

A family of estimators of population mean using multi-auxiliary variate and post-stratification

Gajendra K. Vishwakarma¹, Housila P. Singh¹, Sarjinder Singh²

¹School of Studies in Statistics, Vikram University
Ujjain-456010, M.P., India
vishwagk@rediffmail.com; hpsujn@rediffmail.com

²Department of Mathematics, Texas A&M University-Kingsville
TX 78363, USA
sarjinder@yahoo.com

Received: 2009-09-15 **Revised:** 2010-04-07 **Published online:** 2010-06-01

Abstract. This paper suggests a family of estimators of population mean using multi-auxiliary variate based on post-stratified sampling and its properties are studied under large sample approximation. Asymptotically optimum estimator in the class is identified alongwith its approximate variance formulae. The proposed class of estimators is also compared with corresponding unstratified class of estimators based on estimated optimum value. At the end, an empirical study has been carried out to support the proposed methodology.

Keywords: multi-auxiliary variate, study variate, bias, variance, post-stratified sampling.

1 Introduction

Stratification is one of the most widely used techniques in sample survey design serving the dual purpose of providing samples that are representative of major sub-groups of the population and improving the precision of estimators [1]. Stratified sampling presupposes the knowledge of strata size as well as the availability of a frame for drawing a sample in each stratum [2]. However application of this technique presupposes the knowledge of strata size and the availability of sampling frames within strata. In many socio-economic and agricultural surveys where it is necessary to partition the finite population under consideration, due to its heterogeneity, into different sub-populations (strata), the sampling frame within strata may not be available. However frame for entire population may be available and percentage of population units falling into different strata may be known. Under such circumstances usual stratified sampling can not be used and thus an effort is made to get over the problem through post-stratification which consists in selecting a sample from the whole population by the procedure of simple random sampling without replacement followed by the classification of the selected sample units by strata and then treating it as if it were stratified sample, for instance, see [1, 3–11].

It is further noted that in sample surveys, the information on an auxiliary variate correlated with the principal (study) variate under study is either readily available or may be made available by diverting a part of the survey resources. This information may be utilized to increase the precision of estimators of population mean \bar{Y} of the study variate y . Such an information is the known population mean \bar{X} of the auxiliary variate x . For illustration, the average farm size in a local government area or district may be known while the problem is to estimate the average area under a particular crop per farm. The strata may be formed according to farm size, the percentage of farms falling into different size groups may be known but the identity of farms within a size group may not be known, see [12].

We assume that the population comprises N units, which can be uniquely partitioned into L strata of size N_1, N_2, \dots, N_L such that $\sum_{h=1}^L N_h = N$. The strata weights $W_h = N_h/N$ ($h = 1, 2, \dots, L$) are assumed known. Let (y_{hi}, x_{hi}) ($i = 1, 2, \dots, N_h$) denote the values of variates (y, x) respectively for i -th unit in h -th stratum and \bar{Y}_h and \bar{X}_h denote strata means. A simple random sample of size n is drawn without replacement from the population which results into the configuration $\underline{n} = (n_1, n_2, \dots, n_L)$, n_h denoting the number of units in the sample falling in stratum h , $\sum_{h=1}^L n_h = n$. Assume that n is large enough so that the probability of n_h being zero is small (i.e. $Pr(n_h = 0) = 0$). Based on the foregoing procedure which is known as post-stratification, the usual unbiased post-stratified estimators for population means $\bar{Y} = \sum_{h=1}^L W_h \bar{Y}_h$ and $\bar{X} = \sum_{h=1}^L W_h \bar{X}_h$ of the study variate y and the auxiliary variate x are $\bar{y}_{PS} = \sum_{h=1}^L W_h \bar{y}_h$ and $\bar{x}_{PS} = \sum_{h=1}^L W_h \bar{x}_h$, where $\bar{y}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} y_{hi}$ and $\bar{x}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} x_{hi}$ are the means of the n_h sample units that fall into the h -th stratum whose size N_h is assumed to be known.

For given configuration of sample $\underline{n} = (n_1, n_2, \dots, n_L)$ we have

$$\begin{aligned} \text{Var}(\bar{y}_{PS}|\underline{n}) &= E((\bar{y}_{PS} - \bar{Y})^2|\underline{n}) = \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{hy}^2, \\ \text{Var}(\bar{x}_{PS}|\underline{n}) &= E((\bar{x}_{PS} - \bar{X})^2|\underline{n}) = \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{hx}^2, \\ \text{Cov}\{(\bar{y}_{PS}, \bar{x}_{PS})|\underline{n}\} \\ &= E\{((\bar{y}_{PS} - \bar{Y})(\bar{x}_{PS} - \bar{X}))|\underline{n}\} = \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{hxy}, \end{aligned}$$

see [1], where

$$\begin{aligned} f_h &= \frac{n_h}{N_h}, \quad S_{hy}^2 = \frac{1}{N_h-1} \sum_{i=1}^{N_h} (y_{hi} - \bar{Y}_h)^2, \\ S_{hx}^2 &= \frac{1}{N_h-1} \sum_{i=1}^{N_h} (x_{hi} - \bar{X}_h)^2, \quad S_{hxy} = \frac{1}{N_h-1} \sum_{i=1}^{N_h} (x_{hi} - \bar{X}_h)(y_{hi} - \bar{Y}_h). \end{aligned}$$

Using the results from [13] for $E(n_h^{-1})$, to the terms of order n^{-2} , we have

$$\begin{aligned}\text{Var}(\bar{y}_{PS}) &= \left(\frac{1-f}{n}\right) \sum_{h=1}^L W_h S_{hy}^2 + \left(\frac{N-n}{N-1}\right) \frac{1}{n^2} \sum_{h=1}^L (1-W_h) S_{hy}^2, \\ \text{Var}(\bar{x}_{PS}) &= \left(\frac{1-f}{n}\right) \sum_{h=1}^L W_h S_{hx}^2 + \left(\frac{N-n}{N-1}\right) \frac{1}{n^2} \sum_{h=1}^L (1-W_h) S_{hx}^2, \\ \text{Cov}(\bar{y}_{PS}, \bar{x}_{PS}) &= \left(\frac{1-f}{n}\right) \sum_{h=1}^L W_h S_{hxy} + \left(\frac{N-n}{N-1}\right) \frac{1}{n^2} \sum_{h=1}^L (1-W_h) S_{hxy},\end{aligned}$$

where $f = n/N$ is over all sampling fraction.

It is known that when the auxiliary information is used at the estimation stage, the ratio estimator is the best among a wide class of estimators when the relation between y and x , the variate under study and the auxiliary variate respectively, is a straight line through the origin and the variance of y about this line is proportional to x , see [14]. In such a situation the ratio estimator is as good as regression estimator. In many practical situations, the regression line does not pass through the origin. In these situations, the ratio estimator does not perform equally well as that of regression estimator. Keeping this fact in view and also due to the stronger intuitive appeal statisticians are more inclined towards the use of the ratio and the product estimators and hence a large amount of work has been carried out towards the modification of ratio and product estimators, for instance, see [11, 15–17] etc. These authors have proposed various estimators under simple random sampling without replacement (SRSWOR) and stratified random technique which under some realistic conditions is more efficient than the mean per unit estimator, the ratio and the product estimator are efficient as the linear regression estimator in optimum case. It is to be mentioned that the problem of estimation of population mean \bar{Y} of the study variate y based on post-stratification and auxiliary information has not attracted much attention of survey statisticians, for instance, [12] and [18].

In this paper, following approaches developed by [19] and [20], we have suggested a family of estimators of population mean \bar{Y} of the study variate y based on post-stratification using multi-auxiliary variate and its properties are studied.

