

UNIT - II

Varying Probability Sampling: / ~~Prob. Proportion~~ Proportion to
When units vary in their sizes and the
Variate under study is highly co
(Prob. proportion to size Sampling).

When samples from different sizes sub-
grouped are used and sampling is taken with
the sample probability the chance of selecting
a member from a large group or less than
selecting a member from a smaller group.
This is known as probability proportional to size
Sampling. (when n is in different sizes).

For eg:

If one sample had 20,000 members, the
prob of a member being selected would be
 $1/20,000$ or 0.005%. If another sample
had 10,000 members, the chance of a member
being selected would be $1/10,000$ or 0.01%.

Difference between S.M.S SRS and Varying PS:

In simple random Sampling the
probability of drawing a specified unit

at any given step is the same. In varying Prob. Scheme the probability of drawing a specified unit differ from. how to draw.

Procedures of Selecting a Sample:

The Procedure of Selecting a Sample consist in associating with each unit a number or set of numbers equal to its size. The selection of units is done corresponding to a no. chosen at random from the totality of numbers associated. There are two methods of Selection.

1. Cumulative total method

2. Lohuris' method

* Let the size of the i^{th} unit be ($i=1, 2, 3, \dots, N$) the total being $X = \sum_{i=1}^N x_i$

* We associate the number 1 to x_1 with the first unit, the numbers $(x_1 + i)$ to $(x_1 + x_N)$ with the second unit, and so on.

* A number k is chosen at random from 1 to X and the unit with which this number is associated is selected.

~~in which it is~~

- * The i^{th} unit in the population is being selected with a prob. Proportional to x_i .
- * If a Sample of size n is required, the procedure is repeated n times with replacement of the units selected.
- * This procedure of selection is known as the cumulative total method for the method needs cumulation of the unit sizes.
- * The main difficulty in this procedure is the compulsion to complete successive cumulative totals, which becomes time consuming and costly when the population size is large.

Example:

A village has 10 holdings consisting of 50, 30, 45, 25, 40, 26, 24, 35, 28 and 27 fields respectively. Select a sample of four holdings with the replacement method and with probability proportional to the number of fields in the holdings.

The first step in the selection of holdings is to form cumulative totals as per following

Eg:

S. No. of holdings	Size (x_i)	Cumulative Size	Numbers associated
1	50	50	1-50
2	30	80	51-80
3	45	125	81-125
4	25	150	126-150
5	40	190	151-190
6	26	216	191-216
7	44	260	217-260
8	35	295	261-295
9	28	323	296-323
10	27	350	324-350

To select a holding, a random number not exceeding 350 is drawn with the help of a random number table. Suppose the random number thus selected is 272.

It can be seen from the cumulative totals that the number is associated with the group 261-295, i.e., the 8th holding is selected corresponding to the random number 272.

Similarly, we select three more random numbers. Suppose these numbers are 346, 165 and

Lahiri's method:

Lahiri (1951) suggested an alternative procedure in which cumulations are used completely.

It consists in selecting a number at random between 1 and N and noting down the unit with the corresponding serial number, provisionally. Another random number is then chosen between 1 and M , where M is the maximum size of the N units of the population.

If the second random number is smaller than the size of the unit provisionally selected, the unit is selected into the sample. If not, the entire procedure is repeated until a unit is finally selected. For selecting a sample of n units, the procedure is to be repeated until n units are selected.

Eg: Select a sample of 8 from holding in the previous example by Lahiri's method of PPS, wr.

In this case, $N=10$, $M=50$.

First we have to select a random number which is not greater than 10 and a second random ~~table~~^{number}, which is not greater than 50.

Referring to the random number table, the pair is (10, 13). Hence, the 10th unit is selected in the sample. Similarly, choosing other pairs, we can have (4, 26), (5, 35), (7, 26).

The pair (4, 26) is rejected as 26 is greater than the size value (25) and so another pair is drawn which turns out to be (8, 16).

Hence, the sample will consist of the holdings with serial numbers 10, 5, 7 and 8.

PROCEDURES OF SELECTION OF A PPS SAMPLE WITHOUT REPLACEMENT: $P_i = \frac{f_i}{N}$.

There are several procedures for selecting samples with unequal prob. A brief discussion of some important procedures are as follows.

