that  $\hat{Y}_{st}$  can take any one of  $k$  values  $\hat{Y}_r, r=1, \dots, k$  with prob.  $\frac{1}{k}$ . Therefore  $E(\hat{Y}_{st}) = \sum_{r=1}^k \hat{Y}_r \cdot \frac{1}{k} = \frac{1}{k} \sum_{r=1}^k \frac{N}{n} \sum_{j=1}^n Y_{r+(j-1)k}$   
 $= \frac{N}{nk} \sum_{r=1}^k \sum_{j=1}^n Y_{r+(j-1)k} = \sum_{i=1}^N Y_i = Y$ . Hence  $\hat{Y}_{st}$  is unbiased for the population total  $Y$ .

Again  $V(\hat{Y}_{st}) = E[\hat{Y}_{st} - E(\hat{Y}_{st})]^2 = E[(\hat{Y}_{st} - \bar{Y})^2]$   
 $= \frac{1}{k} \sum_{r=1}^k (\hat{Y}_r - \bar{Y})^2$   
 $E[\hat{Y}_{st}] = \sum_{r=1}^k \hat{Y}_{st} \cdot \frac{1}{k} = \frac{1}{k} \sum_{r=1}^k \frac{N}{n} \sum_{j=1}^n Y_{r+(j-1)k} = \frac{N}{nk} \sum_{r=1}^k \sum_{j=1}^n Y_{r+(j-1)k} = \frac{N}{nk} \sum_{i=1}^N Y_i = \frac{N}{nk} Y = \bar{Y}$

An Alternative Expression for  $V(\hat{Y}_{st})$ :

Th: In linear systematic sampling interval of  $k$ , from a population of size  $N=nk$ , the variance of  $\hat{Y}_{st}$  is given by

$$V(\hat{Y}_{st}) = \frac{N(N-1)}{n} s^2 \{1 + (n-1)\rho\}$$

where  $\rho = \frac{E[(Y_{ij} - \bar{Y})(Y_{i+j-1} - \bar{Y})]}{E[Y_{ij} - \bar{Y}]^2}$  is the

intra-cluster correlation coefficient.

$$= \frac{1}{k} \sum_{r=1}^k (\hat{Y}_r - \bar{Y})^2$$

$$\begin{aligned} \text{if: we have } V(\bar{Y}_{sys}) &= \frac{1}{n} \sum_{i=1}^n (\bar{Y}_i - \bar{Y})^2 = \frac{1}{k} \sum_{j=1}^n \left\{ \sum_{i=1}^n (\bar{Y}_i - \bar{Y})^2 \right\} \\ &= \frac{1}{k} \sum_{j=1}^n \sum_{i=1}^n (\bar{Y}_{ij} - \bar{Y})^2 + \frac{2}{k} \sum_{j=1}^n \sum_{k=1}^n \sum_{i=1}^n (\bar{Y}_{ij} - \bar{Y})(\bar{Y}_{ik} - \bar{Y}) \end{aligned}$$

By definition,  $\rho = \frac{2}{kn(n-1)} \cdot \sum_{j=1}^n \sum_{k=1}^n \sum_{i=1}^n (\bar{Y}_{ij} - \bar{Y})(\bar{Y}_{ik} - \bar{Y}) / \sqrt{V_{sys}}$

$$\Rightarrow \sum_{j=1}^n \sum_{k=1}^n \sum_{i=1}^n (\bar{Y}_{ij} - \bar{Y})(\bar{Y}_{ik} - \bar{Y}) = \frac{kn(n-1)}{2} V_{sys} \rho$$

$$= \frac{1}{2} kn(n-1) \rho \cdot \frac{(N-1) S_y^2}{n} = k \frac{(N-1) S_y^2}{n} \rho$$

Hence,  $V(\bar{Y}_{sys}) = \dots$

$$= k(n-1) S_y^2 + \frac{2}{k} \cdot \frac{kn(n-1)(n-1) \rho S_y^2}{2N}$$

$$= k(n-1) S_y^2 + k(n-1)(n-1) \rho S_y^2 \quad [ \because \frac{n}{N} = k ]$$

$$= k(n-1) S_y^2 \{ 1 + (n-1) \rho \} = N(N-1) S_y^2 \cdot \left\{ \frac{1+n-1}{n} \rho \right\}$$

Conclusion:  $V(\bar{Y}_{sys})$  is systematic sampling be smaller than  $V(\bar{Y}_{srs})$  in srsWOR if  $N(n-1) S_y^2 \cdot \frac{1+n-1}{n} \rho < N(N-1) \frac{S_y^2}{n}$

if  $\rho < \frac{1}{n-1}$

Population with linear trend:

if the values  $Y_1, Y_2, \dots, Y_N$  of the units with labels  $1, \dots, N$  are related by  $Y_i = a + bi, i=1(N),$  i.e. the population constitutes solely a linear trend.