When information on p -auxiliary variates x_1, x_2, \dots, x_p is available. Let W_h ($h = 1, 2, \dots, L$) and $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p$ be the known strata weights and the known population means of the auxiliary variates $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p$ respectively. Suppose the observations (y_{hi}, x_{khi}) , $i = 1, 2, \dots, n_h$, $h = 1, 2, \dots, L$ and $k = 1, 2, \dots, p$ are available. We denote

$$\begin{aligned}\bar{x}_{kPS} &= \sum_{h=1}^L W_h \bar{x}_{kh}, & \bar{x}_{kh} &= \frac{1}{n_h} \sum_{i=1}^{n_h} x_{khi}, \\ \bar{X}_k &= \sum_{h=1}^L W_h \bar{X}_{kh}, & \bar{X}_{kh} &= \frac{1}{N_h} \sum_{i=1}^{N_h} x_{khi}.\end{aligned}$$

Let $\bar{\underline{x}}_{PS}$ denote the column vector of p -elements $\bar{x}_{1PS}, \bar{x}_{2PS}, \dots, \bar{x}_{pPS}$. Superfix T over a column vector denotes the corresponding row vector.

Defining $\varepsilon_0 = (\bar{y}_{PS} - \bar{Y})$, $\varepsilon_k = (\bar{x}_{kPS} - \bar{X}_k)$ and $\varepsilon^T = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p)$, we have for a given configuration of $\underline{n} = (n_1, n_2, \dots, n_L)$, the values of the conditional expectations:

$$E(\varepsilon_0 | \underline{n}) = 0 = E(\varepsilon_k | \underline{n})$$

and if n_h is large, to terms of order n_h^{-1} , the conditional expected values are

$$\begin{aligned} E(\varepsilon_0^2 | \underline{n}) &= \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{h0}^2, \\ E(\varepsilon_k^2 | \underline{n}) &= \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{hk}^2, \\ E(\varepsilon_0 \varepsilon_k | \underline{n}) &= \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{h0k}, \\ E(\varepsilon_k \varepsilon_l | \underline{n}) &= \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{hkl}, \end{aligned} \tag{1}$$

where

$$\begin{aligned} S_{h0k} &= \rho_{h0k} S_{h0} S_{hk} \Rightarrow \rho_{h0k} = \frac{S_{h0k}}{S_{h0} S_{hk}}, \quad S_{hkl} = \rho_{hkl} S_{hk} S_{hl} \Rightarrow \rho_{hkl} = \frac{S_{hkl}}{S_{hk} S_{hl}}, \\ S_{h0}^2 &= \frac{1}{N_h - 1} \sum_{i=1}^{N_h} (y_{hi} - \bar{Y}_h)^2, \quad S_{hk}^2 = \frac{1}{N_h - 1} \sum_{i=1}^{N_h} (x_{hki} - \bar{X}_{hk})^2, \\ S_{h0k} &= \frac{1}{N_h - 1} \sum_{i=1}^{N_h} (y_{hi} - \bar{Y}_h) (x_{hki} - \bar{X}_{hk}), \\ S_{hkl} &= \frac{1}{N_h - 1} \sum_{i=1}^{N_h} (x_{hki} - \bar{X}_{hk}) (x_{hli} - \bar{X}_{hl}). \end{aligned}$$

Putting the above results in matrix notations, we have

$$E(\varepsilon | \underline{n}) = 0, \quad E(\varepsilon \varepsilon^T | \underline{n}) = D, \quad E(\varepsilon_0 \varepsilon | \underline{n}) = A, \tag{2}$$

where

$$\begin{aligned} A^T &= (a_1, a_2, \dots, a_p), \quad a_k = \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{h0k}, \\ D &= [d_{kl}]_{p \times p}, \quad d_{kl} = \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) S_{hkl}. \end{aligned}$$

The unconditional expectations are:

$$E(\varepsilon_0) = E(\varepsilon_k) = 0, \quad \text{for all } k = 1, 2, \dots, p$$

and for large n , to terms of order n^{-1} , the unconditional expected values are given by

$$\begin{aligned} E(\varepsilon_0^2) &= \frac{1-f}{n} \sum_{h=1}^L W_h S_{h0}^2, \\ E(\varepsilon_k^2) &= \frac{1-f}{n} \sum_{h=1}^L W_h S_{hk}^2, \\ E(\varepsilon_0 \varepsilon_k) &= \frac{1-f}{n} \sum_{h=1}^L W_h S_{h0k}, \\ E(\varepsilon_k \varepsilon_l) &= \frac{1-f}{n} \sum_{h=1}^L W_h S_{hkl}. \end{aligned} \tag{3}$$

Putting the above results in matrix notation, we have

$$E(\varepsilon) = 0, \quad E(\varepsilon \varepsilon^T) = A^*, \quad E(\varepsilon_0 \varepsilon) = D^*, \tag{4}$$

where

$$\begin{aligned} A^{*T} &= (a_1^*, a_2^*, \dots, a_p^*), \quad a_k^* = \frac{1-f}{n} \sum_{h=1}^L W_h S_{h0k}, \\ D^* &= [d_{kl}^*]_{p \times p}, \quad \text{and} \quad d_{kl}^* = \frac{1-f}{n} \sum_{h=1}^L W_h S_{hkl}. \end{aligned}$$

2 The suggested family of estimators

Let $\overline{\underline{X}}^T = (\overline{X}_1, \overline{X}_2, \dots, \overline{X}_p)$ denote the row vector of p elements $\overline{X}_1, \overline{X}_2, \dots, \overline{X}_p$. Whatever be the sample chosen let $(\overline{y}_{PS}, \overline{\underline{x}}_{PS}^T)$ assume values in a closed convex subset, Q , of the $(p+1)$ dimensional real space containing the point $(\overline{Y}, \overline{\underline{X}}^T)$. We suggest a family of post-stratified estimators for the population mean using multi-auxiliary variable as:

$$\widehat{\overline{Y}}_G = G(\overline{y}_{PS}, \overline{x}_{1PS}, \overline{x}_{2PS}, \dots, \overline{x}_{pPS}) = G(\overline{y}_{PS}, \overline{\underline{x}}_{PS}^T), \tag{5}$$

where $G(\overline{y}_{PS}, \overline{\underline{x}}_{PS}^T)$ is a function of $\overline{y}_{PS}, \overline{x}_{1PS}, \overline{x}_{2PS}, \dots, \overline{x}_{pPS}$ such that

$$G(\overline{Y}, \overline{\underline{X}}^T) = \overline{Y}, \quad \text{for all } \overline{Y} \tag{6}$$

and such that it satisfies the following conditions:

1. The function $G(\bar{y}_{PS}, \bar{x}_{PS}^T)$ is continuous and bounded in Q ,
2. The first and second order partial derivatives of the function $G(\bar{y}_{PS}, \bar{x}_{PS}^T)$ exist and are continuous and bounded in Q .