General Selection Procedure:

A generalization of the Sampling Scheme, ~~whereas~~^{without replacement (cor.)} there would be to select a PPS sample of size unity and remove the selected unit from the population.

From the remaining units, another PPS Sample of size one is taken as before and the Selected unit removed from the population. This process is repeated until n Selections.

If n units are selected one by one, with Probability Proportional to Size measure x_i , at each draw, without replacing the units selected in the previous draws. The Probability of selecting a first draw for the i^{th} unit is

$$P_j = x_j / X, \quad j=1, \dots, N.$$

$$\text{Where } X = \sum_{j=1}^N x_j.$$

Similarly, the Probability that the i^{th} unit is selected at the second draw when the j^{th} unit has been selected at the first draw, is given by

$$P_{ij} = P_j / (1 - P_j) \quad i \neq j, \text{ and so on.}$$

This Set up of Sampling comprises an ordered set of Sample Values (y_1, y_2, \dots, y_n) with Probabilities (P_1, P_2, \dots, P_n) .

Norain's Scheme of Sample Selection:

The Scheme consists of constructing desired Prob. of selection p_i ($i=1, 2, \dots, N$) such that the Inclusion Probabilities P_i are proportional to the original Probabilities of selection p_i ($i=1, 2, \dots, N$), and Sampling is done without replacement.

The Inclusion Probabilities are given by

$$\pi_i^o = np_i, \quad i=1, 2, \dots, N.$$

The Probabilities of selection and inclusion at the second and subsequent draws are proportional to revised Probabilities p_i' on lines similar to those at the first draw. Let us consider the simple case $N=4$ and $n=2$.

The Problem is to Evaluate π_{ij}^o gives $p_i (i, j=1, 2, 3, 4)$.

The relationship $\sum_{j \neq i} \pi_{ij}^o = \pi_i^o$ provides a system of 4 linear equations with 8 unknowns. The Problem is to choose the values of two arbitrary parameters with the restriction that all the above values are positive. The Computations become tedious for n greater than 2.

Example for General Selection Procedure:

In a village, there are 8 Orchards with 50, 30, 25, 40, 26, 44, 20 and 35 trees, respectively,

Select a Sample of 2 orchards with probability proportional to the number of tree in the orchard and without replacement.

(i) For Selecting a unit by Prob. Proportional to the Number of trees, using Lahiri's method of Selection, Consider the following arrangement:

Orchard number : 1 2 3 4 5 6 7 8

Number of tree : 50 30 25 40 26 44 20 35

Selecting a pair of r. numbers (i,j) ($i \leq 8, j \leq 50$)
 Using the number table, we get the pair $(5, 17)$. Since the
 number ~~x_5~~ of the trees (x_i) for orchard number 5
 is greater than the second number (17) of the
 selected random pair, the 5th orchard is selected
 in the sample.

② For selecting the second unit by probability
 proportional to the number of trees in the orchard,
 we prepare the following arrangement once again
 after deleting the 5th orchard:

Orchard number \Rightarrow . See above.

Number of trees \Rightarrow

As in ① a pair of random numbers (i,j) ($i \leq 7, j \leq 50$)
 has to be selected, using the random number table.
 Referring to the table of random numbers, the
 pair selected is $(6, 18)$. As the size of the 6th unit in
 the above arrangement is greater than the second
 number of the random pair selected, orchard 6 is
 selected into the sample. Thus, the sample selected
 consists of the units at serial numbers 5 and 7 of
 the original list with the number of trees being
 26 and 20, respectively.

Sen-Midzuno Method:

A simple procedure of selecting a sample was suggest by Midzuno and independently by Sen (1952), which consists in selecting the first unit with p_0 and the remaining $(n-1)$ units from $(N-1)$ units of the population by simple random sampling.

For this selection procedure, the inclusion probabilities for individual and pairwise units are given by.

$$\pi_i^0 = p_0 + (1-p_0) \frac{n-1}{N-1} \quad \text{for } i=1, 2, \dots, N.$$

$$= \frac{N-n}{N-1} \quad p_0 + \frac{n-1}{N-1}$$

and $\pi_{ij}^0 = p_0 \frac{n-1}{N-1} + p_j = \frac{n-1}{N-1} + (1-p_0 - p_j)$

$$= \frac{(n-1)(n-2)}{(N-1)(N-2)} \quad \text{for } i \neq j = 1, 2, \dots, N.$$

$$= \frac{(n-1)}{(N-2)} \left[\frac{(N-n)}{(N-2)} (p_0 + p_j) + \frac{(n-2)}{N-2} \right]$$

By extension of the above argument, we can have y_1, y_2, \dots, y_n , a sample of n units.

