

ESTIMATION OF POPULATION MEAN IN TWO PHASE SAMPLING SCHEME

Peeyush Misra

Abstract

For estimating finite population mean, a double sampling regression type estimator using auxiliary information is proposed. Its bias and mean square error are found and its comparative study with some of the well known estimators is done. An empirical study is also given to show numerically that the proposed estimator is more efficient than the others.

1. INTRODUCTION

Often, statisticians make use of the information available on an auxiliary variable with the variable under study for improving the efficiency of an estimator. For further better understanding one may see Cochran (1977), Des Raj (1968), Murthy (1967), Mukhopadhyay (2012), Singh & Chaudhary (1997) and Sukhatme et, al. (1984). It is well known that the auxiliary information in sample surveys results in substantial improvement in the precision of the estimators of the population parameters and we know that sometimes parameters of the auxiliary variables are not known in advance then double or two phase sampling technique is used. In double sampling or two-phase sampling technique, we first take a preliminary large sample of size n' (called first phase sample) from a population of size N and then a sub-sample of size n (called second phase sample) is drawn from the first phase sample of size n' using simple random sampling without replacement at both the phases. At first phase sample of size n' , only the auxiliary variable X be observed but at the second phase sample of size n , the study variable Y and the auxiliary variable X both are observed.

Let $\bar{Y} = \frac{1}{N} \sum_{i=1}^N Y_i$ be the population mean of study variable y and

$\bar{X} = \frac{1}{N} \sum_{i=1}^N X_i$ be the population mean of auxiliary variable x .

Keywords: Auxiliary Variable, Bias, Mean Squared Error and Efficiency.

$$\sigma_Y^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2, \quad \sigma_X^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2 \quad \text{and}$$

$$\rho = \frac{\frac{1}{N} \sum_{i=1}^N (Y_i - \bar{Y})(X_i - \bar{X})}{\sigma_Y \sigma_X}$$

be the population correlation coefficient between y and x .

$$\text{Also let } \mu_{rs} = \frac{1}{N} \sum_{i=1}^N (Y_i - \bar{Y})^r (X_i - \bar{X})^s, \quad C_Y^2 = \frac{\sigma_Y^2}{\bar{Y}^2}, \quad C_X^2 = \frac{\sigma_X^2}{\bar{X}^2} = \frac{\mu_{02}}{\bar{X}^2}, \quad \rho = \frac{\mu_{11}}{\sigma_Y \sigma_X},$$

$$\beta_2 = \frac{\mu_{04}}{\mu_{02}^2}, \quad \beta_1 = \frac{\mu_{03}}{\mu_{02}^3}, \quad \gamma_1 = \sqrt{\beta_1}.$$

Let the first phase sample of size n' be $(x'_1, x'_2, \dots, x'_{n'})$ on x and the second phase sample of size n be $\{(y_1, x_1), (y_2, x_2), \dots, (y_n, x_n)\}$ on variables (y, x) with the first phase sample mean $\bar{x}' = \frac{1}{n'} \sum_{i=1}^{n'} x'_i$ estimator of population mean \bar{X}

and the second phase sample mean $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ and mean $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ respectively on y and x .

For simplicity, it is assumed that N is large enough as compared to n so that finite population correction terms may be ignored. A new double sampling regression type estimator represented by \hat{y} for estimating the population mean is proposed as

$$\hat{y} = \bar{y} + b(\bar{x}' - \bar{x}) + k(s_x^2 - s_x'^2) \quad (1.1)$$

where, b is an estimate of the change in y when x is increased by unity.

2. BIAS AND MEAN SQUARED ERROR OF THE PROPOSED ESTIMATOR

In order to obtain bias and mean square error of the proposed estimator, let us denote by

$$\begin{aligned}
 \bar{y} &= \bar{Y}(1 + e_0) \\
 \bar{x} &= \bar{X}(1 + e_1) \\
 \bar{x}' &= \bar{X}(1 + e_1') \\
 s_{yx} &= e_2 + S_{YX} \\
 s_x^2 &= e_3 + S_X^2 \\
 s_x'^2 &= e_3' + S_X^2
 \end{aligned}
 \tag{2.1}$$