Expanding the function $G(\bar{y}_{PS}, \bar{x}_{PS}^T)$ about the point (\bar{Y}, \bar{X}^T) in a second order Taylor's series, we obtain:

$$\begin{aligned} \hat{Y}_G = & G(\bar{Y}, \bar{X}^T) + (\bar{y}_{PS} - \bar{Y}) \frac{\partial G(\cdot)}{\partial \bar{y}_{PS}} \Big|_{(\bar{Y}, \bar{X}^T)} + (\bar{x}_{PS} - \bar{X})^T G^{(1)}(\bar{Y}, \bar{X}^T) \\ & + \frac{1}{2} \left\{ (\bar{y}_{PS} - \bar{Y})^2 \frac{\partial^2 G(\cdot)}{\partial \bar{y}_{PS}^2} \Big|_{(\bar{y}_{PS}^*, \bar{x}_{PS}^{*T})} \right. \\ & + 2(\bar{y}_{PS} - \bar{Y})(\bar{x}_{PS} - \bar{X})^T \frac{\partial G^{(1)}(\cdot)}{\partial \bar{y}_{PS}} \Big|_{(\bar{y}_{PS}^*, \bar{x}_{PS}^{*T})} \\ & \left. + (\bar{x}_{PS} - \bar{X})^T G^{(2)}(\bar{y}_{PS}^*, \bar{x}_{PS}^{*T})(\bar{x}_{PS} - \bar{X}) \right\}, \end{aligned} \quad (7)$$

where $\bar{y}_{PS}^* = \bar{Y} + \xi(\bar{y}_{PS} - \bar{Y})$, $\bar{x}_{PS}^* = \bar{X} + \xi(\bar{x}_{PS} - \bar{X})$, $0 < \xi < 1$, $G^{(1)}$ denotes the p elements column vector of first partial derivatives of $G(\cdot)$ i.e. $G^{(1)T} = (G_1^{(1)}, G_2^{(1)}, \dots, G_p^{(1)})$ with $G_k^{(1)} = (\partial G(\bar{y}_{PS}, \bar{x}_{PS}) / \partial \bar{x}_{kPS})|_{(\bar{Y}, \bar{X}^T)}$ and $G^{(2)}$ denotes the $p \times p$ matrix of the second partial derivatives of $G(\cdot)$ with respect to \bar{x}_{PS} about the point (\bar{Y}, \bar{X}^T) . Expressing (7) in terms of ε 's and noting that $G(\bar{Y}, \bar{X}^T) = \bar{Y}^T$, we have

$$\begin{aligned} \hat{Y}_G = & \bar{Y} + \frac{\partial G(\cdot)}{\partial \bar{y}_{PS}} \Big|_{(\bar{Y}, \bar{X}^T)} + \varepsilon^T G^{(1)}(\bar{Y}, \bar{X}^T) \\ & + \frac{1}{2} \left\{ \varepsilon_0^2 \frac{\partial^2 G(\cdot)}{\partial \bar{y}_{PS}^2} \Big|_{(\bar{y}_{PS}^*, \bar{x}_{PS}^{*T})} + 2\varepsilon_0 \varepsilon^T \frac{\partial G^{(1)}(\cdot)}{\partial \bar{y}_{PS}} \Big|_{(\bar{y}_{PS}^*, \bar{x}_{PS}^{*T})} \right. \\ & \left. + \varepsilon^T G^{(2)}(\bar{Y}_{PS}^*, \bar{x}_{PS}^{*T}) \varepsilon \right\} \end{aligned} \quad (8)$$

Taking conditional expectation in (8) and noting that second derivatives are bounded. Thus we arrived at the following theorem:

Theorem 1.

$$E(\hat{Y}_G | \underline{n}) = \bar{Y} + o(n_h^{-1}).$$

From Theorem 1, it follows that the bias of the estimator \hat{Y}_G is of the order n_h^{-1} , and hence its contribution to the mean squared error of \hat{Y}_G will be of the order of n_h^{-2} .

Now we prove the following result:

Theorem 2. Up to terms of order n_h^{-1} , the conditional variance of \widehat{Y}_G is minimized for

$$G^{(1)}(\overline{Y}, \overline{X}^T) = -D^{-1}A \quad (9)$$

and the conditional minimum variance is given by

$$\text{Var}(\widehat{Y}_G | \underline{n}) = (1 - R^2) S_0^{*2}. \quad (10)$$

Proof. From (8), we have upto terms of order n_h^{-1} ,

$$\begin{aligned} \text{Var}(\widehat{Y}_G | \underline{n}) &= E\{(\widehat{Y}_G - \overline{Y})^2 | \underline{n}\} \\ &= E\left\{\left(\varepsilon_0 \frac{\partial G(\cdot)}{\partial \overline{y}_{PS}} \Big|_{(\overline{Y}, \overline{X}^T)} + \varepsilon^T G^{(1)}(\overline{Y}, \overline{X}^T)\right)^2 \Big| \underline{n}\right\} \\ &= E\left\{(\varepsilon_0 + \varepsilon^T G^{(1)}(\overline{Y}, \overline{X}^T))^2 | \underline{n}\right\}, \end{aligned}$$

from (6) which implies that $\frac{\partial G(\cdot)}{\partial \overline{y}_{PS}} \Big|_{(\overline{Y}, \overline{X}^T)} = 1$,

$$\begin{aligned} \text{Var}(\widehat{Y}_G | \underline{n}) &= \left[E(\varepsilon_0^2 | \underline{n}) + 2E(\varepsilon_0 \varepsilon^T | \underline{n}) G^{(1)}(\overline{Y}, \overline{X}^T) \right. \\ &\quad \left. + (G^{(1)}(\overline{Y}, \overline{X}^T))^T E(\varepsilon \varepsilon^T | \underline{n}) (G^{(1)}(\overline{Y}, \overline{X}^T)) \right] \\ &= S_0^{*2} + 2A^T G^{(1)}(\overline{Y}, \overline{X}^T) + (G^{(1)}(\overline{Y}, \overline{X}^T))^T D (G^{(1)}(\overline{Y}, \overline{X}^T)) \quad (11) \end{aligned}$$

which is minimized for

$$G_{opt}^{(1)} = -D^{-1}A = \delta_0 \quad (\text{say}), \quad (12)$$

where $S_0^{*2} = \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h}\right) S_{h0}^2$.

Thus the resulting conditional variance is given by

$$\min.\text{Var}(\widehat{Y}_G | \underline{n}) = (1 - R^2) S_0^{*2}, \quad (13)$$

where $R^2 = \frac{A^T D^{-1} A}{S_0^{*2}}$ and R is the multiple correlation coefficient between \overline{y}_{PS} and the vector \overline{x}_{PS} . Hence proved the Theorem 2. \square

The conditional variance of any estimator of the class (5) can be obtained from (11). From (11) the conditional minimum variance (i.e. $\min.\text{Var}(\widehat{Y}_G | \underline{n})$) is not larger than the conditional variance of the unbiased estimator \overline{y}_{PS} , since $A^T D^{-1} A > 0$.

Taking unconditional expectation in (8) and noting that second derivatives are bounded, we have:

Theorem 3.

$$E(\widehat{Y}_G) = \overline{Y} + o(n^{-1}).$$

Theorem 3 shows that the bias of the estimator \widehat{Y}_G is of the order n^{-1} , and hence its contribution to the mean square error (MSE) of \widehat{Y}_G will be of the order n^{-2} . Thus, to the first order of approximation, the unconditional variance of \widehat{Y}_G will be the same.

Theorem 4. *Upto terms of order n^{-1} , the unconditional $\text{Var}(\widehat{Y}_G)$ is minimized for*

$$G^{(1)}(\overline{Y}, \overline{X}^T) = -D^{*-1}A^* \quad (14)$$

and the unconditional minimum variance of \widehat{Y}_G is given by

$$\min.\text{Var}(\widehat{Y}_G) = (1 - R^{*2})S_0^{**2}, \quad (15)$$

where $S_0^{**2} = \frac{1-f}{n} \sum_{h=1}^L W_h S_{h0}^2$.