The Probability of including these n units in the sample is given by

$$P_{ij} \dots q = \frac{1}{(N-1)} (p_i + p_j + \dots + p_n).$$

ESTIMATION IN PPS SAMPLING WITHOUT REPLACEMENT: Total and its Sampling Variance:

ORDERED ESTIMATORS:-

Those estimators which are based on the order of units selected in the sample and do not require calculations of inclusion probabilities.

Das (1951) and Des Raj (1956) have proposed such estimators those make use of conditional Probabilities without calculating P_i and P_{ij} which are generally difficult to compute for many sampling schemes. We shall consider the estimator proposed by Des Raj, for the case when $n=2$.

DesRay's Order estimator:

Consider, Case (i) $n=2$ and the result is generalized.

* Let the initial probability of units i and P_i :
 $i = 1, 2, \dots, N$ where $P_i = x_i / X$.

* The first draw is made with Prob. P_i , the second draw is taken with conditional Prob. $P_j P_i(1-P_i)$.

* Suppose y_1 and y_2 are the value of units drawn at the 1st and 2nd draw respectively and P_1 and P_2 are their corresponding initial probability.

Define the estimators

$$z_1 = y_1 / P_1$$

$$z_2 = y_2 / [y_2(1-P_1) / P_2]$$

Then, $\hat{Y} = (z_1 + z_2) / 2$

$$= [y_1(1+P_1) / P_1 + y_2(1-P_1) / P_2]$$

$$= \frac{y_1}{P_1} [1+P_1] + \frac{y_2}{P_2} [1-P_1]$$

2.

Theorem:

In PPS SRSWOR, the estimator \hat{Y}_D is an unbiased estimator of sampling Variance:

$$V(\hat{Y}_D) = \left(1 - \frac{1}{N} \sum_i p_i^2\right) \left[\frac{1}{2} \sum_i^N (y_i/p_i - \bar{Y})^2 p_i \right] - \frac{1}{4} \sum_i^N (y_i/p_i - \bar{Y})^2 p_i^2.$$

Proof:

Either $Z_1 = y_1/p_1$.

$$(i.e) E(Z_1) = \sum (y_i/p_i) = \sum y_i = Y.$$

Similarly for y_1 drawn at the first draw

$$\text{we have, } E\left(y_2 (1-p_1)/p_2 | y_1\right) = \sum y_2 \frac{(1-p_1)}{p_2} \frac{p_1}{(1-p_1)}.$$

where \sum is taken over all the values of y_2

Expect y_1 .

$$\text{Hence } E\left(y_2 \frac{(1-p_1)}{p_2} | y_1\right) = Y - y_1$$

$$\Rightarrow E(Z_2) = E\left(E[Z_2 | y_1]\right)$$

$$= y_1 + Y - y_1$$

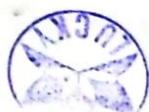
$$= Y.$$

In General, $E[Z_0] = Y$ thus we can be
Show that

To find out Variance.

$$V[Z] = \sum_i \sum_j p_i p_j \left(\frac{y_i}{p_i} - \frac{y_j}{p_j} \right)^2$$

~~Lucky Spinning Co., LTD.~~



also $V(Z_2) = E_1 V_2(Z_2) + V_1 E_2(Z_2)$

Since,

$$E_2(Z_2) = Y, V_1 E_2(Z_2) = 0.$$

Hence we have

$$V(Z_2) = \sum_i \sum_j p_i p_j \left(\frac{y_i}{p_i} - \frac{y_j}{p_j} \right) \cdot (1 - p_i - p_j)$$

thus,

$$V(Y_0) = \frac{1}{4} \sum_i \sum_j p_i p_j \left(\frac{y_i}{p_i} - \frac{y_j}{p_j} \right) (2 - p_i - p_j)$$

.....
Final

Theorem ②.

