

A Family of Estimators of Population Mean with Two Phase Sampling Subject to Auxiliary Information is Attribute

MONIKA SAINI

(Acceptance Date 7th August, 2013)

Abstract

In this paper, a problem of estimating the finite population mean in two phase sampling when information regarding the population proportion possessing certain attributes is considered. Under simple random sampling without replacement scheme, the expressions for mean square error of the proposed estimators have been obtained, up to the first order of approximation. The conditions under which the proposed estimators are more efficient than the mean per unit estimator have also been obtained. The gains in efficiency over the existing ones have been illustrated numerically.

Key words: Simple random sampling, Auxiliary attributes, Percent relative efficiency, Study variable, Two phase sampling.

1. Introduction

The estimation of the population mean is an unrelenting issue in sampling theory and several efforts have been made to improve the precision or accuracy of an estimator of unknown population parameter of interest when study variable y is highly correlated with the auxiliary variable x . There are many situations when auxiliary information is qualitative in nature that is auxiliary information is available in the form of an attribute, which is highly correlated with study variable. For example

- (a) Sex and height of the person.
- (b) Amount of milk produced and a particular

breed of the cow.

- (c) Amount of yield of wheat crop and a particular variety of wheat etc.^{1,2}

In such situations, taking the advantage of point bi-serial correlation between the study variable y and the auxiliary attributes the estimators of population parameter of interest can be constructed by using prior knowledge of the population parameter of auxiliary attribute.

Now consider a finite population which consists of N identifiable units U_i ($1 \leq i \leq N$). Assume that a sample of size n drawn by using

simple random sampling without replacement (SRSWOR) from a population of size N. Let y_i and ϕ_i denote the observations on the variable y and ϕ respectively for i^{th} unit ($i=1,2,\dots,N$). Suppose there is a complete dichotomy in the population with respect to the presence or absence of an attribute, say ϕ , and it is assumed that attribute ϕ takes only two values 0 and 1 according as

$$\begin{aligned} \phi_i &= 1, \text{ if } i^{th} \text{ unit of the population} \\ &\text{possesses attribute } \phi \\ &= 0 \text{ otherwise.} \end{aligned}$$

Let $A = \sum_i^N \phi_i$ and $a = \sum_i^n \phi_i$ denote the total number of units in the population and sample respectively possessing attribute ϕ . Let $P = \frac{A}{N}$ and $p = \frac{a}{n}$ denote the proportion of units in the population and sample respectively possessing attribute ϕ . Let $s_y^2 = \frac{1}{n-1} \sum_i^n (y_i - \bar{y})^2$ and $s_\phi^2 = \frac{1}{n-1} \sum_i^n (\phi_i - p)^2$ be sampling variance corresponding to the population variances $S_y^2 = \frac{1}{N-1} \sum_i^N (y_i - \bar{Y})^2$ and $S_\phi^2 = \frac{1}{N-1} \sum_i^N (\phi_i - P)^2$ respectively, where $\bar{y} = \frac{1}{n} \sum_i^n y_i$ and $\bar{Y} = \frac{1}{N} \sum_i^N y_i$

$$\text{Let } s_{y\phi} = \frac{1}{n-1} \sum_i^n (y_i - \bar{y})(\phi_i - p)$$

and $\hat{\rho}_{y\phi} = \frac{s_{y\phi}}{s_y s_\phi}$ be the sample point bi-serial covariance and bi-serial correlation between variable of interest y and auxiliary attribute ϕ corresponding to the population point bi-serial covariance and bi-serial correlation $S_{y\phi} = \frac{1}{N-1} \sum_i^N (y_i - \bar{Y})(\phi_i - P)$ is the population covariance

$$\text{and } \rho_{y\phi} = \frac{S_{y\phi}}{S_y S_\phi}, \text{ respectively.}$$

It is assumed that population proportion P and population variance S_ϕ^2 for the auxiliary attribute ϕ are unknown. In such situation, we can estimate P and S_ϕ^2 from the sample by using a two phase sampling procedure as per jhaji^{1,2}, Kiregyra⁵ and swain⁶. We use simple random sampling without replacement at both phases as described below:

- (i) We draw a sample s^* of fixed size n^* from the population and observe ϕ and estimate P as well as S_ϕ^2 .
- (ii) Given s^* , we draw a sample s ($s \subseteq s^*$) of fixed size n and observe y.

$$\begin{aligned} \text{Let } p^* &= \frac{1}{n^*} \sum_i^{n^*} \phi_i, s_\phi^2 = \frac{1}{n-1} \sum_i^n (\phi_i - p)^2 \\ \text{and } s_\phi^{*2} &= \frac{1}{n^*-1} \sum_i^{n^*} (\phi_i - p^*)^2. \end{aligned}$$

In the present paper, some unbiased estimators for estimating the population mean of the variable under study, which make use estimated auxiliary information regarding the population proportion certain attribute, are proposed. The expressions for MSE have been obtained. A comparison between all suggested estimators with other known estimators using real data set is considered.