Theorem: For populations possessing linear trend,  $V(\bar{Y}_{sys}) < V(\bar{Y}_{srs})$  where  $\bar{Y}_{sys}$  and  $\bar{Y}_{srs}$  are the usual estimators under linear systematic sampling and simple random sampling, respectively.

Proof:  $V(\bar{Y}_{srs}) = N^2 \cdot \frac{N-1}{nN} \cdot \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$

Let  $Y_i = a + bi, i=1(N),$  then  $\bar{Y} = a + b \left( \frac{N+1}{2} \right)$

Now,  $\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n \left( a + bi - a - b \frac{N+1}{2} \right)^2 = b^2 \sum_{i=1}^n \left( i - \frac{N+1}{2} \right)^2$

$$= b^2 \left[ \sum_{i=1}^n i^2 - N \left( \frac{N+1}{2} \right)^2 \right] = b^2 \frac{N(N^2-1)}{12}$$

$$\text{var. } V(\hat{Y}_{sys}) = \frac{N^2(N-1)}{Nn} \cdot \frac{1}{N-1} \beta^2 \frac{N(N^2-1)}{12} = N^2 \beta^2 \frac{(K-1)(NK+1)}{12} \quad (2)$$

using  $N = nk$ .

$$\text{Again, } \sum_{r=1}^n [\hat{Y}_r - Y]^2 = \sum_{r=1}^n N^2 \beta^2 \left[ r - \frac{K+1}{2} \right]^2 = N^2 \beta^2 \frac{K(K^2-1)}{12}$$

$$\left[ \text{Since } \hat{Y}_r = \frac{N}{n} \sum_{j=1}^n Y_{r+j-1}K = \frac{N}{n} \sum_{j=1}^n \left\{ \alpha + \beta(r+j-1)K \right\} \right. \\ \left. = N \left( \alpha + \beta r + \beta K \frac{n-1}{2} \right) \text{ and } \bar{Y} = \frac{N}{n} \sum_{i=1}^n Y_i = N \alpha + \beta \frac{N(n+1)}{2} \right]$$

$$\text{Then } \hat{Y}_r - Y = \beta N \left[ r + \frac{nk - k - nk - 1}{2} \right] = N \beta \left[ r - \frac{K+1}{2} \right]$$

$$\text{Therefore, } V(\hat{Y}_{sys}) = N^2 \beta^2 \frac{K(K^2-1)}{12}$$

$$\text{Note that } \frac{V(\hat{Y}_{sys})}{V(\hat{Y}_{srs})} = \frac{K^2-1}{(K-1)(NK+1)} = \frac{K+1}{NK+1} < 1, \text{ for all } n > 1.$$

Hence, the linear systematic sampling is more precise than SRS in the presence of linear trend.

### Estimation of Variance in Systematic Sampling :

#### Problems in estimation :

From the results of a Simple Random Sample with  $n > 1$ , we calculate an unbiased estimator of the variance of the sample mean, the estimator being unbiased whatever the form of the popl. Since a systematic sample can be regarded as a simple random sample with  $n = 1$ ; hence the above useful property does not hold for the systematic sample. In other words, since the

variance of  $\hat{Y}_{sys}$  in systematic sampling is  $V(\hat{Y}_r)$  where  $\hat{Y}_r$  is the value of  $\hat{Y}_{sys}$  corresponding to the random start  $r$ , an unbiased estimator of  $V(\hat{Y}_{sys})$  can not be obtained by making just one observation on  $\hat{Y}_r$ , that is, by taking just one systematic sample.