so that ignoring finite population correction, for simplicity we have

$$E(e_0) = E(e_1) = E(e_1') = E(e_2) = E(e_3) = E(e_3') = 0
 \tag{2.2}$$

$$E(e_0^2) = \frac{\mu_{20}}{n\bar{Y}^2} = \frac{1}{n} C_Y^2$$

$$E(e_1^2) = \frac{\mu_{02}}{n\bar{X}^2} = \frac{1}{n} C_X^2$$

$$E(e_1'^2) = \frac{\mu_{02}}{n'\bar{X}^2} = \frac{1}{n'} C_X^2$$

$$E(e_3^2) = \left(\frac{\beta_2(x) - 1}{n} \right) S_X^4 = \frac{\mu_{02}^2}{n} \left(\frac{\mu_{04}}{\mu_{02}^2} - 1 \right)$$

$$E(e_3'^2) = \left(\frac{\beta_2(x) - 1}{n'} \right) S_X^4 = \frac{\mu_{02}^2}{n'} \left(\frac{\mu_{04}}{\mu_{02}^2} - 1 \right)$$

$$E(e_0 e_1) = \frac{\mu_{11}}{n\bar{Y}\bar{X}} = \frac{1}{n} \rho C_Y C_X$$

$$E(e_0 e_1') = \frac{\mu_{11}}{n'\bar{Y}\bar{X}} = \frac{1}{n'} \rho C_Y C_X$$

$$E(e_0 e_3) = \frac{\mu_{12}}{n\bar{Y}}$$

$$\begin{aligned}
E(e_0 e'_3) &= \frac{\mu_{12}}{n' \bar{Y}} \\
E(e_1 e'_1) &= \frac{\mu_{02}}{n' \bar{X}^2} = \frac{1}{n'} C_X^2 \\
E(e_1 e_2) &= \frac{\mu_{12}}{n \bar{X}} \\
E(e_1 e_3) &= \frac{\mu_{03}}{n \bar{X}} \\
E(e_1 e'_3) &= \frac{\mu_{03}}{n' \bar{X}} \\
E(e'_1 e_2) &= \frac{\mu_{12}}{n' \bar{X}} \\
E(e'_1 e_3) &= \frac{\mu_{03}}{n' \bar{X}} \\
E(e'_1 e'_3) &= \frac{\mu_{03}}{n' \bar{X}} \\
E(e_3 e'_3) &= \left(\frac{\beta_2(x) - 1}{n'} \right) S_X^4 = \frac{\mu_{02}^2}{n'} \left(\frac{\mu_{04}}{\mu_{02}^2} - 1 \right) \tag{2.3}
\end{aligned}$$

The proposed double sampling regression type estimator represented by $\hat{\bar{y}}$ for estimating the population mean given in (1.1) is

$$\hat{\bar{y}} = \bar{y} + b(\bar{x}' - \bar{x}) + k(s_x^2 - s_x'^2) \tag{2.4}$$

In terms of e_i 's, $i=0,1,2,3$; the above proposed double sampling regression type estimator up to terms of order $O(1/n)$ reduces to

$$\begin{aligned}
\hat{\bar{y}} - \bar{Y} &= \bar{Y} e_0 + \beta \bar{X} e'_1 - \beta \bar{X} e_1 + e_3 k - e_3' k + \frac{\beta \bar{X}}{S_X^2} (e_1 e_3 - e_1' e_3) \\
&\quad + \frac{\bar{X}}{S_X^2} (e_1' e_2 - e_1 e_2) \tag{2.5}
\end{aligned}$$

where, $\beta = \frac{S_{YX}}{S_X^2}$.

Taking expectation on both the sides of (2.5), the bias of \hat{y} up to terms of order $O(1/n)$ is given by

$$\text{Bias}(\hat{y}) = \{E(\hat{y}) - \bar{Y}\} = \frac{1}{S_X^2} \left(\frac{1}{n} - \frac{1}{n'} \right) (\beta\mu_{03} - \mu_{12}) \tag{2.6}$$

