

Estimation of Population Mean in Two Phase Sampling using Attribute Auxiliary Information

Nadeem Shafique Butt¹, Muhammad Qaiser Shahbaz²

¹ College of Statistical and Actuarial Sciences, University of the Punjab, Pakistan
nadeemshafique@hotmail.com

² Department of Mathematics, COMSATS Institute of Information Technology, Pakistan
qshahbaz@gmail.com

Abstract. A new estimator for population mean has been proposed in two phase sampling by using information of multiple auxiliary attributes. The minimum variance of the proposed estimator has been obtained.

1. Introduction

The auxiliary information has always been a source of improvement in estimation of certain population characteristics. Several estimators have been developed in single and two phase sampling which utilizes information on auxiliary variables as well as auxiliary attributes. The classical estimators which use information on auxiliary variables are the ratio and regression estimators as given in Hansen, Hurwitz, & Madow (1953). The classical regression estimator of population mean is given as:

$$\bar{y}_{lr} = \bar{y} + \beta(\bar{X} - \bar{x}) \quad (1.1)$$

The value of β for which the mean square error of (1.1) is minimum is $\beta = \frac{S_{xy}}{S_x^2}$. The minimum mean square error of (1.1) is given as

$$MSE(\bar{y}_{lr}) = \theta S_y^2 (1 - \rho_{yx}^2) \quad (1.2)$$

Where $\theta = n^{-1} - N^{-1}$ and ρ_{yx} is the correlation coefficient between X and Y. The estimator (1.1) in case of auxiliary attribute is discussed by Naik & Gupta (1996), and the estimator in this case is given as:

$$t_{1(1)} = \bar{y} + b(p_1 - P_1) \quad (1.3)$$

where P_1 is sample proportion for auxiliary variables. The mean square error of (1.3) is:

$$MSE(t_{1(1)}) = \theta \left(1 - \rho_{pb_1}^2\right) S_y^2 \quad (1.4)$$

where $\rho_{pb_1}^2$ is the squared point bi-serial correlation coefficient. Jhajj, Sharma, & Grover (2006) has proposed a family of estimators in single and two phase sampling using information on a single auxiliary attributes. The proposed family is based upon a general function and is given as:

$$t_{2(1)} = g_{\omega}(\bar{y}, v_1) \tag{1.5}$$

where $v_1 = \frac{p_1}{P_1}$ and $g_{\omega}(\bar{y}, v_1)$ is a parametric function of \bar{y} and v_1 such that $g_{\omega}(\bar{Y}, 1) = \bar{Y}$, for all \bar{Y} .

The mean square error of each estimator; to the terms of order $1/n$; of this family is,

$$MSE(T_{2(1)}) \approx \theta(1 - \rho_{pb_1}^2) S_y^2. \tag{1.6}$$

The mean square error of the proposed family depends upon the point bi-serial correlation coefficient.

Shabbir & Gupta (2007) have also proposed an estimator for population mean in single phase sampling using information of single auxiliary attribute. The estimator is given as:

$$t_{3(1)} = [d_1 \bar{y} + d_2 (P_1 - p_1)] \frac{P_1}{p_1}, \text{ for } p_1 > 0 \tag{1.7}$$

where d_1 and d_2 are unknown constants. The mean square error of (1.7) is:

$$MSE(t_{3(1)}) \approx \frac{\theta(1 - \rho_{pb_1}^2) S_y^2}{1 + \theta(1 - \rho_{pb_1}^2) C_y^2} \tag{1.8}$$

In this paper we have proposed a modified regression type estimator using information on several auxiliary attributes.

2. Notations

In this section we define the notations used for the development of the estimator and its variance. Let δ be a vector of auxiliary attributes with covariance matrix S_{δ} , τ be another auxiliary attribute and Y be the variable of interest.

Let $s_{\tau\delta}$ be the vector of covariances between τ and δ , $s_{y\delta}$ be the vector of covariances between Y and δ . Using these notations we define $\gamma = S_{\delta}^{-1} s_{\tau\delta}$ as