[ It has been pointed out earlier that  $\pi_i = \frac{1}{K}$  and  $\pi_{ij} = \begin{cases} \frac{1}{K}, & \text{if } U_i \text{ and } U_j \text{ are in the same cluster or group.} \\ 0, & \text{if } U_i \text{ and } U_j \text{ are in the different clusters or groups; that is,} \end{cases}$

in linear systematic sampling the second order inclusion probabilities are not positive for all pairs of units in the population. This makes unbiased estimation of the variance of the estimator.

A way out is to make use of the method interpenetrating Subsamples.

Interpenetrating Subsamples:

This technique, particularly useful for the study of correlated errors, was proposed by Mahalanobis (1946). To present it in the simplest terms, a random sample of  $n$  units is divided at random into  $k$  subsamples, each subsample containing  $m = \frac{n}{k}$  units. The field work and processing of the sample are planned so that there is no correlation between the errors of measurement of any two units in different subsamples. For instance, suppose that the correlation with which we have to deal arises solely from biases of the interviewers. If each of  $k$  interviewers is assigned to a different subsample and if there is no correlation between errors of measurement for different interviewers, we have an example of the technique.

Consider the mathematical model for errors of measurement:

Let  $y_{ijr}$  be the value obtained in the  $r$ th replication of the  $j$ th member within the  $i$ th subsample (interviewer). Then

$y_{ijr} = \mu_{ij} + d_{ijr}$ , where  $\mu_{ij}$  is the true mean of the unit and  $d_{ijr}$  is the response deviation on the unit, or the fluctuating component of the measurement error.

[Since the  $i$ th subsample is a random subsample, it is itself a simple random sample of size  $m$ . Hence, the variance of its mean

is  $V(\bar{y}_{ir}) = \frac{1-f}{m} S_p^2 + \frac{\sigma_d^2}{m} [1 + (m-1)\rho_w]$ , assuming

$Cov(\bar{d}_r, \bar{d}_r') = 0$ , where  $\rho_w$  is the correlation coefficient between the  $d_{ijr}$  obtained by the same interviewer (subsample).

If  $\rho_w$  is ignored,  $V(\bar{y}_{ir}) = \frac{1}{m} [S_p^2 + \sigma_d^2 \{1 + (m-1)\rho_w\}]$

Since errors are independent in the different subsample

$V(\bar{y}_r) = \frac{1}{k} V(\bar{y}_{ir}) = \frac{1}{k} [S_p^2 + \sigma_d^2 \{1 + (m-1)\rho_w\}]$

From the sample results we can compute an ANOVA of the  $km$  observations into components between "subsamples (interviewers)" with  $(k-1)$  d.f.'s and "within interviewers" with  $k(m-1)$  d.f.'s.

(47)

ANOVA Table [ON a single unit basis]

Sources of Variation	d.f	m. s.	E(m.s.)
Between interviewers (Subsamples)	$k-1$	$S_b^2 = \frac{m}{k-1} \sum_{r=1}^k (\bar{y}_{1r} - \bar{y}_r)^2$	$S_{\mu'}^2 + \sigma_1^2 \{1 + (m-1)\rho_{ww}\}$
Within Subsamples (Interviewers)	$k(m-1)$	$S_w^2 = \frac{1}{k(m-1)} \sum_{r=1}^k \sum_{j=1}^m (y_{jr} - \bar{y}_{jr})^2$	$S_{\mu'}^2 + \sigma_1^2 (1 - \rho_{ww})$
Total	$km-1$		

Table gives the important result:  $\frac{S_b^2}{n}$  is an unbiased estimator of  $\frac{1}{n} E(S_b^2) = V(\bar{y}_r)$ , ignoring f.p.c. Thus interpenetrated subsamples provide an estimator of  $V(\bar{y}_r)$  that takes proper account of both the simple response variance and correlated component. Note that  $\frac{S_b^2}{n} = \frac{m}{n(k-1)} \sum_{r=1}^k (\bar{y}_{1r} - \bar{y}_r)^2 = \frac{1}{k(k-1)} \sum_{r=1}^k (\bar{y}_{1r} - \bar{y}_r)^2$

Here  $m$  subsamples are interpenetrating in the sense that each is a probability sample over the population.