Proof. From (8), we have upto terms of order n^{-1} ,

$$\begin{aligned} \text{Var}(\widehat{Y}_G) &= E(\widehat{Y}_G - \overline{Y})^2 = E\left(\varepsilon_0 \frac{\partial G(\cdot)}{\partial \overline{y}_{PS}} \Big|_{(\overline{Y}, \overline{X}^T)} + \varepsilon^T G^{(1)}(\overline{Y}, \overline{X}^T)\right)^2 \\ &= E(\varepsilon_0 + \varepsilon^T G^{(1)}(\overline{Y}, \overline{X}^T))^2, \end{aligned}$$

from (6) which implies that $\frac{\partial G(\cdot)}{\partial \overline{y}_{PS}} \Big|_{(\overline{Y}, \overline{X}^T)} = 1$,

$$\begin{aligned} \text{Var}(\widehat{Y}_G) &= \left[E(\varepsilon_0^2) + 2E(\varepsilon_0 \varepsilon^T) G^{(1)}(\overline{Y}, \overline{X}^T) \right. \\ &\quad \left. + (G^{(1)}(\overline{Y}, \overline{X}^T))^T E(\varepsilon \varepsilon^T) (G^{(1)}(\overline{Y}, \overline{X}^T)) \right]. \end{aligned}$$

Using the results (3) and (4) in the above expression we get the unconditional variance over all possible distribution, for large n , to the terms of order n^{-1} , as:

$$\begin{aligned} \text{Var}(\widehat{Y}_G) &= S_0^{**2} + 2A^{*T} G^{(1)}(\overline{Y}, \overline{X}^T) \\ &\quad + ((G^{(1)}(\overline{Y}, \overline{X}^T))^T D^* (G^{(1)}(\overline{Y}, \overline{X}^T))) \end{aligned} \quad (16)$$

which is minimized for

$$G_{opt}^{(1)} = -D^{*-1}A^* = G_0^{(1)} \quad (\text{say}), \quad (17)$$

where $S_0^{**2} = \frac{1-f}{n_h} \sum_{h=1}^L W_h S_{h0}^2$.

Thus the resulting unconditional variance of \widehat{Y}_G is given by

$$\min.\text{Var}(\widehat{Y}_G) = (1 - R^{*2})S_0^{**2}, \quad (18)$$

where $R^{*2} = \frac{A^{*T} D^{*-1} A^*}{S_0^{**2}}$ and R^* is the multiple correlation coefficient between \overline{y}_{PS} and the vector \overline{x}_{PS} . Hence proved the Theorem 4. \square

The unconditional variance of any estimator of the class (5) can be obtained from (16). From (16), the $\min \text{Var}(\widehat{Y}_G)$ is not large than the unconditional variance of the unbiased estimator y_{PS} , since $A^{*T} D^{*-1} A^* > 0$.

Let $G^{(1)}(\overline{Y}, \overline{X}^T) = -\alpha G_0^{(1)} = -\alpha D^{*-1} A^*$, is a departure from the optimum value ($\alpha > 0$ is a constant), we have

$$\begin{aligned} \text{Var}(\widehat{Y}_G) &= [S_0^{**2} - 2\alpha A^{*T} D^{*-1} A^* + \alpha^2 A^{*T} D^{*-1} A^*] \\ &= [S_0^{**2} - \alpha(2 - \alpha) A^{*T} D^{*-1} A^*]. \end{aligned} \quad (19)$$

It is well known that the unconditional variance of the usual unbiased estimator \overline{y}_{PS} is

$$\text{Var}(\widehat{Y}_G) = \frac{1-f}{n} \sum_{h=1}^L W_h S_{h0}^2 = S_0^{**2}. \quad (20)$$

Thus for any $G^{(1)}(\overline{Y}, \overline{X}^T)$, we find from (18) and (19) that

$$\text{Var}(\overline{y}_{PS}) - \text{Var}(\widehat{Y}_G) = \alpha(2 - \alpha) A^{*T} D^{*-1} A^* \quad (21)$$

which shows that the proposed class of estimators \widehat{Y}_G would be better than usual unbiased estimator \overline{y}_{PS} as for as $0 < \alpha < 2$.

Remark 1. It is to be mentioned that optimum estimators in the class are not unique but all of them have the same variance given either by (13) or (18). We also note that in practice the value of $\delta_0 = -D^{-1} A$ at (12) or $G_0^{(1)} = -D^{*-1} A^* = \delta_0^*$ at (17) may not be known. However, they can be estimated by

$$\widehat{\delta}_0 = \widehat{\delta}_0^* - \widehat{D}^{-1} \widehat{A}, \quad (22)$$

where

$$\begin{aligned} \widehat{D} &= [\widehat{d}_{kl}]_{p \times p}, \quad \widehat{d}_{kl} = \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) s_{hkl}, \\ \widehat{A}^T &= (\widehat{a}_1, \widehat{a}_2, \dots, \widehat{a}_p), \quad \widehat{a}_k = \sum_{h=1}^L W_h^2 \left(\frac{1-f_h}{n_h} \right) s_{h0k}, \\ s_{h0k} &= \frac{1}{n_h - 1} \sum_{i=1}^{n_h} (y_{hi} - \overline{y}_h)(x_{hki} - \overline{x}_{hk}), \\ s_{hkl} &= \frac{1}{n_h - 1} \sum_{i=1}^{n_h} (x_{hki} - \overline{x}_{hk})(x_{hli} - \overline{x}_{hl}), \quad \overline{y}_h = \frac{1}{n_h} \sum_{i=1}^{n_h} y_{hi}, \\ \overline{x}_{kh} &= \frac{1}{n_h} \sum_{i=1}^{n_h} x_{khi}. \end{aligned}$$

In such a case we may define a class of estimators based on estimated optimum value $\hat{\delta}_0$ as:

$$\hat{\bar{Y}}_G^* = G^*(\bar{y}_{PS}, \bar{x}_{PS}^T, \hat{\delta}_0^T), \quad (23)$$

where $G^*(\bar{y}_{PS}, \bar{x}_{PS}^T, \hat{\delta}_0^T)$ is a function of $(\bar{y}_{PS}, \bar{x}_{PS}^T, \hat{\delta}_0^T)$ such that:

$$\begin{aligned} G^*(\bar{Y}, \bar{X}^T, \delta^T) &= \bar{Y}, \\ \frac{\partial G^*(\cdot)}{\partial \bar{y}_{PS}} \bigg|_{(\bar{Y}, \bar{X}^T, \delta^T)} &= 1, \\ \frac{\partial G^*(\cdot)}{\partial \bar{x}_{PS}^T} \bigg|_{(\bar{Y}, \bar{X}^T, \delta^T)} &= \delta_0 = -D^{-1}A, \\ \frac{\partial G^*(\cdot)}{\partial \hat{\delta}_0^T} \bigg|_{(\bar{Y}, \bar{X}^T, \delta^T)} &= 0. \end{aligned} \quad (24)$$

Under (24) the class of estimators $\hat{\bar{Y}}_G^*$ at (23) is expected to have, to the first order of approximation, the conditional and unconditional variances respectively as

$$\text{Var}(\hat{\bar{Y}}_G^* | \underline{n}) = \min. \text{Var}(\hat{\bar{Y}}_G^* | \underline{n}) = (1 - R^2) S_0^{*2} \quad (25)$$

and

$$\text{Var}(\hat{\bar{Y}}_G^*) = \min. \text{Var}(\hat{\bar{Y}}_G^*) = (1 - R^{*2}) S_0^{**2}. \quad (26)$$

3 Comparison with corresponding unstratified multivariate estimators

We assume that information on p auxiliary variates x_1, x_2, \dots, x_p is available. A simple random sample of size n is drawn from the given finite population of size N . Let y_i and x_i denote the values of the variates y and x_k of the i -th unit of the sample, $k = 1, 2, \dots, p$; $i = 1, 2, \dots, n$. Defining:

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^n \bar{y}_i, & \bar{x}_k &= \frac{1}{n} \sum_{i=1}^n x_{ki}, & \bar{Y} &= \sum_{i=1}^N y_i, & \bar{X}_k &= \frac{1}{N} \sum_{i=1}^N x_{ki}, \\ S_0^2 &= \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2, & S_k^2 &= \frac{1}{N-1} \sum_{i=1}^N (x_{ki} - \bar{X}_k)^2, \\ S_{0k} &= \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})(x_{ki} - \bar{X}_k). \end{aligned}$$

Further ρ_{0k} and ρ_{kl} denote the correlation coefficients between the variates y and x_k and between the x_k and x_l .