vector of regression coefficients between τ and δ , and $\gamma = \mathbf{S}_{\delta}^{-1} \mathbf{s}_{y\delta}$ as vector of regression coefficients between Y and δ . We also define $\beta_{y\tau,\delta} = S_{\tau y,\delta} / S_{\tau,\delta}^2$ as partial regression coefficient between Y and τ keeping the δ at constant level. Also $S_{y\tau,\delta} = S_{y\tau} - \mathbf{s}'_{\tau\delta} \mathbf{S}_{\delta}^{-1} \mathbf{s}_{\tau\delta}$ is partial covariance between Y and τ after removing the effect of δ , $\mathbf{S}_{y,\delta}^2 = S_y^2 - \mathbf{s}'_{\tau\delta} \mathbf{S}_{\delta}^{-1} \mathbf{s}_{\tau\delta}$ is the partial variance of Y , and $S_{\tau,\delta}^2 = S_{\tau}^2 - \mathbf{s}'_{\tau\delta} \mathbf{S}_{\delta}^{-1} \mathbf{s}_{\tau\delta}$ is the partial variance of τ . We also define $\rho_{y\tau,\delta}^2 = S_{y\tau,\delta}^2 / (S_{\tau,\delta}^2 S_{y,\delta}^2)$ as partial correlation coefficient between Y and τ after removing the effect of δ , $\rho_{y\tau,\delta}^2$ as squared multiple bi-serial correlation coefficient between Y and combined effects of τ and δ , $\rho_{y,\delta}^2$ as squared multiple correlation coefficient between Y and combined effects of δ .

Using the above notations we proposed the new estimators in the section 3.

3. The Proposed Estimator

We propose following unbiased estimator of population mean in two phase sampling using information of several auxiliary attributes:

$$t_{nss(A)} = \bar{y}_2 + k \left[p_{\tau_1} + \gamma' (\mathbf{p}_{\delta} - \mathbf{p}_{\delta_1}) - \left\{ p_{\tau_2} + \eta' (\mathbf{p}_{\delta} - \mathbf{p}_{\delta_2}) \right\} \right] \quad (3.1)$$

Using

$\bar{y}_2 = \bar{Y} + \bar{e}_{y_2}$, $p_{\tau_1} = p_{\tau} + \bar{e}_{\tau_1}$, $p_{\tau_2} = p_{\tau} + \bar{e}_{\tau_2}$, $\mathbf{p}_{\delta_1} = \mathbf{p}_{\delta} + \bar{\mathbf{e}}_{\delta_1}$ and $\mathbf{p}_{\delta_2} = \mathbf{p}_{\delta} + \bar{\mathbf{e}}_{\delta_2}$ in (3.1) we have:

$$t_{nss(A)} - \bar{Y} = \bar{e}_{y_2} + k \left[(\bar{e}_{\tau_1} - \bar{e}_{\tau_2}) - \gamma' \bar{\mathbf{e}}_{\delta_1} + \eta' \bar{\mathbf{e}}_{\delta_2} \right]$$

Squaring and applying expectation, the mean square error of (3.1) is given as:

$$S = MSE(t_{nss}) = \theta_2 s_y^2 + k^2 \left[(\theta_2 - \theta_1) s_{\tau}^2 + \theta_1 \gamma' \mathbf{S}_{\delta} \gamma + \theta_2 \eta' \mathbf{S}_{\delta} \eta + 2(\theta_1 - \theta_2) \eta' \mathbf{s}_{\tau\delta} - 2\theta_1 \gamma' \mathbf{S}_{\delta} \eta \right] + 2k \left[(\theta_1 - \theta_2) s_{y\tau} - \theta_1 \gamma' \mathbf{s}_{y\delta} + \theta_2 \eta' \mathbf{s}_{y\delta} \right] \quad (3.2)$$

The optimum values of γ , η and k are obtained by minimizing (3.2). These values are obtained by solving following three equations, obtained by partially differentiating (3.2) and setting the derivative to zero

$$2k \left[(\theta_2 - \theta_1) s_\tau^2 + \theta_1 \boldsymbol{\gamma}' \mathbf{S}_\delta \boldsymbol{\gamma} + \theta_2 \boldsymbol{\eta}' \mathbf{S}_\delta \boldsymbol{\eta} + 2(\theta_1 - \theta_2) \boldsymbol{\eta}' \mathbf{s}_{\tau\delta} - 2\theta_1 \boldsymbol{\gamma}' \mathbf{S}_\delta \boldsymbol{\eta} \right] + 2 \left[(\theta_1 - \theta_2) s_{y\tau} - \theta_1 \boldsymbol{\gamma}' \mathbf{s}_{y\delta} + \theta_2 \boldsymbol{\eta}' \mathbf{s}_{y\delta} \right] = 0 \quad (1)$$

$$k \mathbf{S}_\delta (\boldsymbol{\gamma} - \boldsymbol{\eta}) - \mathbf{s}_{y\delta} = 0 \quad (2)$$

$$k \mathbf{S}_\delta (\theta_2 \boldsymbol{\eta} - \theta_1 \boldsymbol{\gamma}) - k (\theta_2 - \theta_1) \mathbf{s}_{\tau\delta} + \theta_2 \mathbf{s}_{y\delta} = \mathbf{0} \quad (3)$$

