

## A Chain Regression Estimator in Two Phase Sampling Using Multi-auxiliary Information

B.K. PRADHAN

Department of Statistics, Utkal University, Bhubaneswar-751004, India  
bkpradhanuu@rediffmail.com

**Abstract.** In this paper the population mean of the study variable  $y$  is estimated in a two phase sampling setup using three auxiliary variables with chain regression concept when the population mean of one of the auxiliary variables is unknown and other auxiliary population mean are known.

2000 Mathematics Subject Classification: 62D05

Key words and phrases: Chain regression, multi-auxiliary information, two phase sampling.

### 1. Introduction

Information on variables correlated with the main variable under study is popularly known as auxiliary information which may be fruitfully utilised either at planning stage or at design stage or at the information stage to arrive at improved estimator compared to those, not utilising auxiliary information. Use of auxiliary information for forming ratio and regression method of estimation were introduced during the 1930's with a comprehensive theory provided by Cochran [1]. Assuming knowledge of multi-auxiliary variables, multivariate ratio estimator was suggested by Olkin [5], multivariate difference estimator by Raj [7], multiple regression estimator by Shukla [11], weighted regression estimator by Srivastava [13, 14, 15] and Ratio-cum-product estimator by Singh [12]. Extension of these estimators to different sampling designs were taken up by Tripathy [17, 18]. Further contribution are due to Rao and Mudholkar [8], Wright [19] and many others.

When information on any auxiliary variable  $x$  highly correlated with  $y$  is readily available on all units of the population, it is well known that ratio and regression estimators provide more efficient estimates of population mean of  $y$ , envisaging advance information on population mean  $\bar{X}$  of  $x$ . However, in certain practical situation when  $\bar{X}$  is unknown, information on auxiliary variables  $Z$  and  $W$  are readily available on all the units of the population, which may also be incorporated in the method of estimation.

## 2. Two phase sampling set up

Consider a finite population  $U$  of size  $N$  indexed by quadruplet characters  $(y, x, z, w)$ . Our purpose is to estimate the population mean  $\bar{Y}$  of a study variable  $y$  in the presence of three auxiliary variables  $x, z$  and  $w$ , when the population mean  $\bar{X}$  of  $x$  is unknown but information on  $z$  and  $w$  are available on all the units of the population.

Let us now consider a two phase sampling where in the first phase a large sample  $S'(S' \subset U)$  of fixed size  $n'$  is drawn following SRSWOR and observe three auxiliary variables  $x, z$  and  $w$  to estimate  $\bar{X}$ , while in the second phase a sub-sample  $S \subset S'$  of fixed size  $n$  is drawn by SRSWOR to observe the characteristic  $y$  under study.

## 3. Use of one auxiliary variable

The two phase regression estimators in this case will be

$$(3.1) \quad \bar{t}_{1(\text{Reg})} = \bar{y}_n + b_{yx}(\bar{x}_{n'} - \bar{x}_n)$$

where  $b_{yx}$  is the sample regression coefficient of  $y$  on  $x$  calculated from data based on  $S$  and

$$\bar{x}_{n'} = \frac{1}{n'} \sum_{i \in S'} x_i \quad \bar{x}_n = \frac{1}{n} \sum_{i \in S} x_i \quad \text{and} \quad \bar{y}_n = \frac{1}{n} \sum_{i \in S} y_i.$$

The mean square error (MSE) of  $\bar{t}_{1(\text{Reg})}$  by first order approximation is

$$(3.2) \quad MSE(\bar{t}_{1(\text{Reg})}) = \left( \frac{1}{n'} - \frac{1}{N} \right) (1 - \rho_{yx}^2) S_y^2 + \left( \frac{1}{n'} - \frac{1}{N} \right) \rho_{yx}^2 S_y^2$$

where

$$S_y^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2$$

and  $\rho_{yx}$  is the correlation coefficient between  $y$  and  $x$ .

## 4. Use of second auxiliary variable

Swain [16], Kiregyera [2], Mukherjee *et al.* [4], Sahoo *et al.* [9] and Mishra *et al.* [3] used a second auxiliary variable  $z$  closely related to  $x$  to suggest different improved estimators assuming that the information on  $z$  is available on all the units of the population.

