

ALTERNATIVE SAMPLING STRATEGY BASED UPON COEFFICIENT OF VARIATION WHEN AUXILIARY INFORMATION IS AVAILABLE

SPERSH BHATT, M.K.PANDE

Abstract: It is well know that in order to improve the efficiencies of the estimates probability sampling is preferred over non probability sampling. If the difference in the size of the units is large enough to affect the study, we make use of PPS sampling where the probability of selecting a unit is proportional to the size measure of the unit. Sometimes we may be confronted with situations where information on a character closely related to the main variable is available from a previous study or other secondary sources. Various authors have utilized this auxiliary information by taking the initial probability of selection equal to the size measure of the auxiliary information. This scheme however fails to give best results when the population under consideration is skewed. The paper presents an alternative without replacement sampling strategy obtained by utilizing auxiliary information to modify the initial probability of selection of the units. A computer program was developed in Visual Basic to find out the probabilities of selection and the variance of the sampling strategy proposed using the Horvitz Thompson estimator of population total. The empirical comparison of the proposed strategy with the existing Midzuno-Sen strategy shows that the proposed scheme performs better than the Midzuno-Sen strategy when the population is skewed.

Keywords: coefficient of variation, probability of selection, relative efficiency, sampling strategy, sampling design

Introduction: It is well known that to avoid personal bias, random sampling is preferred over non random sampling. Attaching equal probability of selection to different units yields the method of simple random sampling. When unequal probabilities of selection are attached to different units in the population, it is called unequal probability sampling. If a sample from a finite population is drawn, usually the values of some character 'x' closely related to the main character of interest is available for all units of the population. The variable 'x' which is suitably normed, is often taken as a measure of the size of the unit. This occurs in socio-economic, agricultural and industrial surveys which are accompanied with the knowledge of past data. A unit with higher values of 'x' shall contribute more to the population total of main variable, than those with smaller sizes. One expects that, a selection procedure which gives higher selection probabilities to bigger units than to smaller units, should be more efficient than simple random sampling.

Consider a finite population 'U' of distinguishable units labeled 1,2,3,.....N. The collection of all possible samples is called the sample space denoted by 'S'. With each sample 's' a probability p(s) is attached which is the probability of drawing the sample 's'.

We thus have

- (1) $p(s) \geq 0$
- (2) $\sum_{s \in S} p(s) = 1$

Here the sample from 'U' is an ordered sequence of labels from 'U' and represented by $S = (i_1, i_2, \dots, i_n)$ where i_k is the label of the unit drawn at the k^{th} draw and $1 \leq k \leq n$. The labels represent the units drawn with or without replacement in 'n' consecutive draws, hence the labels need not be distinct from each other. The size of the sample is 'n' and 'r' is the effective sample size(which is the number of distinct labels in 'S').

Let P_i denote the probability that the i^{th} unit is selected in

the sample from the population .
By the addition law of probability

$$P_i = \sum_{s \in S} p(s)$$

where summation is taken over all possible samples containing the i^{th} unit of the population. It is further assumed that $p(s)$ is such that $P_i > 0$ for $i = 1, 2, \dots, N$.

The collection $S = \{s\}$ with a probability measure $P = \{p(s)\}$, defined on 'S', such that $p(s) \geq 0$ and $\sum_{s \in S} p(s) = 1$ is called the sampling design and is

denoted by D(S,P). A sampling procedure in which P_i (the probability of including the unit i in a sample of size n) is πp_i . These are referred to as π -ps methods. Here p_i is the probability of selecting the i^{th} unit of the population into the sample at the first draw. To estimate the population mean or total with such procedures, the commonly used estimator is the Horvitz-Thompson (H-T) estimator.