Define $\underline{x}^T = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p)$, $\varepsilon_0^* = (\bar{y} - \bar{Y})$ and $\varepsilon_k^* = (\bar{x}_k - \bar{X}_k)$ such that:

$$\begin{aligned} E(\varepsilon_0^*) &= 0, \quad E(\varepsilon_k^*) = 0 \quad \text{for all } k = 1, 2, \dots, p, \\ E(\varepsilon_0^{*2}) &= \frac{1-f}{n} S_0^2, \quad E(\varepsilon_k^{*2}) = \frac{1-f}{n} S_k^2, \\ E(\varepsilon_0^* \varepsilon_k^*) &= \frac{1-f}{n} b_k, \quad E(\varepsilon_k^* \varepsilon_l^*) = \frac{1-f}{n} q_{kl}, \end{aligned}$$

where $(kl) = 1, 2, \dots, p$, $b_k = S_{0k} = \rho_{0k} S_0 S_k$, $q_{kl} = \rho_{kl} S_k S_l$.

Putting the above results in matrix notations, we have

$$E(\varepsilon_0^*) = 0, \quad E(\varepsilon_0^* \varepsilon^*) = \frac{1-f}{n} b, \quad E(\varepsilon_0^* \varepsilon^{*T}) = \frac{1-f}{n} Q,$$

where, $b^T = (b_1, b_2, \dots, b_p)$, $Q = [q_{kl}]_{p \times p}$.

Let $\underline{X}^T = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)$ denote the row vector of p elements $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p$.

Whatever be the sample chosen, let $(\bar{y}, \underline{x}^T)$ assume values in a closed convex subset, W , of the $(p+1)$ dimensional real space containing the point $(\bar{Y}, \underline{X}^T)$. Following [21] one may define a class of estimator of population mean \bar{Y} as

$$\hat{Y}_G^{(1)} = G(\bar{y}, \underline{x}^T) = G(\bar{y}, \bar{x}_1, \bar{x}_2, \dots, \bar{x}_p), \quad (27)$$

where $G(\bar{y}, \underline{x}^T)$ is a function of $\bar{y}, \bar{x}_1, \bar{x}_2, \dots, \bar{x}_p$ such that

$$G(\bar{Y}, \underline{X}^T) = \bar{Y}, \quad \text{for all } \bar{Y}$$

and such that it satisfies the following conditions:

1. The function $G(\bar{y}, \underline{x}^T)$ is continuous and bounded in W .
2. The first and second order partial derivatives of the function $G(\bar{y}, \underline{x}^T)$ exist and are continuous and bounded in W .

To the first degree of approximation, the variance of $\hat{Y}_G^{(1)}$ is given by

$$\begin{aligned} \text{Var}(\hat{Y}_G^{(1)}) &= \frac{1-f}{n} \left[S_0^2 + 2b^T G^{(1)}(\bar{Y}, \underline{X}^T) \right. \\ &\quad \left. + (G^{(1)}(\bar{Y}, \underline{X}^T))^T Q (G^{(1)}(\bar{Y}, \underline{X}^T)) \right] \end{aligned} \quad (28)$$

which is minimized when

$$G^{(1)}(\bar{Y}, \underline{X}^T) = -Q^{-1}b = \eta_0 \quad (\text{say}), \quad (29)$$

where $G^{(1)}(\bar{Y}, \bar{X}^T)$ denotes the p elements column vector of the partial derivatives of $G(\bar{y}, \bar{x}^T)$ with respect to \bar{x}^T about the point (\bar{Y}, \bar{X}^T) .

Thus the resulting minimum variance of $\hat{Y}_G^{(1)}$ is given by

$$\min.\text{Var}(\hat{Y}_G^{(1)}) = \frac{1-f}{n}(1-R^{**2})S_0^2, \quad (30)$$

where $R^{**2} = (b^T Q^{-1} b)/S_0^2$ and R^{**} is the multiple correlation coefficient between y and (x_1, x_2, \dots, x_p) .

Following [18], and from (18) and (23) it can be shown that the proposed class of estimators \hat{Y}_G in post-stratified sampling is unconditionally more efficient than the corresponding unstratified class of estimators $\hat{Y}_G^{(1)}$.

Remark 2. In practice, the exact optimum value η_0 of $G^{(1)}(\bar{Y}, \bar{X}^T)$ at (29) is not known, it is available to replace it by its consistent estimate of η_0 from the sample data at hand. Thus following the procedure outlined in [21], we define a class of estimators for population mean \bar{Y} (based on estimated optimum values) as:

$$\hat{Y}_G^{(1)} = G(\bar{y}, \bar{x}^T, \hat{\eta}_0^T), \quad (31)$$

where

$$\begin{aligned} \hat{\eta}_0 &= -\hat{Q}^{-1}\hat{b} \text{ with } \hat{Q} = [\hat{q}_{kl}]_{p \times p}, \quad \hat{q}_{kl} = s_{kl}, \quad \hat{b} = (\hat{b}_1, \hat{b}_2, \dots, \hat{b}_k), \quad \hat{b}_k = s_{0k}, \\ s_{0k} &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})(x_{ki} - \bar{x}_k), \quad s_{kl} = \frac{1}{n-1} \sum_{i=1}^n (x_{ki} - \bar{x}_k)(x_{li} - \bar{x}_l). \end{aligned}$$

Now $G(\bar{y}, \bar{x}^T, \hat{\eta}_0^T)$ is a function of $(\bar{y}, \bar{x}^T, \hat{\eta}_0^T)$ such that:

$$\begin{aligned} G(\bar{Y}, \bar{X}^T, \eta_0^T) &= \bar{Y}, \\ \left. \frac{\partial G(\cdot)}{\partial \bar{y}} \right|_{(\bar{Y}, \bar{X}^T, \eta_0^T)} &= 1, \\ \left. \frac{\partial G(\cdot)}{\partial \bar{x}_{PS}} \right|_{(\bar{Y}, \bar{X}^T, \eta_0^T)} &= \eta_0, \\ \left. \frac{\partial G(\cdot)}{\partial \hat{\eta}_0} \right|_{(\bar{Y}, \bar{X}^T, \eta_0^T)} &= 0. \end{aligned} \quad (32)$$

Under the condition (32), it can be shown to the first degree of approximation that the variance of $\hat{Y}_G^{(1)}$ is

$$\text{Var}(\hat{Y}_G^{(1)}) = \min.\text{Var}(\hat{Y}_G^{(1)}) = \frac{1-f}{n}(1-R^{**2})S_0^2. \quad (33)$$

From (26) and (33), we

$$\text{Var}(\hat{\hat{Y}}_G^*) < \text{Var}(\hat{\hat{Y}}_G^{(1)}),$$

if

$$(1 - R^{*2})S_0^{**2} = \frac{1 - f}{n}(1 - R^{**2})S_0^2. \quad (34)$$

Thus the proposed class of estimators $\hat{\hat{Y}}_G^*$ based on estimated optimum values in post-stratified sampling would be better than the corresponding non-stratified class of estimators $\hat{\hat{Y}}_G^{(1)}$ based on estimated optimum values in simple random sampling without replacement (SRSWOR), if the condition (34) holds true.

4 Empirical study

In the empirical study, we consider the relative efficiency of the post-stratified sampling estimator $\hat{\hat{Y}}_G (= \hat{\theta}_1, \text{ say})$ and non-stratified estimator $\hat{\hat{Y}}_G^{(1)} (= \hat{\theta}_2, \text{ say})$ with respect to the simple sample mean estimator without using an auxiliary information $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i (= \hat{\theta}_0, \text{ say})$. The percent relative efficiency of the estimator $\hat{\theta}_j, j = 1, 2$ with respect to the estimator $\hat{\theta}_0$ is computed as:

$$RE(\hat{\theta}_0, \hat{\theta}_j) = \frac{V(\hat{\theta}_0)}{V(\hat{\theta}_j)} \times 100 \% = RE(0, j), \quad (\text{say}). \quad (35)$$

We consider the problem of estimation of Forced Expiratory Volume (FEV) of the Pulmonary Disease persons, based on a dataset of 654 persons available on the CD with the book by [22], by using their age and height at the estimation stage, and using the other variables gender and smoking status as post-stratification variables. Thus the population of 654 persons has been divided into four post-strata. Post-stratum 1 consists of non-smoking females (0, 0), post-stratum 2 consists of non-smoking males (0, 1), post-stratum 3 consists of smoking females (1, 0), and post-stratum 4 consists of smoking males (1, 1). The descriptive parameters of the three variables: *FEV*, *Age* and *Height* (*HT*) in the four post-strata are given in Table 1. The population correlation coefficients between the three variables in the four post-strata are given in Table 2.