Kiregyera [2] has suggested a regression type estimator

$$(4.1) \quad \bar{t}_{2(\text{Reg})} = \bar{y}_n + b_{yx}[\bar{x}_{n'} + b_{xz}(\bar{Z} - \bar{z}_{n'}) - \bar{x}_n]$$

a technique earlier suggested by Swain [16] which yields

$$(4.2) \quad MSE(\bar{t}_{2(\text{Reg})}) = \left( \frac{1}{n'} - \frac{1}{N} \right) (1 - \rho_{yx}^2) S_y^2 + (\rho_{yx}^2 + \rho_{yx}^2 \rho_{xz}^2 - 2\rho_{yx} \rho_{yz} \rho_{xz}) S_y^2.$$

Sahoo *et al.* [9] considered a chain regression type estimator

$$(4.3) \quad \bar{t}_{3(\text{Reg})} = \bar{y}_n + b_{yx}(\bar{x}_{n'} - \bar{x}_n) + b_{yz}(\bar{Z} - \bar{z}_{n'}).$$

The MSE of  $\bar{t}_{3(\text{Reg})}$  to the first order of approximation is

$$(4.4) \quad MSE(\bar{t}_{3(\text{Reg})}) = \left(\frac{1}{n} - \frac{1}{N}\right) (1 - \rho_{yx}^2) S_y^2 + \left(\frac{1}{n'} - \frac{1}{N}\right) (\rho_{yx}^2 - \rho_{yz}^2).$$

Another regression type estimator suggested by Mukherjee *et al.* [4] is

$$(4.5) \quad \bar{t}_{4(\text{Reg})} = \bar{y}_n + b_{yx.z} [\bar{x}_{n'} + b_{xz}(\bar{Z} - \bar{z}_{n'}) - \bar{x}_n] + b_{yz.x}(\bar{Z} - \bar{z}_n).$$

Using a generalized method, a regression type estimator suggested by Mishra and Rout [3] is

$$(4.6) \quad \bar{t}_{5(\text{Reg})} = \bar{y}_n + b_{yx.z}(\bar{x}_{n'} - \bar{x}_n) + b_{yz.x}(\bar{Z} - \bar{z}_n) + (b_{yz} - b_{yz.x})(\bar{Z} - \bar{z}_{n'}).$$

In fact, on simplification  $\bar{t}_{4(\text{Reg})} = \bar{t}_{5(\text{Reg})}$ . It may be seen that

$$(4.7) \quad MSE(\bar{t}_{4(\text{Reg})}) = MSE(\bar{t}_{5(\text{Reg})}) = \left(\frac{1}{n} - \frac{1}{N}\right) (1 - \rho_{y.xz}^2) S_y^2 + \left(\frac{1}{n'} - \frac{1}{N}\right) (1 - \rho_{yz}^2) S_y^2.$$

### 5. Suggested estimators

Since  $\widehat{X}$  based on  $n'$  unit is an unbiased estimator of  $\bar{X}$ , a regression type estimator

$$(5.1) \quad \widehat{X} = \bar{x}_{n'} + \beta_{xz.w}(\bar{Z} - \bar{z}_{n'}) + \beta_{xw.z}(\bar{W} - \bar{w}_{n'})$$

is found by considering the estimator

$$(5.2) \quad \widehat{X} = \lambda_1 \bar{x}_{n'} + \lambda_2 \bar{z}_1 + \lambda_3 \bar{Z} + \lambda_4 \bar{w}_{n'} + \lambda_5 \bar{W}$$

which is unbiased for  $\bar{X}$  by considering  $\lambda_1 = 1$ ,  $\lambda_2 = -\lambda_3$  and  $\lambda_4 = -\lambda_5$  and then by minimizing the variance  $\widehat{X}$  given by

$$(5.3) \quad \widehat{X} = \bar{x}_{n'} + \lambda_2(\bar{z}_{n'} - \bar{Z}) + \lambda_4(\bar{w}_{n'} - \bar{W}).$$

Here

$$\bar{x}_{n'} = \frac{1}{n'} \sum_{i \in S'} x_i, \quad \bar{z}_{n'} = \frac{1}{n'} \sum_{i \in S'} z_i, \quad \bar{w}_{n'} = \frac{1}{n'} \sum_{i \in S'} w_i$$

and  $\beta_{xz.w}$  and  $\beta_{xw.z}$  are usual partial regression coefficients.

Let us consider a chain regression type estimator of  $\bar{Y}$  given by

$$(5.4) \quad \bar{t}^* = \bar{y}_n + \lambda_1^*(\widehat{X} - \bar{x}_n) + \lambda_2^*(\bar{Z} - \bar{z}_n) + \lambda_3^*(\bar{W} - \bar{w}_n).$$

where  $\bar{y}_n$ ,  $\bar{x}_n$ ,  $\bar{z}_n$  and  $\bar{w}_n$  are the sample mean based on  $n$  observations of the second phase and  $\lambda_1^*$ ,  $\lambda_2^*$  and  $\lambda_3^*$  are suitable constants.