The unbiased H-T estimator for population total Y can also be written as

$$\hat{Y}_{HT} = \sum_{i=1}^N \frac{Y_i \delta_i}{P_i}$$

where $\delta_i = \begin{cases} 1, & \text{if } i \in S \\ 0 & \text{otherwise} \end{cases}$

the variance of the H-T estimator for population total Y is given b

$$V(\hat{Y}_{HT}) = \sum_{i=1}^N \frac{1-P_i}{P_i} Y_i^2 + 2 \sum_{i < j} \frac{(P_{ij} - P_i P_j)}{P_i P_j} Y_i Y_j$$

where $i, j = 1, 2, 3, \dots, N$.

Here P_{ij} is the probability of including the units i and j in the sample and

$$P_{ij} = \sum_{s \in S} p(s)$$

Yates and Grundy(1953) provided an alternative estimator

of the population total Y , which is given by

$$V(\hat{Y}_{HT})_{YG} = \sum_{i \neq j}^N (P_i P_j - P_{ij}) \left(\frac{Y_i}{P_i} - \frac{Y_j}{P_j} \right)^2$$

Some estimators of variance of the Horvitz-Thompson estimator have been given by Yates and Grundy and Sen (1953), Jessen(1969) and Ramakrishnan(1971).

Midzuno(1952)developed a sampling strategy in which the unit at the first draw is selected with unequal probability of selection. At all subsequent draws they are selected with equal probability and without replacement.

In the Midzuno-Sen scheme of probability proportional to size (pps) sampling the probability that the ith unit is included in the sample is given by

$$\frac{(N-n)}{(N-1)} P_i + \frac{(n-1)}{(N-1)}$$

and the probability that both ith and jth units are included in the sample is given by

$$\frac{n-1}{N-1} \left[\frac{(N-n)}{(N-2)} (P_i + P_j) + \frac{n-2}{N-2} \right]$$

In the above scheme the probability of selection for a specific size 's' is given by

$$\frac{\sum_{i \in s} x_i}{\sum_{s \in S} \left(\sum_{i \in s} x_i \right)}$$

Midzuno's scheme made π -ps has been considered by Rao(1963), Sankaranarayanan(1969), Chaudhuri(1974), Mukhopadhyay(1974) among others.

The Midzuno scheme though easy to implement is known to be less efficient in comparison to other unequal probability schemes. On the other hand Sampford scheme is known to be usually a good performer in the class of unequal probability schemes. This scheme however suffers from the drawback that it is rather difficult to implement particularly when $n > 2$.

Section 2 of the paper presents the methodology of obtaining the proposed strategy and the empirical study used for comparing the proposed strategy with the conventional Midzuno-Sen scheme.

Section 3 of the paper gives the tables and graphs giving the variance comparison of the Horvitz Thompson estimator under the proposed scheme and the Midzuno scheme.

Methodology and empirical study: The coefficient of variation(C.V), for auxiliary information ,is given by s.d./mean, where s.d.(standard deviation) is a measure of dispersion and mean is an average. It is proposed to take Range as the measure of dispersion and geometric mean, given by

$$\bar{x}_g = \left(\prod_{i \in s} X_i \right)^{\frac{1}{n}} \text{ as average.}$$

Thus we may take

$$C.V = \sum_{i \in s} \frac{\max(X_i) - \min(X_i)}{\bar{x}_g}$$

We now may take the probability of selection for a specified sample of size s based upon the C.V as

$$p(s) = \frac{C.V.}{\sum_{s \in S} [C.V.]} \dots \dots \dots (\alpha)$$

It can be easily shown that the Horvitz-Thompson estimator under the above scheme is unbiased for the population total.

The empirical comparison for variance under the Midzuno scheme and the proposed one based on Coefficient of Variation has been done using a computer program developed in Visual Basic.

P_i and P_{ij} have been calculated on the basis of (α) and then the following Yates –Grundy formula for variance is used

$$V(\hat{Y}_{HT})_{YG} = \sum_{i \neq j}^N (P_i P_j - P_{ij}) \left(\frac{Y_i}{P_i} - \frac{Y_j}{P_j} \right)^2$$

To compare the two schemes 10 natural populations have been considered. These have been taken from Murthy (1977). Here Y stands for the number of cultivators in 1961 and X for area in 1951.