Now squaring both sides of (2.5) and taking expectation, the mean square error of \hat{y} up to terms of order $O(1/n)$ is given by

$$\begin{aligned} \text{MSE}(\hat{y}) &= \{E(\hat{y}) - \bar{Y}\}^2 \\ &= \bar{Y}^2 E(e_0^2) + \beta^2 \bar{X}^2 E(e_1^2) + \beta^2 \bar{X}^2 E(e_1'^2) + k^2 E(e_3^2) + k^2 E(e_3'^2) \\ &\quad - 2\beta\bar{Y}\bar{X}E(e_0e_1) + 2\beta\bar{Y}\bar{X}E(e_0e_1') + 2\bar{Y}kE(e_0e_3) - 2\bar{Y}kE(e_0e_3') \\ &\quad - 2\beta^2\bar{X}^2E(e_1e_1') - 2\beta\bar{X}kE(e_1e_3) + 2\beta\bar{X}kE(e_1e_3') + 2\beta\bar{X}kE(e_1'e_3) \\ &\quad - 2\beta\bar{X}kE(e_1'e_3') - 2k^2E(e_3e_3') \end{aligned}$$

using values of the expectation given in (2.2) and (2.3), we have

$$\begin{aligned} \text{MSE}(\hat{y}) &= \frac{\mu_{20}}{n} + \beta^2\mu_{02} \left(\frac{1}{n} - \frac{1}{n'} \right) - 2\beta\mu_{11} \left(\frac{1}{n} - \frac{1}{n'} \right) + \left(\frac{1}{n} - \frac{1}{n'} \right) \\ &\quad \left[\{\beta_2(x) - 1\} S_X^4 k^2 - 2(\beta\mu_{03} - \mu_{12})k \right] \end{aligned} \tag{2.7}$$

which attains the minimum for the optimum value

$$k = \frac{(\beta\mu_{03} - \mu_{12})}{\{\beta_2(x) - 1\} S_X^4} \tag{2.8}$$

Substituting the value of k given by (2.8) in (2.7), we get the minimum mean square error of \hat{y} to be

$$\text{MSE}(\hat{y})_{\min} = \frac{\mu_{20}}{n} + \beta^2\mu_{02} \left(\frac{1}{n} - \frac{1}{n'} \right) - 2\beta\mu_{11} \left(\frac{1}{n} - \frac{1}{n'} \right)$$

$$-\left(\frac{1}{n} - \frac{1}{n'}\right) \frac{(\beta\mu_{03} - \mu_{12})^2}{\{\beta_2(x) - 1\} S_X^4} \quad (2.9)$$

3. EFFICIENCY COMPARISON

- (i) **General estimator of mean in case of SRSWOR:** The general estimator of mean in case of SRSWOR is $\hat{y}_{wor} = \bar{y}$ with

$$MSE(\hat{y}) = \frac{\mu_{20}}{n} \quad (3.1)$$

It is clear that the proposed estimator is more efficient than the estimator \hat{y}_{wor} based on simple random sampling when no auxiliary information is used.

- (ii) **Usual double sampling regression estimator:** The usual double sampling regression estimator is $\bar{y}_{ld} = \bar{y} + b(\bar{x}' - \bar{x})$ with

$$MSE(\bar{y}_{ld}) = \frac{\mu_{20}}{n} - \left(\frac{1}{n} - \frac{1}{n'}\right) \frac{\mu_{11}^2}{\mu_{02}} \quad (3.2)$$

It is clear that the proposed estimator is more efficient than the usual double sampling regression estimator where the auxiliary information already is in use.

4. EMPIRICAL STUDY

To illustrate the performance of the proposed estimator, let us consider the following data

Population I: Cochran (1977, Page Number- 181)

y : Paralytic Polio Cases 'placebo' group

x : Paralytic Polio Cases in not inoculated group

$\mu_{02} = 71.8650173$, $\mu_{20} = 9.889273356$, $\mu_{11} = 19.4349481$, $\mu_{12} = 346.3174191$,

$\mu_{03} = 1453.077703$, $\mu_{40} = 424.1846721$, $\mu_{21} = 94.21286383$,

$\mu_{22} = 3029.312542$, $\mu_{30} = 47.34479951$, $\mu_{04} = 46132.5679$, $\bar{y} = 2.588235294$,

$$\bar{x} = 8.370588235, S_x = 8.477323711, S_y = 3.144721507, \rho = 0.729025009,$$

$$\beta_2(y) = 4.337367369, \beta_2(x) = 8.932490454, C_x = 1.012751251,$$

$$C_y = 1.215006037, \beta = 0.270436839, n = 34, n' = 50 \text{ (say).}$$

$$MSE(\hat{\bar{y}}_{wor}) = 0.290860981, MSE(\bar{y}_{ld}) = 0.241393443 \text{ and } MSE(\hat{\bar{y}})_{\min} = 0.240893525.$$

PRE of the proposed estimator $\hat{\bar{y}}$ over $\hat{\bar{y}}_{wor} = 120.7425483$.