In PPS SWR the estimator \hat{y}_D is an unbiased estimator of Population total y_D and its Sampling Variance is given by

$$V(\hat{y}_D) = \frac{1}{n^2} \sum_{i=1}^n p_i p_{ij} \left[\frac{y_i}{p_i} - \frac{\bar{y}_D}{\bar{p}_D} \right]^2 \left\{ r_{ij}^{(1)} + \dots + r_{ij}^{(n-1)} \right\}$$

Where $r_{ij}^{(k)}$ is the prob. that y_0 & y_j are not included in the sequence.

Proof:

$$E(Z_i) = y \text{ and } E[Z_i | y_1, y_2, \dots, y_{i-1}] = y, \quad \text{where } i=2, 3, \dots, n$$

hence

$$E(Z_i) = y \text{ for } i=2, 3, \dots, n$$

It follows that

$$\hat{y}_D = \sum_{i=1}^n Z_i / n \text{ is an unbiased estimator.}$$

To derive Variance :-

$E(Z_i Z_j) = y^2$; it shows that Z_i and Z_j are uncorrelated.

* the expression for $V(\hat{y}_D)$ is somewhat complex, it can be modified to a simpler form

$$\text{as } V(\hat{y}_D) = V \left[\sum_{i=1}^n \frac{Z_i}{n} \right] = \frac{1}{n^2} \sum_{i=1}^n V(Z_i).$$

Unordered Estimators:

The unordered estimators which do not depend on a order in which the units of drawn within the sample.

Hornitz and Thompson (1952), Knoottho (1957) and Baen (1958), have shown that these uncorrelated estimators are more efficient than correlated estimators.

Hornitz and Thompson Estimator:-

Hornitz and Thompson (1952) suggested an estimator which is an unbiased estimator of population total. Let the initial probability of selection of the unit U_i is p_i where

$$p_i = \frac{x_i}{X} \quad \text{for } i=1, 2, \dots, n.$$

The Prob. that unit U_i & U_j is included in the sample would be given by

$$\Pi_{ij} = p_i + \sum_{j \neq i} p_j p_i / (1 - p_j).$$

$$= p_i \left[1 + \sum_{j \neq i} p_j / (1 - p_j) \right].$$

From that the Prob. that both a units U_i and U_j are included in the sample is given by

$$\Pi_{ij} = p_i p_j / (1-p_i) + p_i p_j / (1-p_j)$$

$$= p_i p_j \left[\frac{1}{1-p_i} + \frac{1}{1-p_j} \right]$$

Suppose that y_e be the value of the i^{th} unit
with Π_i , the prob. of inclusion in the
sample. Horvitz-Thompson estimator is defined by

$$\hat{Y}_{HT} = \sum_i y_e / \Pi_i$$

Theorem:

In PPSWOR \hat{Y}_{HT} is unbiased and its Sampling
Variance is given by.

$$V_{HT}(\hat{Y}_{HT}) = \sum_i \frac{(1-\Pi_i)}{\Pi_i} y_i^2 + \sum_i \sum_{j \neq i} \frac{(\Pi_{ij} - \Pi_i \Pi_j)}{\Pi_i \Pi_j} y_i y_j$$

where Π_{ij} is the prob. of inclusion of both the i^{th}
and j^{th} units in the sample.

Proof:

The most general form of linear estimator

$$\hat{Y} = \sum_i a_i y_i$$

Where, a_i is a random variate that takes the
values 1 if the i^{th} unit is drawn and 0 otherwise.

c_i are constants attached to the units U_i ,
 $i=1, 2, \dots, N$.

Obviously a_i follows a binomial distribution
a sample of size 1 with prob π_i .

$$\text{Hence, } E(a_i) = \pi_i$$

$$V(a_i) = \pi_i \times (1 - \pi_i)$$

~~LUCKY SPINNING CO. LTD.~~



Since $a_i \cdot a_j = 1$ only if both the units are
distinct and appear in the sample, covariance of

$$\text{Cov}(a_i, a_j) = \pi_{ij} - \pi_i \pi_j$$

$$\text{Now, } E(\cdot) = E\left(\sum_i^N a_i (i y_i)\right)$$

$$E(\hat{y}) = E\left[\sum_{i=1}^N a_i e_i y_i\right]$$

$$E(\hat{y}) = y.$$

If π_i is unbiased

$$\therefore \pi_i = 1/\bar{\pi}_p.$$

Hence $\hat{y}_{HT} = \sum_{i=1}^N y_i / \bar{\pi}_p$ is an unbiased estimator.