In order to have a closer look at the data structure in four different post-strata, we have also devoted Fig. 1 to display the three variables.

To investigate various situations, we apply power transformations on the study variable in all the four strata as $Y_i = (FEV)^T$ for different choice of values of T in the range of 0.1 to 2.5 with a step of 0.1. The other two variables: $X_{1i} = (Age)$ and $X_{2i} = (Height)$ were used at the estimation stage. We decided to select a sample of size being 10 % of the total population size, and later we post-stratified the sample based on gender and smoking status into four different homogeneous groups. A total of

sample size $n = 65.4$ (can be rounded to 65) was selected from the entire population of $N = 654$ persons. Out of 65.4 persons, 27.9 persons were found to be from stratum-1, 3.9 persons were from stratum-2, 31.0 persons were from stratum-3 and 2.6 persons were from stratum-4. We used R-code given in the Appendix to produce the results shown in Table 3.

Table 1. Descriptive parameters of FEV , Age and HT .

Stratum	N_h	Mean	St. Dev	Minimum	Q_1	Median	Q_3	Maximum	Skewness	Kurtosis
<i>FEV</i>										
1	279	2.3792	0.6393	0.7910	1.8770	2.4170	2.8660	3.8160	-0.07	-0.69
2	39	2.9659	0.4229	2.1980	2.6770	3.0740	3.2080	3.8350	-0.25	-0.39
3	310	2.7344	0.9741	0.7960	1.9565	2.5475	3.3578	5.7930	0.70	-0.13
4	26	3.7430	0.8890	1.6940	3.3420	3.878	4.4300	4.8720	-0.89	0.13
<i>Age</i>										
1	279	9.366	2.693	3.00	8.00	9.00	11.00	18.00	0.42	0.48
2	39	13.256	2.245	10.00	11.00	13.00	15.00	19.00	0.65	0.42
3	310	9.687	2.778	3.00	8.00	10.00	11.00	19.00	0.43	0.33
4	26	13.923	2.465	9.00	12.00	14.00	16.00	18.00	-0.13	-0.78
<i>HT</i>										
1	279	59.605	4.739	46.00	57.00	60.50	63.00	71.00	-0.60	-0.10
2	39	64.551	2.291	60.00	63.00	65.00	66.00	69.50	0.09	-0.68
3	310	61.519	6.268	47.00	57.00	61.00	67.00	74.00	-0.11	-0.90
4	26	68.058	3.232	58.00	67.00	68.00	69.75	72.00	-1.63	3.86

Table 2. Pearson correlation coefficient values for four strata.

Correlations	Stratum-1 (0, 0)		Correlations	Stratum-2 (0, 1)	
	<i>FEV</i>	<i>Age</i>		<i>FEV</i>	<i>Age</i>
<i>Age</i>	0.767	—	<i>Age</i>	-0.047	—
<i>HT</i>	0.843	0.776	<i>HT</i>	0.251	-0.092
	Stratum-3 (1, 0)			Stratum-4 (1, 1)	
	0.822	—	<i>Age</i>	0.394	—
<i>HT</i>	0.883	0.842	<i>HT</i>	0.750	0.352

Table 3. Relative efficiency of the post-stratified and non-stratified estimators with respect to the sample mean estimator.

T	Correlations	Stratum-1	Stratum-2	Stratum-3	Stratum-4	$RE(0, 1)$	$RE(0, 2)$
0.2	ρ_{yx_1}	0.76839	-0.05719	0.82536	0.42073		
	ρ_{yx_2}	0.85847	0.24024	0.90837	0.78591	558.13	514.11
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		

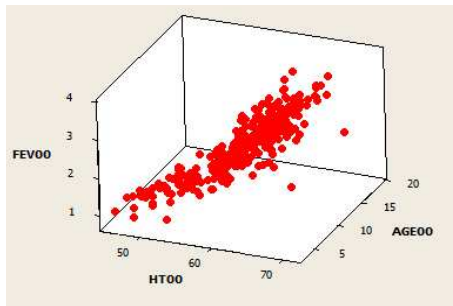
T	Correlations	Stratum-1	Stratum-2	Stratum-3	Stratum-4	$RE(0, 1)$	$RE(0, 2)$
0.3	ρ_{yx_1}	0.76675	-0.05889	0.82367	0.42475	557.63	518.16
	ρ_{yx_2}	0.85921	0.23834	0.91010	0.79138		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
0.4	ρ_{yx_1}	0.76623	-0.05933	0.82311	0.42574	556.72	518.43
	ρ_{yx_2}	0.85923	0.23785	0.91040	0.79273		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
0.5	ρ_{yx_1}	0.76604	-0.05947	0.82291	0.42607	556.34	518.44
	ρ_{yx_2}	0.85931	0.23768	0.91048	0.79317		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
0.6	ρ_{yx_1}	0.76597	-0.05953	0.82283	0.42620	556.17	518.44
	ρ_{yx_2}	0.85931	0.23762	0.91051	0.79335		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
0.7	ρ_{yx_1}	0.76593	-0.05956	0.82279	0.42626	556.10	518.44
	ρ_{yx_2}	0.85931	0.23759	0.91053	0.79343		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
0.8	ρ_{yx_1}	0.76592	-0.05957	0.82278	0.42628	556.06	518.43
	ρ_{yx_2}	0.85931	0.23757	0.91054	0.79347		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
0.9	ρ_{yx_1}	0.76591	-0.05957	0.82277	0.42630	556.05	518.43
	ρ_{yx_2}	0.85931	0.23757	0.91054	0.79348		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.0	ρ_{yx_1}	0.76591	-0.05957	0.82277	0.42630	556.05	518.43
	ρ_{yx_2}	0.85931	0.23757	0.91054	0.79348		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.1	ρ_{yx_1}	0.76592	-0.05957	0.82277	0.42629	556.06	518.43
	ρ_{yx_2}	0.85931	0.23757	0.91053	0.79347		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.2	ρ_{yx_1}	0.76593	-0.05956	0.82279	0.42626	556.09	518.44
	ρ_{yx_2}	0.85931	0.23758	0.91053	0.79344		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.3	ρ_{yx_1}	0.76598	-0.05954	0.82281	0.42622	556.14	518.44
	ρ_{yx_2}	0.85931	0.23760	0.91052	0.79339		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.4	ρ_{yx_1}	0.76599	-0.05951	0.82285	0.42616	556.23	518.44
	ρ_{yx_2}	0.85931	0.23764	0.91050	0.79330		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.5	ρ_{yx_1}	0.76606	-0.05946	0.82293	0.42604	556.37	518.44
	ρ_{yx_2}	0.85930	0.23770	0.91047	0.79313		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.6	ρ_{yx_1}	0.76618	-0.05936	0.82306	0.42582	556.63	518.44
	ρ_{yx_2}	0.85929	0.23780	0.91042	0.79284		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		