PRE of the proposed estimator $\hat{\bar{y}}$ over $\bar{y}_{ld} = 100.2075264$.

Population II: Mukhopadhyay (2012, Page Number - 104)

y : Quality of raw materials (in lakhs of bales)

x : Number of labourers (in thousands)

$$\mu_{02} = 9704.4475, \mu_{20} = 90.95, \mu_{11} = 612.725, \mu_{12} = 93756.3475,$$

$$\mu_{03} = 988621.5173, \mu_{40} = 35456.4125, \mu_{21} = 11087.635, \mu_{22} = 2893630.349,$$

$$\mu_{30} = 1058.55, \mu_{04} = 341222548.2, \bar{y} = 41.5, \bar{x} = 441.95, S_x = 98.51115419,$$

$$S_y = 9.536770942, \rho = 0.652197067, \beta_2(y) = 4.286367314,$$

$$\beta_2(x) = 3.623231573, C_x = 0.22290113, C_y = 0.229801709,$$

$$\beta = 0.063138576, n = 20, n' = 35 \text{ (say).}$$

$$MSE(\hat{\bar{y}}_{wor}) = 4.5475, MSE(\bar{y}_{ld}) = 3.718501766 \text{ and } MSE(\hat{\bar{y}})_{\min} = 3.633327695.$$

PRE of the proposed estimator $\hat{\bar{y}}$ over $\hat{\bar{y}}_{wor} = 125.1607447$.

PRE of the proposed estimator $\hat{\bar{y}}$ over $\bar{y}_{ld} = 102.3442441$.

Population III: Murthy (1967, Page Number - 398)

y : Number of absentees

x : Number of workers

$$\mu_{02} = 1299.318551, \mu_{20} = 42.13412655, \mu_{11} = 154.6041103,$$

$$\mu_{12} = 5086.694392, \mu_{03} = 32025.12931, \mu_{40} = 11608.18508,$$

$$\mu_{21} = 1328.325745, \mu_{22} = 148328.4069, \mu_{30} = 425.9735118,$$

$$\mu_{04} = 4409987.245, \bar{y} = 9.651162791, \bar{x} = 79.46511628, S_x = 36.04606151,$$

$$S_y = 6.491080538, \rho = 0.660763765, \beta_2(y) = 6.53877409,$$

$$\beta_2(x) = 2.612197776, C_x = 0.453608617, C_x = 0.672569791,$$

$$\beta = 0.118988612, n = 43, n' = 50 \text{ (say)}.$$

$$MSE(\hat{\bar{y}}_{wor}) = 0.979863408, MSE(\bar{y}_{ld}) = 0.919969037 \text{ and } MSE(\hat{\bar{y}})_{\min} = 0.918021175.$$

$$\text{PRE of the proposed estimator } \hat{\bar{y}} \text{ over } \hat{\bar{y}}_{wor} = 106.7364713.$$

$$\text{PRE of the proposed estimator } \hat{\bar{y}} \text{ over } \bar{y}_{ld} = 100.2121805.$$

Population IV: Singh and Chaudhary (1997, Page Number - 176)

y : Total number of guava trees

x : Area under guava orchard (in acres)

$$\mu_{02} = 12.50056686, \mu_{20} = 187123.9172, \mu_{11} = 1377.39858,$$

$$\mu_{12} = 4835.465464, \mu_{03} = 37.09863123, \mu_{40} = 1.48935E+11,$$

$$\mu_{21} = 712662.4414, \mu_{22} = 8747904.451, \mu_{30} = 100476814.5,$$

$$\mu_{04} = 540.1635491, \bar{y} = 746.9230769, \bar{x} = 5.661538462, S_x = 3.535614072,$$

$$S_y = 432.5782209, \rho = 0.900596235, \beta_2(y) = 4.253426603,$$

$$\beta_2(x) = 3.456733187, C_x = 0.624497051, C_y = 0.579146949,$$

$$\beta = 110.1868895, n = 13, n' = 30 \text{ (say)}.$$

$MSE(\hat{\bar{y}}_{wor}) = 14394.14747$, $MSE(\bar{y}_{ld}) = 7778.476942$ and $MSE(\hat{\bar{y}})_{min} = 7715.002111$.