2. Suggested Estimators :

Under the given sampling design, we propose the following estimators of population mean as

$$(i) \hat{Y}_{d1} = \bar{y} + \left(\frac{p^*}{p}\right) - 1 \quad (2.1)$$

To obtain the characteristics of proposed estimators to the first degree of approximation,

we define $\delta_y = \frac{\bar{y}}{\bar{Y}} - 1$, $\delta_1 = \frac{p}{P} - 1$ and

$$\delta_2 = \frac{p^*}{P} - 1$$

So we have $E(\delta_i) = 0$, $i = (y, 1, 2)$

$$\text{and } E(\delta_y^2) = \left(\frac{1-f}{n}\right) \frac{S_y^2}{\bar{Y}^2}, E(\delta_1^2) = \left(\frac{1}{n} - \frac{1}{N}\right) \frac{S_p^2}{P^2},$$

$$E(\delta_2^2) = \left(\frac{1}{n} - \frac{1}{n'}\right) \frac{S_p^2}{P^2} \quad E(\delta_1 \delta_2) = \left(\frac{1}{n} - \frac{1}{n'}\right) \rho_{y\theta} \frac{S_y S_\theta}{\bar{Y} P}$$

Expressing (2.1) in terms of δ 's, we have

$$\hat{Y}_{d1} = \bar{Y} (1 + \delta_y) + (1 + \delta_2) (1 + \delta_1)^{-1} - 1 \quad (2.2)$$

Expanding the right hand side of (2.2) and retaining terms up to first powers of δ 's, we have

$$\hat{Y}_{d1} = \bar{Y} (1 + \delta_y) + (1 + \delta_2)(1 + \delta_1 + \dots) - 1 \quad (2.3)$$

Taking conditional expectation on (2.3), we obtain

$$\hat{Y}_{d1} = \bar{Y}$$

Now to find out the sampling variance, we have to the first order of approximation

$$\text{Var}(\hat{Y}_{d1}) = E(\hat{Y}_{d1} - E(\hat{Y}_{d1}))^2$$

$$\text{Var}(\hat{Y}_{d1}) = \left(\frac{1}{n} - \frac{1}{N}\right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n'}\right) \left\{ \frac{S_\theta^2}{P^2} - 2 \frac{\rho_{y\theta} S_y S_\theta}{P} \right\} \quad (2.4)$$

$$(ii) \hat{Y}_{d2} = \bar{y} - e^{(p-p^*)} + 1 \quad (2.5)$$

Expanding the right hand side of (2.5), it is clear

that $E(\hat{Y}_{d2}) = E_1 E_2(\hat{Y}_{d2}) = \bar{Y}$ so that \hat{Y}_{d2} is conditionally unbiased estimator of the population mean.

Now conditional variance for \hat{Y}_{d2} is

$$\begin{aligned} V(\hat{Y}_{d2}) &= V_1 E_2(\hat{Y}_{d2}) + E_1 V_2(\hat{Y}_{d2}) \\ &= \left(\frac{1}{n} - \frac{1}{N}\right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n'}\right) \{S_\theta^2 + S_\theta^2 - 2\rho_{y\theta} S_y S_\theta\} \\ &\quad \text{(Neglecting the higher terms)} \end{aligned}$$

After simplification the above equation can be written as

$$V(\hat{Y}_{d2}) = \left(\frac{1}{n} - \frac{1}{N}\right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n'}\right) \{S_\theta^2 - 2\rho_{y\theta} S_y S_\theta\} \quad (2.6)$$

$$(iii) \hat{Y}_{d3} = \bar{y} - e^{(p^*-p)} + 1 \quad (2.7)$$

Expanding the right hand side of (2.7), it is clear that $E(\hat{Y}_{d3}) = E_1 E_2(\hat{Y}_{d3}) = \bar{Y}$ so that

\hat{Y}_{d3} is conditionally unbiased estimator of the population mean.

Now conditional variance for \hat{Y}_{d3} is

$$\begin{aligned} V(\hat{Y}_{d3}) &= V_1 E_2(\hat{Y}_{d3}) + E_1 V_2(\hat{Y}_{d3}) \\ &= \left(\frac{1}{n'} - \frac{1}{N}\right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n'}\right) \{S_y^2 + S_\theta^2 + 2\rho_{y\theta} S_y S_\theta\} \\ &\quad \text{(Neglecting the higher terms)} \end{aligned}$$

After simplification the above equation can be written as

$$V(\hat{Y}_{d3}) = \left(\frac{1}{n} - \frac{1}{N}\right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n'}\right) \{S_\theta^2 + 2\rho_{y\theta} S_y S_\theta\} \quad (2.8)$$

3. Efficiency Comparisons :

In this section, the conditions for which the proposed estimators \hat{Y}_{d1} , \hat{Y}_{d2} and \hat{Y}_{d3} are better than the mean per unit.