T	Correlations	Stratum-1	Stratum-2	Stratum-3	Stratum-4	$RE(0, 1)$	$RE(0, 2)$
1.7	ρ_{yx_1}	0.76640	-0.05919	0.82329	0.42542	557.05	518.38
	ρ_{yx_2}	0.85927	0.23800	0.91031	0.79230		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.8	ρ_{yx_1}	0.76681	-0.05884	0.82373	0.42463	557.73	518.11
	ρ_{yx_2}	0.85920	0.23839	0.91007	0.79123		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
1.9	ρ_{yx_1}	0.76755	-0.05814	0.82451	0.42301	558.47	516.96
	ρ_{yx_2}	0.85897	0.23918	0.90946	0.78903		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
2.0	ρ_{yx_1}	0.76878	-0.05666	0.82574	0.41945	557.30	511.93
	ρ_{yx_2}	0.85810	0.24082	0.90766	0.78420		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
2.1	ρ_{yx_1}	0.76992	-0.05338	0.82653	0.41105	542.98	490.27
	ρ_{yx_2}	0.85451	0.24439	0.90143	0.77282		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
2.2	ρ_{yx_1}	0.76560	-0.04576	0.81962	0.38930	471.43	410.63
	ρ_{yx_2}	0.83905	0.25223	0.87710	0.74336		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
2.3	ρ_{yx_1}	0.72779	-0.02781	0.76664	0.32980	296.36	247.49
	ρ_{yx_2}	0.77699	0.26831	0.78138	0.66352		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
2.4	ρ_{yx_1}	0.57491	0.00564	0.55799	0.20157	148.39	130.84
	ρ_{yx_2}	0.59403	0.28686	0.50344	0.50005		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		
2.5	ρ_{yx_1}	0.31371	-0.00472	0.28826	0.06937	105.93	103.92
	ρ_{yx_2}	0.33364	0.26138	0.20492	0.34442		
	$\rho_{x_1x_2}$	0.77642	-0.09219	0.84230	0.35209		

For $T = 0.2$, the values of the population correlation coefficients between FEV and Age are 0.76839, -0.05719, 0.82536 and 0.42073 in the first, second, third and fourth post-stratum, respectively. The values of the population correlation coefficients between FEV and $Height$ are 0.85847, 0.24024, 0.90837 and 0.78591 in the first, second, third and fourth stratum, respectively. In the same way, the values of the populations correlation coefficients between Age and $Height$ are 0.77642, -0.09219, 0.84230 and 0.35209 in the first, second, third and fourth stratum respectively. In this particular situation, the percent relative efficiency of the post-stratified sampling estimator $\hat{\theta}_1$ with respect to the simple sample mean estimator $\hat{\theta}_0$ remains 558.13 % and that of the non-stratified estimator $\hat{\theta}_2$ remains 514.11 %. In the same way, the results in Table 3 are readable for other values of T . It is to be noted that so long as the value of T is less than or equal to 2.0, the percent relative efficiency of the post-stratified estimator remains around 557 % and that of the non-stratified estimator remains around 511 %. As soon as the value of T becomes 2.3, the relative efficiency of the post-stratified estimator drastically reduces

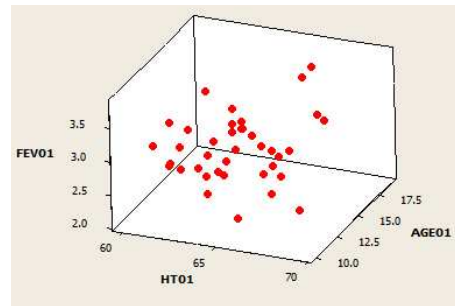
to 296.36 % and that of non-stratified estimator reduces to 247.49 %. For higher value of T equal to 2.5, the relative efficiency of the post-stratified sampling estimator reduces to 105.93 % and that of non-stratified sampling estimator reduces to 103.92 %. Thus, we conclude that the proposed post stratified sampling estimator can be used to estimate population mean of a study variable in the presence of multi-auxiliary variables more efficiently than a non-stratified sampling estimator.

3D Scatterplot of FEV00 vs AGE00 vs HT00



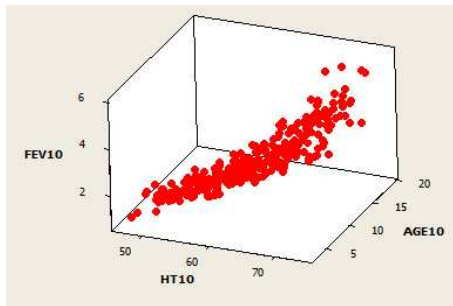
(a)

3D Scatterplot of FEV01 vs AGE01 vs HT01



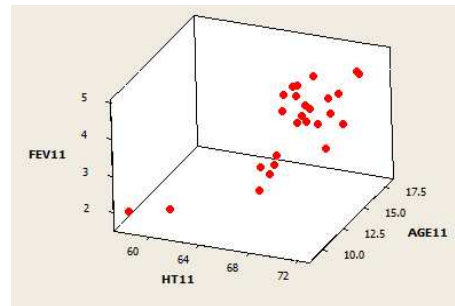
(b)

3D Scatterplot of FEV10 vs AGE10 vs HT10



(c)

3D Scatterplot of FEV11 vs AGE11 vs HT11



(d)

Fig. 1. Pictorial representation of four post-strata: (a) stratum-1 (Female = 0, Smoking = 0); (b) stratum-2 (Female = 0, Smoking = 1); (c) stratum-3 (Male = 1, Smoking = 0); (d) stratum-4 (Male = 1, Smoking = 1).

Appendix

```
#R-Code used in the simulation study (File Name: post2.r)
names<-c(0,0,0,0,0);
inpl1<-read.fwf("c:\\rc\\out00",c(14,10,8,5,5),header=FALSE,
sep="\t", as.is=FALSE,skip=0,col.names=names);
```

```

inp10<-read.fwf("c:\\rc\\out01",c(14,10,8,5,5),header=FALSE,
sep="\t",as.is=FALSE,skip=0,col.names=names);
inp01<-read.fwf("c:\\rc\\out10",c(14,10,8,5,5),header=FALSE,
sep="\t",as.is=FALSE,skip=0,col.names=names);
inp00<-read.fwf("c:\\rc\\out11",c(14,10,8,5,5),header=FALSE,
sep="\t",as.is=FALSE,skip=0,col.names=names);
y1<-c(inp11[[1]])
x11<-c(inp11[[2]])
x12<-c(inp11[[3]])
sex11<-c(inp11[[4]])
sk11<-c(inp11[[5]])
y2<-c(inp10[[1]])
x21<-c(inp10[[2]])
x22<-c(inp10[[3]])
sex10<-c(inp10[[4]])
sk10<-c(inp10[[5]])
y3<-c(inp01[[1]])
x31<-c(inp01[[2]])
x32<-c(inp01[[3]])
sex01<-c(inp01[[4]])
sk01<-c(inp01[[5]])
y4<-c(inp00[[1]])
x41<-c(inp00[[2]])
x42<-c(inp00[[3]])
sex00<-c(inp00[[4]])
sk00<-c(inp00[[5]])
np1<-length(y1)
np2<-length(y2)
np3<-length(y3)
np4<-length(y4)
np<-np1+np2+np3+np4
print(c(np1,np2,np3,np4,np))
w1<-np1/np
w2<-np2/np
w3<-np3/np
w4<-np4/np
ns<-0.10*np
ns1<-ns*np1/np
ns2<-ns*np2/np
ns3<-ns*np3/np
ns4<-ns*np4/np
print(c(ns1,ns2,ns3,ns4,ns))
f1<-(1-ns1/np1)/ns1
f2<-(1-ns2/np2)/ns2
f3<-(1-ns3/np3)/ns3
f4<-(1-ns3/np3)/ns3
t<-0.1