Solving the above equations simultaneously, the optimum values of $\boldsymbol{\alpha}$, $\boldsymbol{\beta}$ and k are:

$$\boldsymbol{\gamma} = \mathbf{S}_\delta^{-1} \mathbf{s}_{x\delta}, \quad \boldsymbol{\eta} = \boldsymbol{\gamma} - k^{-1} \boldsymbol{\gamma} \quad \text{and} \quad k = \boldsymbol{\eta}'_{y\tau\delta} = \frac{s_{y\tau\delta}}{s_{\tau\tau\delta}}$$

Using the optimum values in (3.2) and simplifying, the mean square error of proposed estimator is:

$$MSE(t_{nss}) = s_{y\delta}^2 \left[\theta_2 (1 - \rho_{\tau y\delta}^2) + \theta_1 \rho_{\tau y\delta}^2 \right] \quad (3.3)$$

Further, by using the fact that $s_{y\delta}^2 = S_y^2 (1 - \rho_{\tau y\delta}^2)$ and utilizing the relationship that $1 - \rho_{y\tau\delta}^2 = (1 - \rho_{y\delta}^2)(1 - \rho_{\tau y\delta}^2)$, the mean square error of proposed estimator can be written as:

$$MSE(t_{nss}) = s_y^2 \left\{ \theta_2 (1 - \rho_{y\delta}^2) + \theta_1 \rho_{\tau y\delta}^2 (1 - \rho_{y\delta}^2) \right\} \quad (3.4)$$

From (3.4) we can see that the mean square error of (3.1) depends upon the squared multiple and partial correlation coefficients. The estimator and its mean square error for multiphase sampling can be analogously written from (3.1) and (3.4). Specifically if a sample of size n_h is taken at h^{th} phase and a sample of n_q is taken at q^{th} phase with $n_q < n_h$, the estimator of the population mean is:

$$t_{nss(A)} = \bar{y}_2 + k \left[p_{\tau_h} + \boldsymbol{\gamma}' (\mathbf{p}_\delta - \mathbf{p}_{\delta_h}) - \left\{ p_{\tau_q} + \boldsymbol{\eta}' (\mathbf{p}_\delta - \mathbf{p}_{\delta_q}) \right\} \right] \quad (3.5)$$

The mean square error of (3.5) can be written from (3.4) as:

$$MSE(t_{nss}) = s_y^2 \left\{ \theta_q (1 - \rho_{y\delta}^2) + \theta_h \rho_{\tau y\delta}^2 (1 - \rho_{y\delta}^2) \right\} \quad (3.6)$$

For practical applicability, the proposed estimator can be easily modified by using the sample estimates in place of population parameters. The consistent estimate of population mean can be straight-away written as:

$$t_{nss} = \bar{y}_2 + b_{y\tau\delta} (p_{\tau_1} - p_{\tau_2}) + b_{y\tau\delta} \mathbf{b}'_{\tau\delta} (\mathbf{p}_{\delta_2} - \mathbf{p}_{\delta_1}) + \mathbf{b}'_{y\delta} (\mathbf{p}_\delta - \mathbf{p}_{\delta_2}) \quad (3.7)$$

The estimated standard error of (3.1) is given as:

$$S.E(t_{nss}) = s_y \sqrt{\theta_2 (1 - r_{y.\tau.\delta}^2) + \theta_1 r_{y\tau.\delta}^2 (1 - r_{y.\delta}^2)}$$

(3.8)

Using (3.7) and (3.8), the confidence interval for true population mean can be constructed.

References:

1. Hansen, M. H., Hurwitz, W. N., & Madow, W. G. (1953). *Sample Survey Methods and Theory* (Vol. II): John Wiley.
2. Jhajj, H. S., Sharma, M. K., & Grover, L. K. (2006). A family of estimators of population mean using information of auxiliary attribute. *Pak. J. Stat.*, 22(1), 43-50.
3. Naik, V., & Gupta, P. (1996). A note on estimation of mean with known population proportion of an auxiliary character. *Jour. Ind. Soc. Agr. Stat*, 48(2), 151-158.
4. Shabbir, J., & Gupta, S. (2007). On estimating the finite population mean with known population proportion of an auxiliary variable. *Pak. J. Statist.*, 23(1), 1-9.