$$V_{HT}\left(\hat{y}_{HT}\right) = V_{HT}\left(\sum_{i=1}^N a_i e_i y_i\right)$$

$$V_{HT}\left(\hat{y}_{HT}\right) = \sum_{i=1}^N \frac{y_i^2}{\bar{\pi}_p^2} V(a_i) + \sum_{i=1}^N \sum_{j \neq i} \frac{y_i y_j}{\bar{\pi}_p \bar{\pi}_p} \text{cov}(a_i, a_j).$$

$$= \sum_{i=1}^N \pi_i (1 - \pi_i) \frac{y_i^2}{\bar{\pi}_p^2} + \sum_{i=1}^N \sum_{j \neq i} \frac{(\pi_{ij} - \pi_i \pi_j)}{\bar{\pi}_p \bar{\pi}_p} y_i y_j$$

Note:

① The $V(\hat{y}_{HT})$ depends on $\bar{\pi}_p$ and π_i which are calculated from Sampling Procedure.

② $\bar{\pi}_p = ny_p/y$, $V(\hat{y})$ becomes zero.

③ If there is some way of choosing values of π_i so that they are very close to y_p/y , you will have a very small variance.

Test:

1) An unbiased estimator of $V(\hat{y}_{HT})$ is given by

$$V(\hat{y}_{HT}) = \sum_{i=1}^n (1 - \pi_i) \frac{y_i^2}{\bar{\pi}_p^2} + \sum_{i=1}^n \sum_{j \neq i} \frac{(\pi_{ij} - \pi_i \pi_j)}{\bar{\pi}_p \bar{\pi}_p} \frac{y_i y_j}{\bar{\pi}_p \bar{\pi}_p}.$$

Provided that none of the π_{ij} in the population is zero.

Another expression for $V(\hat{Y}_{HT})$ derived by Yates and Gouraudy (1923) is given by,

$$V_{Y_{HT}}(\hat{Y}_{HT}) = \sum (T_{ij}T_{ij} - \bar{T}_{ij}\bar{T}_{ij}) \left(\frac{Y_{ij}}{T_{ij}} - \frac{\bar{Y}_{ij}}{\bar{T}_{ij}} \right)^2 \text{ and}$$

If it is unbiased estimator is given by

$$V_{Y_{UO}}(\hat{Y}_{HT}) = \sum_{i,j} \frac{(T_{ij}\bar{T}_{ij} - \bar{T}_{ij}\bar{T}_{ij})}{\bar{T}_{ij}} \left(\frac{Y_{ij}}{T_{ij}} - \frac{\bar{Y}_{ij}}{\bar{T}_{ij}} \right)^2$$

Moorthy's unordered estimator:-

Moorthy (1957) suggested that an unordered estimator can be obtained by weighting all possible ordered estimator with their respective prob's. In sampling n units w.r.t. from a finite popn, there will be $\binom{N}{n}$ unordered samples. Each unordered sample of size n can be ordered in $M (= n!)$ ways: (ie) an unordered sample correspond to M ordered samples. Consider a scheme of selection in which the prob of selecting the sample s_i is $p(s_i)$. Then the prob of getting the unordered sample s is the sum of prob of getting the ordered samples (s) corresponding to (s) (ie)

$P_s = \sum p(s_i)$. Let y_{si} be an estimator of popn parameter α based on the ordered sample s_i an unordered

estimate α is given by $\hat{Y}_m = \sum_{i=1}^n y_{si} p'_{si}$; where $p'_{si} = P_{ri}/P_s$.