The optimum values of  $\lambda_1^*$ ,  $\lambda_2^*$  and  $\lambda_3^*$  are obtained by minimising  $V(\bar{t}^*)$  and we find

$$(5.5) \quad \lambda_1^* = \beta_{yx.zw}, \quad \lambda_2^* = \beta_{yz.xw} \text{ and } \lambda_3^* = \beta_{yw.xz}$$

where  $\beta_{yx.zw}$ ,  $\beta_{yz.xw}$  and  $\beta_{yw.xz}$  are usual partial regression coefficients. When the partial regression coefficients are known,  $\bar{t}^*$  is an unbiased estimator of  $\bar{Y}$  with

$$(5.6) \quad V(\bar{t}^*) = \left(\frac{1}{n} - \frac{1}{N}\right) (1 - \rho_{y.xzw}^2) S_y^2 + \left(\frac{1}{n'} - \frac{1}{N}\right) (1 - \rho_{y.zw}^2) \rho_{yx.zw}^2 S_y^2$$

where  $\rho_{y.xzw}$  and  $\rho_{y.zw}$  are usual multiple correlation coefficients and  $\rho_{yx.zw}$  is the usual partial correlation coefficient.

The estimators under consideration require advance knowledge of the population regression coefficients and partial regression coefficients, which are usually unknown. However, in practice the consistent estimators  $b_{yx.zw}$ ,  $b_{yz.xw}$  and  $b_{yw.xz}$  of the population parameters  $\beta_{yx.zw}$ ,  $\beta_{yz.xw}$  and  $\beta_{yw.xz}$  may be substituted for the purpose. Although the estimators will turn out to be biased, this bias would be negligible in large samples and the approximate mean square errors to  $O(1/n)$  will be equivalent to those derived and for large sample, the difference would be minimal.

## 6. Comparison of efficiency

Sahoo *et al.* [9] has established that  $\bar{t}_{1(\text{Reg})}$  and  $\bar{t}_{2(\text{Reg})}$  are less efficient than  $\bar{t}_{3(\text{Reg})}$ . Mishra and Rout [3] has proved that

$$(6.1) \quad MSE(\bar{t}_{5(\text{Reg})}) < MSE(\bar{t}_{3(\text{Reg})}).$$

Now, from (4.7) and (5.6), we find

$$(6.2) \quad MSE(\bar{t}_{5(\text{Reg})}) - V(\bar{t}^*) = \left[ \left( \frac{1}{n} - \frac{1}{N} \right) A + \left( \frac{1}{n'} - \frac{1}{N} \right) B \right] S_y^2$$

where  $A = \rho_{y.xzw}^2 - \rho_{y.xz}^2$  and  $B = (1 - \rho_{yz}^2)\rho_{yx.z}^2 - (1 - \rho_{y.zw}^2)\rho_{yx.zw}^2$ . On simplification, we find

$$(6.3) \quad A + B = (1 - \rho_{yz}^2)\rho_{yw.z}^2 \geq 0.$$

Since

$$\left( \frac{1}{n} - \frac{1}{N} \right) > \left( \frac{1}{n'} - \frac{1}{N} \right),$$

we have from (6.3)

$$(6.4) \quad \left( \frac{1}{n} - \frac{1}{N} \right) A + \left( \frac{1}{n'} - \frac{1}{N} \right) B \geq \left( \frac{1}{n'} - \frac{1}{N} \right) (A + B) \geq 0.$$

Hence

$$(6.5) \quad MSE(\bar{t}_{5(\text{Reg})}) \geq V(\bar{t}^*).$$

The inequality (6.5) shows that  $\bar{t}^*$  is an improved regression estimator compared to  $\bar{t}_{1(\text{Reg})}$ ,  $\bar{t}_{2(\text{Reg})}$ ,  $\bar{t}_{3(\text{Reg})}$  and  $\bar{t}_{5(\text{Reg})}$ .

## 7. Numerical illustration

Percent relative efficiency of different estimators compared to mean per unit estimator are presented in Table 2.

Table 1. Description of Population

	Population I	Population II
Source	“Spray congealing: Particle size relationships using a centrifugal wheel automizer” by Scott, Robinson, Pauls and Lantz (1964)	“Measurement of four characters of: Flucus Religiosament” by Pradhan, B.K. (2000)
y	Mean surface-volume particle size of product	Length of petiole
x	Feed rate per unit whetted wheel periphery (gm/sec/cm)	Length of pamina(blade) of the leaf
z	Peripheral wheel velocity(cm/sec)	Width of the leaf at its widest paint
w	Feed Viscosity (poise)	Width of leaf half way along the blade
size	N=35	N=160
$\rho_{yx}$	0.712296	0.5423
$\rho_{yz}$	-0.8070192	0.6166
$\rho_{yw}$	-0.1623959	0.2704
$\rho_{xz}$	-0.2633457	0.8568
$\rho_{xw}$	-0.0781118	0.7424
$\rho_{zw}$	0.1335984	0.8027