```

```

for (i in 1:25) {
t<-t+0.1
print(c('t=',t))
y1<-(y1)^t
x11<-(x11)
x12<-(x12)
ry1x11<-cov(y1,x11)/sqrt(var(y1)*var(x11))
ry1x12<-cov(y1,x12)/sqrt(var(y1)*var(x12))
rx11x12<-cov(x11,x12)/sqrt(var(x11)*var(x12))
print(c('ry1x11=',ry1x11,'ry1x12=',ry1x12,'rx11x12=',rx11x12))
y2<-(y2)^t
x21<-(x21)
x22<-(x22)
ry2x21<-cov(y2,x21)/sqrt(var(y2)*var(x21))
ry2x22<-cov(y2,x22)/sqrt(var(y2)*var(x22))
rx21x22<-cov(x21,x22)/sqrt(var(x21)*var(x22))
print(c('ry2x21=',ry2x21,'ry2x22=',ry2x22,'rx21x22=',rx21x22))
y3<-(y3)^t
x31<-(x31)
x32<-(x32)
ry3x31<-cov(y3,x31)/sqrt(var(y3)*var(x31))
ry3x32<-cov(y3,x32)/sqrt(var(y3)*var(x32))
rx31x32<-cov(x31,x32)/sqrt(var(x31)*var(x32))
print(c('ry3x31=',ry3x31,'ry3x32=',ry3x32,'rx31x32=',rx31x32))
y4<-(y4)^t
x41<-(x41)
x42<-(x42)
ry4x41<-cov(y4,x41)/sqrt(var(y4)*var(x41))
ry4x42<-cov(y4,x42)/sqrt(var(y4)*var(x42))
rx41x42<-cov(x41,x42)/sqrt(var(x41)*var(x42))
print(c('ry4x41=',ry4x41,'ry4x42=',ry4x42,'rx41x42=',rx41x42))
s02<-f1*w1^2*var(y1)+f2*w2^2*var(y2)+f3*w3^2*var(y3)+
f4*w4^2*var(y4)
a<-matrix(0,1,2)
d<-matrix(0,2,2)
a[1]<-f1*w1^2*cov(y1,x11)+f2*w2^2*cov(y2,x21)+
f3*w3^2*cov(y3,x31)+ f4*w4^2*cov(y4,x41)
a[2]<-f1*w1^2*cov(y1,x12)+f2*w2^2*cov(y2,x22)+
f3*w3^2*cov(y3,x32)+ f4*w4^2*cov(y4,x42)
d[1,1]<-f1*w1^2*cov(x11,x11)+f2*w2^2*cov(x21,x21)+
f3*w3^2*cov(x31,x31)+ f4*w4^2*cov(x41,x41)
d[1,2]<-f1*w1^2*cov(x11,x12)+f2*w2^2*cov(x21,x22)+
f3*w3^2*cov(x31,x32)+ f4*w4^2*cov(x41,x42)
d[2,1]<-d[1,2]
d[2,2]<-f1*w1^2*cov(x12,x12)+f2*w2^2*cov(x22,x22)+
f3*w3^2*cov(x32,x32)+ f4*w4^2*cov(x42,x42)
#print(d)

```

```
#print(a)
invd<-solve(d)
#print(invd)
out<-a%*%invd%*%t(a)
rsq<-out/s02
y<-c(y1,y2,y3,y4)
x1<-c(x11,x21,x31,x41)
x2<-c(x12,x22,x32,x42)
vary<-(1-ns/np)*var(y)/ns
rel<-vary*100/(s02*(1-rsq))
print(c("re (post stattification)=",rel))
b<-matrix(0,1,2)
q<-matrix(0,2,2)
b[1]<-cov(y,x1)
b[2]<-cov(y,x2)
q[1,1]<-cov(x1,x1)
q[1,2]<-cov(x1,x2)
q[2,1]<-q[1,2]
q[2,2]<-cov(x2,x2)
#print(q)
#print(b)
invq<-solve(q)
#print(invq)
out1<-b%*%invq%*%t(b)
rsq1<-out1/var(y)
re2<-100/(1-rsq1)
print(c("re (no stratification)=",re2))
}
```

Acknowledgements

The authors are thankful to the Journal Secretary Dr. Romas Baronas and a referee for fruitful comments on the original version of the manuscript.

References

1. D. Holt, T.M.F. Smith, Post-stratification, *J. Roy. Stat. Soc. A Sta.*, **142**, pp. 33–46, 1979.
2. P.V. Sukhatme, B.V. Sukhatme, S. Sukhatme, C. Ashok, *Sampling Theory of Surveys with Applications*, Iowa State University Press, Ames, Iowa, 1984.
3. W. Fuller, Estimation employing post-strata, *J. Am. Stat. Assoc.*, **61**, pp. 1172–1183, 1966.
4. D.C. Doss, H.O. Hartly, G.R. Somayajulu, An exact small sample theory for post-stratification *J. Stat. Plan. Infer.*, **3**, pp. 235–249, 1979.
5. P. Jagers, A. Oden, L. Trulsson, Post-stratification and ratio estimation: Usages of auxiliary information in survey sampling and opinion polls, *Int. Stat. Rev.*, **53**, pp. 221–238, 1985.

6. P. Jagers, Post-stratification against bias in sampling, *Int. Stat. Rev.*, **55**, pp. 159–167, 1986.
7. M.C. Agrawal, K.B. Panda, An efficient estimator in post-stratification, *Metron*, **51**, pp. 179–188, 1993.
8. M.C. Agrawal, K.B. Panda, On efficient estimation in post-stratification, *Metron*, **53**, pp. 107–115, 1995.
9. M. Ruiz Espejo, D. Pineda, On variance estimation for post-stratification: a review, *Metron*, **55**, pp. 209–220, 1997.
10. H.P. Singh, M. Ruiz Espejo, Improved post-stratified estimation, *B. Int. Statist. Inst.*, **60**, pp. 341–342, 2003.
11. H.P. Singh, G.K. Vishwakarma, An efficient variant of the product and ratio estimators in stratified random sampling, *Statistics in Transition*, **7**(6), pp. 1311–1325, 2006.
12. A.F. Ige, T.P. Tripathi, Estimation of population mean using post-stratification and auxiliary information, *Abacus*, **18**(2), pp. 265–276, 1989.
13. F. Stephan, The expected value and variance of the reciprocal and other negative powers of a positive Bernoullian variate, *Ann. Math. Stat.*, **16**, pp. 50–61, 1945.
14. W.G. Cochran, *Sampling Techniques*, 3rd ed., Wiley Eastern Ltd., New York, 1977.
15. H.P. Singh, A generalized class of estimators of ratio, product and mean using supplementary information on an auxiliary character in PPSWR sampling scheme, *Gujarat Stat. Rev.*, **13**(2), pp. 1–30, 1986.
16. H.P. Singh, M. Ruiz Espejo, On linear regression and ratio-product estimation of a finite population mean, *Statistician*, **52**(1), pp. 59–67, 2003.
17. H.P. Singh, G.K. Vishwakarma, A family of estimators of population mean using auxiliary information in stratified sampling, *Commun. Stat.-Theor. M.*, **37**(7), pp. 1038–1050, 2008.
18. R.K. Tuteja, S. Bahl, T.P. Tripathi, Sampling strategies for population mean based on post-stratification and multivariate auxiliary information, in: *Proceedings-III, International Symposium on Optimization and Statistics, Dec. 19–21, held at Aligarh Muslim University, India*, pp. 56–61, 1995.
19. S.K. Srivastava, A class of estimators using auxiliary information in sample surveys, *Can. J. Stat.*, **8**, pp. 253–254, 1980.
20. H.P. Singh, G.K. Vishwakarma, Estimation of mean using auxiliary information and post-stratification, *Commun. Stat.-Theor. M.*, 2009 (accepted).
21. S.K. Srivastava, H.S. Jhaggi, A class of estimators of the population mean using multi-auxiliary information, *Calcutta Stat. Assoc.*, **32**, pp. 47–56, 1983.
22. B. Rosner, *Fundamentals of Biostatistics*, Thomson-Brooks/Cole, 6rd ed., 2006.