Theorem(5):-

In pps SwoR the unordered estimator \hat{Y}_m is an unbiased estimator of α . and it's sampling variance is given by

$$V(\hat{Y}_m) = \sum_{i=1}^n \sum_{j=1}^n p_{sj} \left[\sum_{i=1}^n y_{si} p'_{si} \right]^2 - \left(\sum_{i=1}^n \sum_{j=1}^n y_{sj} p'_{sj} \right)^2$$

Result :-

1) In pps SwoR the ordered estimator y_{si} is an unbiased estimator of α and its sampling variance is given by

$$V(y_{si}) = \sum_{i=1}^n \sum_{j=1}^n p_{sj} y_{sj}^2 - \left(\sum_{i=1}^n \sum_{j=1}^n y_{sj} p_{sj} \right)^2$$

The variance of an unordered estimator \hat{Y}_m is less than or equal to that of the ordered estimator, because

$$V(y_{si}) - V(\hat{Y}_m) = \sum_{i=1}^n \sum_{j=1}^n p_{sj} y_{sj}^2 - \sum_{i=1}^n \sum_{j=1}^n \left[\sum_{k=1}^n y_{ki} p'_{ki} \right]^2$$

2) An unbiased estimator of the sampling variance of \hat{Y}_m is given by

$$\begin{aligned} V(\hat{Y}_m) &= \frac{1}{(P_s)^2} \sum_{i=1}^n \sum_{j=1}^n \left[P_s P_{ri} (p_{ri} - p_{sj} p_{sj}) \right] \\ &\quad - p_i p_j \left(\frac{y_i}{p_i} - \frac{y_j}{p_j} \right)^2 \end{aligned}$$

Where, $p(s_i, j)$ is a conditional prob. of getting the s th sample given that the units i & j are already selected in the first two draws.

PPS Systematic Sampling:

This PPS Systematic Sampling is always better than PPS SWR. Madow (1949) suggested that generalization can be made out sys. sampling with equal probs.

The method consists of arranging the units at random and frequency cumulative totals

$$T_i = \sum_{j=1}^i x_j \quad (i = 1, 2, \dots, N)$$

where x_j is size of j th unit. In selecting n units by PPS Systematic Sampling, a random number R is chosen from 1 to K , where $K = \frac{T_N}{n}$. The units corresponding to the no's $R + jk$ ($j = 0, 1, 2, \dots, (n-1)$) ~~constitute~~ constitute the sample.

(Ex),

Consider the example give in cumulative total here, $T_N = \sum_{j=1}^N x_j = 330$, and $K = \frac{330}{8} = 41.25$
 $K \approx 42$

Suppose the r.n choosen below 1 to 82 is ~~ex~~ 51.

Then the units corresponding to the sys sample numbers.

$$R+jk$$

$$R = 51, j = 0, 1, 2 \dots (n-1)$$

$$R+k = 51 + 82 = 133$$

$$R+2k = 51 + 2 \times 82 = 215$$

$$R+3k = 51 + 3 \times 82 = 297.$$

51, 133, 215, 297 are 2, 4, 6 and 9. It is clear from the procedure that i^{th} unit is selected in the sample if

$$T_{p-1} < R+jk \leq T_p \text{ for } j = 0, 1, 2 \dots (n-1)$$

and the prob. of inclusion of i^{th} unit is npi .

This method can be applied even when $\frac{T_n}{n}$ is not an integer. The Sampling interval k can be taken as integer nearest to $\frac{T_n}{n}$.

An unbiased estimator of Popln total Y , derived by Hartley and Rao (1962) is given by

$$\hat{Y}_{HR} = \sum_{i=1}^{n_1} \frac{y_i}{n_i} = \frac{n_1 y_1}{n pi}.$$

The main drawback of this method, as in sysr. Sampling its that, an unbiased variance estimator on the basis of single sample is

not possible Hartley and Rao (1962) have further considered this selection procedure when the units are arranged at random and asymptotic approach has been used in developing formulae for variance and estimated variance. For large values of N and for value of n relatively small compared to N , the ~~approx~~ approx sampling variance of the estimator can be written as

$$\sqrt{(\hat{Y}_{HR})} \approx \frac{N}{2} \left(\frac{y_i}{p_i} - \bar{y} \right)^2 p_i [1 - (n-1)p_i]$$

which shows that even the units are arranged from random, selecting them with PPS Sys Sampling is more efficient than PPS SWR for the terms $[1 - (n-1)p_i]$ are acting as reduction terms.

Further this variance can be estimated

$$\text{by } \text{re}(\hat{Y}_{HR}) = \frac{1}{n^2(n-1)} \sum_i^n \sum_{j \neq i}^n \left[1 - n(p_i + p_j) + \sum_j^n p_j^2 \right] \left(\frac{y_i}{p_i} - \frac{y_j}{p_j} \right)^2$$