PRE of the proposed estimator $\hat{\bar{y}}$ over $\hat{\bar{y}}_{wor} = 186.5734742$.

PRE of the proposed estimator $\hat{\bar{y}}$ over $\bar{y}_{ld} = 100.8227455$.

Population V: Singh and Chaudhary (1997, Page Number: 154-155)

y : Number of milch animals in survey

x : Number of milch animals in census

$\mu_{02} = 431.5847751$, $\mu_{20} = 270.9134948$, $\mu_{11} = 247.3944637$,

$\mu_{12} = 3119.839406$, $\mu_{03} = 5789.778954$, $\mu_{40} = 154027.4827$,

$\mu_{21} = 2422.297374$, $\mu_{22} = 210594.3138$, $\mu_{30} = 2273.46265$,

$\mu_{04} = 508642.4447$, $\bar{y} = 1133.294118$, $\bar{x} = 1140.058824$, $S_x = 20.77461853$,

$S_y = 16.45945002$, $\rho = 0.723505104$, $\beta_2(y) = 2.098635139$,

$\beta_2(x) = 2.730740091$, $C_x = 0.018222409$, $C_y = 0.014523547$,

$\beta = 0.573223334$, $n = 17$, $n' = 30$ (say).

$MSE(\hat{\bar{y}}_{wor}) = 15.93609$, $MSE(\bar{y}_{ld}) = 12.32127$ and $MSE(\hat{\bar{y}})_{min} = 12.31813399$.

PRE of the proposed estimator $\hat{\bar{y}}$ over $\hat{\bar{y}}_{wor} = 129.3709579$.

PRE of the proposed estimator $\hat{\bar{y}}$ over $\bar{y}_{ld} = 100.0254189$.

5. CONCLUSIONS

From (2.9) it is clear that the proposed double sampling estimator is more efficient than the estimator $\hat{\bar{y}}_{wor}$ based on simple random sampling when no auxiliary information is used and is also more efficient than the usual double sampling regression estimator \bar{y}_{ld} of mean where the auxiliary information already is in use.

From (2.8), the mean squared error of the estimator $\hat{\bar{y}}$ is minimized for the optimum value

$$k = \frac{(\beta\mu_{03} - \mu_{12})}{\{\beta_2(x) - 1\}S_X^4} \quad (5.1)$$

The optimum value involving some unknown parameters may not be known in advance for practical purposes; hence the alternative is to replace the unknown parameters of the optimum value by their unbiased estimators giving estimator depending upon estimated optimum value.

Acknowledgement

The author is thankful to the referees and the editor in chief for their valuable suggestions over the earlier draft of the paper.

REFERENCES

- [1] Cochran, W.G. (1977): Sampling Techniques, 3rd edition, John Wiley and Sons, New York.
- [2] Des Raj (1968): Sampling Theory, McGraw- Hill, New York.
- [3] Murthy, M (1967): Sampling Theory and Methods, 1st edition, Calcutta Statistical Publishing Society, Kolkata, India.
- [4] Mukhopadhyay, Parimal (2012): Theory and Methods of Survey and Sampling, 2nd edition, PHI Learning Private Limited, New Delhi, India.
- [5] Singh, Daroga and Chaudhary, F. S. (1997): Theory and Analysis of Sampling Survey Designs, New Age International Publishers, New Delhi, India.
- [6] Sukhatme, P. V., Sukhatme, B. V., Sukhatme, S. And Asok, C. (1984): Sampling Theory of Surveys with Applications, 3rd Edition, Ames, Iowa (USA) and Indian Society of Agricultural Statistics, New Delhi, India.

Peeyush Misra

Department of Statistics,

D.A.V.(P.G.) College,

Dehradun- 248001,

Uttarakhand (India)

Email: dr.pmisra.dav@gmail.com