The variance of is given by

$$\text{Var}(\bar{y}) = \left(\frac{1}{n} - \frac{1}{N}\right) S_y^2 \quad (3.1)$$

Suggested Estimators vs. Mean per unit :

From equation (3.1) and (2.4)

$$V(\bar{y}) - V(\hat{Y}_{d1}) > 0 \text{ if } \rho_{y\phi} > \frac{1}{2P} \frac{S_\phi}{S_y}$$

Thus we arrive at the following theorem:

Theorem 3.1: The estimator (\hat{Y}_{d1})

is more efficient than \bar{y} if $\rho_{y\phi} > \frac{1}{2P}$

Corollary 3.1: If S_ϕ and S_y are approximately the same then (\hat{Y}_{d1}) is more efficient than (\bar{y}) whenever

$$\rho_{y\phi} > \frac{1}{2P}$$

From equation (3.1) and (2.6), we obtain

$$V(\bar{y}) - V(\hat{Y}_{d2}) > 0 \text{ if } \rho_{y\phi} > \frac{1}{2} \frac{S_\phi}{S_y}$$

Thus we arrive at the following theorem:

Theorem 3.2: The estimator (\hat{Y}_{d2})

is more efficient than \bar{y} if $\rho_{y\phi} > \frac{1}{2} \frac{S_\phi}{S_y}$.

Corollary 3.2: If S_ϕ and S_y are approximately the same then (\hat{Y}_{d2}) is more efficient than (\bar{y}) whenever

$$\rho_{y\phi} > \frac{1}{2}$$

From equation (3.1) and (2.8), we obtain

$$V(\bar{y}) - V(\hat{Y}_{d3}) > 0 \text{ if } \rho_{y\phi} < -\frac{1}{2} \frac{S_\phi}{S_y}$$

Thus we arrive at the following theorem:

Theorem 3.3: The estimator (\hat{Y}_{d3})

is more efficient than \bar{y} if.

$$\rho_{y\phi} < -\frac{1}{2}$$

Corollary 3.3: If S_ϕ and S_y are approximately the same then (\hat{Y}_{d3}) is more efficient than (\bar{y}) whenever

$$\rho_{y\phi} < -\frac{1}{2}$$

4. Numerical illustration :

To get a rough idea about the gain in efficiency for the proposed estimators. We compare the performance of various estimators considered here using the two data sets as previously used by Shabbir and Gupta⁴.

Population³ (Source: Sukhatme and Sukhatme(1970), pp. 256).

y = Number of villages in the circles.
 ϕ = A circle consisting more than five villages.
 $N = 89, \bar{Y} = 3.36, P = 0.124, \rho_{y\phi} = 0.766,$
 $n' = 45, n = 23, S_y = 2.019 S_\phi = 0.332,$
 $S_{y\phi} = 0.513$

The percent relative efficiency (PRE's) of the suggested estimators with respect to the mean per unit estimator have been computed and compiled in table 1.

Table 1. Suggested Estimators vs. Mean per unit

Population (Source: Sukhatme and Sukhatme(1970), pp. 256)	Estimator PREs			
	\bar{y}	\hat{Y}_{d1}	\hat{Y}_{d2}	\hat{Y}_{d3}
	100	120	133	104

5. Conclusion

We have developed new estimators and obtained the minimum MSE equations for the proposed estimators. Theoretically, we have demonstrated that all proposed estimators are always more efficient than the mean square error. In addition, we support this theoretical result numerically using the data used by Shabbir and Gupta⁴.

References

1. Jhajj, H.S., Sharma, M.K. and Grover, L.K., An efficient class of chain estimators of population variance under sub-sampling scheme, *Journal of Japan Statistical Society* 35(2), 273-286 (2005).
2. Jhajj, H.S., Sharma, M.K. and Grover, L.K., A family of estimators of population mean using information on auxiliary attribute. *Pakistan journal of Statistics*, Vol. 22(1), pp. 43-50 (2006).
3. P.V. Sukhatme, B.V. Sukhatme, *Sampling Theory of Surveys with Applications*, Iowa State University Press, Ames, USA (1970).
4. Shabbir, J. and Gupta, S., Estimation of the finite population mean in two phase sampling when auxiliary variables are attributes, *Hecettepe Journal of Mathematics and Statistics*, Vol. 39(1), pp. 121-129 (2010).
5. Kiregyra, B., Regression type estimators using two auxiliary variables and the models of double sampling from finite populations, *Metrika*, 31(3-4), 215-226 (1984).
6. Swain, A.K.P.C., A note on the use of multiple auxiliary variables in sample surveys, *Trabajos de Estadística*, 30, 135-141 (1970).