Five cases have been considered for N=7 and n=3, two cases for N=8 and n=3 and three cases for N=9 and n=3. Further details regarding the natural populations and the P_i and P_{ij} values for these populations, as computed for the sampling design, using (α) maybe obtained from the author as the detailed description is not possible due to bravity.

Tables and figures:

Table 1Description of Natural Populations for N=7, n=3

Sl.No.	Natural Population No.	Source	Variance M-S	Variance C.V.	% Relative efficiency of the estimator of C.V. scheme over M-S scheme
1.	1	Murthy,(1977),Pg.127	7857190.04	3153547.76	249.16
2.	2	Ibid, Page 127	8321296.48	4735136.92	175.74
3.	3	Ibid, Page 127	3341035.17	1397001.59	239.16
4.	4	Ibid, Page 127	4038465.53	2836361.92	142.38
5.	5	Ibid, Page 127	8960843.40	5935520.55	150.97

Table 2 Description of Natural Populations for N=8, n=3					
Sl. No.	Natural Population No.	Source	Variance M-S	Variance C.V.	% Relative efficiency of the estimator of C.V. scheme over M-S scheme
1.	1	Murthy,(1977), Pg. 129	1205843.96	1012225.28	119.13
2.	2	Ibid, Page 129-130	1849844.39	1653255.34	111.89

Table 3 Description of Natural Populations for N=9, n=3					
Sl. No.	Natural Population No.	Source	Variance M-S	Variance C.V.	% Relative efficiency of the estimator of C.V. scheme over M-S scheme
1.	1	Murthy,(1977),Pg.127	18053182.83	4805651.92	375.67
2.	2	Ibid, Page 127	24616710.38	6703196.98	367.24
3.	3	Ibid, Page 127	11123097.95	2804140.52	396.67

Conclusion: From the tables given above it can be easily concluded that the proposed scheme performs better in all the cases when the population is skewed.

References:

- Chaudhuri, A. (1974): On some properties of the sampling scheme due to Midzuno, Bull. Cal. Stat. Assoc., 18, 1-24.
- Horvitz, D.G. and Thompson, D.J. (1952): A generalisation of sampling without replacement from a finite universe. Jour. Amer. Stat. Assoc., 47, 663-85.
- Jessen,R. J. (1969): Some methods of probability non-replacement sampling. Jour. Amer. Stat. Assic.
- Midzuno,H. (1952): On the sampling system with probability proportional to sum of sizes. Ann. Inst. Stat. Math, 3, 99-107.
- Mukhopadhyay, P (1974): π ps sampling schemes to base HTE. Cal. Stat. Assoc. Bull, 23 ,21-44.
- Murthy, M.N.(1977):Sampling theory and Methods, 2nd ed., Statistical Publication Society Calcutta.
- Ramakrishnan, M. K. (1971): Generalisation of Yates and Grundy estimates. Research Rep. 92/MKR/I, University of Sheffield, Manchester, Sheffield.
- Rao, J. N. K. (1963): On three simple procedures of unequal probability sampling without replacement. Jour.
- Amer. Stat. Assoc., 59, 202-15.[9]Sampford,M.R.(1967):On Sampling without replacement with unequal probabilities of selection. Biometrika, 54, 499-513.
- Sankarnarayanan,K.(1969): An IPPS sampling scheme using Lahiri's methodof selection. Jour. Ind. Soc. Agr. Stat., 21, 2, 59-66
- Sen, A.R. (1953): On the estimate of the variance in sampling with varying probabilities. J. Ind Soc., Agr.Stat., 5, 119-27.
- Yates,F. and Grundy, P. M. (1953): Selection without replacement within size. Jour. Roy. Stat. Soc. B, 15, 243-61.