In linear systematic sampling,  $V(\hat{Y}_{LSS}) = V(\hat{y}_r) = V(N\bar{y}_r) = N^2 V(\bar{y}_r)$ . Hence, by interpenetrating subsamples, an estimator of  $V(\hat{Y}_{LSS})$  is  $V(\hat{Y}_{LSS}) = N^2 \frac{S_b^2}{n} = N \cdot k \cdot S_b^2 = N^2 \frac{1}{k(k-1)} \sum_{r=1}^k (\bar{y}_{1r} - \bar{y}_r)^2$

\* \* \*

Comparison of Linear Systematic with Stratified random Sampling:

Linear Systematic Sampling stratifies the population into  $n$  strata, which consist of the first  $k$  units, the second  $k$  units, and so on. We might therefore expect the linear systematic sample to be about as precise as the corresponding stratified random sample with one unit per stratum. The difference is that with the systematic sample the units occur at the same relative position in the stratum, whereas with the stratified random sample the position in the stratum is determined separately by randomisation within each stratum. The systematic sample is spread more evenly over the population and this fact has, sometimes, made systematic sampling considerably more precise than stratified random sampling.

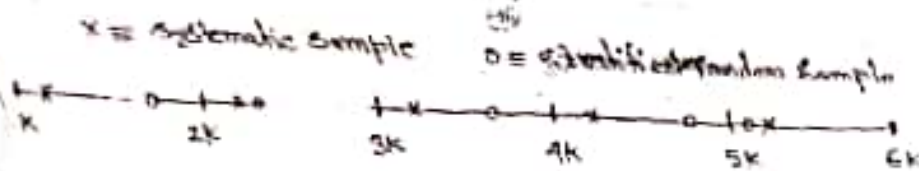


Fig. Systematic and Stratified Random Sampling.

The performance of linear systematic sampling in relation to that of stratified or SRS is greatly dependent on the properties of the population. For some populations and some values of  $n$ ,  $V(\bar{y}_{st}) = \frac{S^2(N-1)}{nN} \{1 + (n-1)\rho\}$  may even increase when a large sample is taken — even a small positive correlation may have a large effect because of the multiplier  $(n-1)$ .

Thus it is difficult to give advice about the situations in which systematic sampling is to be recommended — a knowledge of the structure of the population is necessary for its most effective use.

### Circular Systematic Sampling :

Since  $N$  is not in general an integral multiple of  $k$ , different systematic samples from the same finite population may vary by one unit in size. Thus, with  $N=23$ ,  $k=5$ , the members of the units in the five systematic samples are shown in the table:

Systematic Sample number

I	II	III	IV	V
1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	-	-

The first 3 samples have  $n=5$  and the last two have  $n=4$ . This fact introduces a disturbance into the theory of systematic sampling. (The disturbance is probably negligible if  $n > 50$  and will be ignored, for simplicity.) In some cases, some units of the population will never appear in the sample thereby the estimation of the pop'n total (mean) becomes impossible. For example, when  $N=30$ ,  $n=7$ ,  $k=4$ , the units labels 29 and 30 will never appear as sampled units.

These problems can be overcome by introducing a method, known as circular systematic sampling (CSS).

This method consists in choosing a random start from 1 to  $N$  and selecting the unit corresponding to the random start and thereafter every  $k$ th unit in cyclical manner till a sample of size  $n$  units is obtained.  $k$  being the integer nearest to  $\frac{N}{n}$ . That is, if  $r$  is the number selected at random from 1 to  $N$ , the sample consists of the units corresponding to the numbers  $r + jk$  if  $r + jk \leq N$  and  $(r + jk - N)$ , if  $r + jk > N$ , for  $j = 0, 1, \dots, n-1$ . It is to be noted that, if the sampling interval is taken as the integer closest to  $\frac{N}{n}$ , it is not always possible to get a sample of the given size as shown in the following example. Let  $N = 15$ ,  $n = 6$  and  $k = 3$ . The sample with the random start 3 has only five distinct units namely 3, 6, 9, 12, 15.

Under circular systematic sampling, the conventional expansion estimator  $\hat{Y}_{css} = \frac{N}{n} \sum_{j=1}^n y_j^i$ , where  $y_j^i$  being the value of the  $j$ th unit in CSS corresponding to the random start  $i$ , is unbiased for population total  $Y = \sum_{i=1}^N Y_i$  (whenever  $N$  and  $k$  are relatively prime) and its variance is given by

$$V(\hat{Y}_{css}) = \frac{1}{N} \sum_{i=1}^N [\hat{Y}_{ci} - Y]^2$$