Table 2. Relative efficiency of different estimators of population variance with respect to  $S_y^2$  under comparison

Estimator	Auxiliary variables used	Percent Relative Efficiency of Population I: $N = 35$ , $n' = 12$ , $n = 8$	Percent Relative Efficiency of Population II: $N = 160$ , $n' = 50$ , $n = 20$
$\bar{y}_n$	None	100	100
$\bar{t}_{1(\text{Reg})}$	$X$	128.08	125.26
$\bar{t}_{2(\text{Reg})}$	$X, Z$	159.03	145.75
$\bar{t}_{3(\text{Reg})}$	$X, Z$	243.36	147.31
$\bar{t}_{4(\text{Reg})} = \bar{t}_{5(\text{Reg})}$	$X, Z$	378.98	161.47
$\bar{t}_{(\text{Reg})}^*$	$X, Z, W$	436.27	204.28

**Remark 7.1.**  $\bar{t}_{(\text{Reg})}^*$  has substantial gain in efficiency compared to  $\bar{t}_{5(\text{Reg})}$ ,  $\bar{t}_{4(\text{Reg})}$ ,  $\bar{t}_{3(\text{Reg})}$ ,  $\bar{t}_{2(\text{Reg})}$ ,  $\bar{t}_{1(\text{Reg})}$  and  $\bar{y}_n$  for the population under consideration.

**Acknowledgment.** The author wishes to express his sincere gratitude to the referee for his valuable suggestions in improving the manuscript.

## References

- [1] W. G. Cochran, Sampling theory when the sampling-units are of unequal sizes, *J. Amer. Statist. Assoc.* **37** (1942), 199–212.
- [2] B. Kiregyvera, Regression-type estimators using two auxiliary variables and the model of double sampling from finite populations, *Metrika* **31**(3–4) (1984), 215–226.
- [3] G. Mishra and K. Rout, A regression estimator in two-phase sampling in presence of two auxiliary variables, *Metron* **55**(1–2) (1997), 177–186.
- [4] R. Mukerjee, T. J. Rao and K. Vijayan, Regression type estimators using multiple auxiliary information, *Austral. J. Statist.* **29**(3) (1987), 244–254.
- [5] I. Olkin, Multivariate ratio estimation for finite populations, *Biometrika* **45** (1958), 154–165.
- [6] B. K. Pradhan, Some problems of estimation in multi-phase sampling, Ph.D. Dissertation (unpublished), Utkal University, Orissa, India (2000).
- [7] Des Raj, On a method of using multi-auxiliary information in sample surveys, *J. Amer. Statist. Assoc.* **60** (1965), 270–277.
- [8] P. S. R. S. Rao and G. S. Mudholkar, Generalized multivariate estimator for the mean of finite populations, *J. Amer. Statist. Assoc.* **62** (1967), 1009–1012.
- [9] J. Sahoo, L. N. Sahoo, and S. Mohanty, A regression approach to estimation in two phase sampling using two auxiliary variables, *Current Science* **65**(1) (1993), 73–75.
- [10] M. W. Scott, M. J. Rabinson, J. F. Pauls, and R. J. Lantz, Spray congealing: Particle size relationships using a centrifugal wheel atomizer, *J. Pharmaceutical Sci.*, **53**(6) (1964), 670–675.
- [11] G. K. Shukla, Multivariate regression estimate, *J. Indian Statist. Assoc.* **3** (1965), 202–211.
- [12] M. P. Singh, Ratio cum product method of estimation, *Metrika* **12** (1967), 34–42.
- [13] S. K. Srivastava, An estimate of the mean of a finite population using several auxiliary variables, *J. Indian Statist. Assoc.* **3** (1965), 189–194.
- [14] S. K. Srivastava, An estimator using auxiliary information in sample surveys, *Calcutta Statist. Assoc. Bull.* **16** (1967), 121–132.
- [15] S. K. Srivastava, A generalized estimator for the mean of a finite population using multi-auxiliary information, *J. Amer. Statist. Assoc.* **66** (1971), 404–407.
- [16] A. K. P. C. Swain, A note on the use of multiple auxiliary variables in sample surveys, *Trabajos de Estadística*, **21** (1970), 135–141.
- [17] T. P. Tripathy, Contribution to sampling theory using multivariate information, Ph.D. Thesis (unpublished), Punjab University, Patiala, India, (1970).
- [18] T. P. Tripathy, On double sampling for multivariate ratio and difference methods of estimation, *J. Indian Soc. Ag. Stat.* **28** (1976), 35–54.
- [19] R. L. Wright, Finite population sampling with multivariate auxiliary information, *J. Amer. Statist. Assoc.* **78**(384) (1983), 879–884.