

BIBLIOTECA DI SCIENZE STATISTICHE

SERVIZIO BIBLIOTECARIO NAZIONALE

BID P000850386

ACQ. 684 / '02 INV. 82184

COLL. \_\_\_\_\_ CLASS. 5-666.WP.2002/8

**Optimal estimation for  
finite population  
parameters in two phase  
sampling**

G. Diana, C. Tommasi

2002.9

UNIVERSITA' DI PADOVA  
BIBLIOTECA DI SCIENZE STATISTICHE  
Via C. Battisti, 241 - 35121 PADOVA

**Dipartimento di Scienze Statistiche  
Università degli Studi  
Via C. Battisti 241-243  
35121 Padova**

luglio 2002

BIBLIOTECA DI SCIENZE STATISTICHE

SEZIONE DI STATISTICA

NUMERO 2584

ANNO 1971

CLASSE 5 (6) (7) (8) (9)

COLL. 2

Optimal estimation for  
finite population  
parameters in two phase  
sampling

G. Ghata, C. Tamara

1971

UNIVERSITA' DI PADOVA  
BIBLIOTECA DI SCIENZE STATISTICHE  
VIA S. BENEDETTO 171 - 35100 PADOVA

Dipartimento di Scienze Statistiche  
Università degli Studi  
Via C. Battisti 241-243  
35131 Padova

luglio 1972

# Optimal estimation for finite population parameters in two phase sampling

Giancarlo Diana      Chiara Tommasi

## Abstract

In this paper we propose a general approach for estimating a finite population parameter in double sampling. When two dependent samples are drawn, several estimators were proposed to estimate the population mean, ratio and variance. While there are few proposals in double sampling with independent samples. We treat both cases, i.e. dependent and independent samples, showing that all the proposed estimators can be obtained as particular cases of a unique general class. The minimum variance bound for any estimator in this class is provided (at the first order of approximation). Furthermore, a chain regression type estimator which reaches this minimum is found.

**keywords:** two-phase sampling – dependent and independent samples – auxiliary variable – regression type estimator.

## 1 Introduction

The estimation of a parameter of interest when at least one auxiliary variable is available, was widely discussed. It is well known that if there are not information about the population characteristics of the auxiliary variable then a double sampling scheme may be used. We have two possibilities: the second phase sample may be drawn either from the first phase units (i.e. dependent samples) or independent of them from the population (i.e. independent samples). In the first case several estimators were proposed to estimate the population mean and the ratio of two means, while only few papers concern the population variance. For instance, the population mean estimation is treated, among others, by Srivastava *et al.* (1990), Singh and Singh (1991), Upadhyaya *et al.* (1992), Tracy and Singh (1999) and Singh and Espejo (2000). While Chand (1975), Srivastava *et al.* (1988,1989), Khare (1991), Singh *et al.* (1994), Singh and Singh (1994), Prasad *et al.* (1996) are some of

the authors who deal with the estimation of ratio of two means. On the other hand, only Singh (1991) and Gupta *et al.* (1992) work about the population variance in a two-phase sampling.

On the other hand, the double sampling with independent samples was not so fully treated. Only Khare (1991) and Singh and Singh (1994) estimate the ratio of two means using this sampling scheme, when only one auxiliary variable is available.

In this paper we propose a general approach for estimating a finite population parameter in double sampling, both for dependent and independent samples. Specifically, in Section 2 we give a general class of estimators which includes all the previously quoted estimators. Furthermore, we provide the minimum attainable MSE for the estimators in this class, at the first order of approximation. We call it the minimum variance bound, since MSE and variance are the same at this order of approximation. We suggest also a method to get an estimator which reaches this minimum. Such estimator has a nice interpretation when we are estimating a population mean, it is a chain regression type estimator. For this reason we call it the "best" estimator in the class, even if it is not unique. All the results of Section 2 are valid whatever the parameter of interest and the used sampling scheme. Section 3 specifies such general results to the estimation of a population mean, when dependent samples are drawn. Then, in the same context, Section 4 is devoted to the estimation of the ratio and the product of two means and the population variance. These subjects are discussed exhaustively in Diana and Tommasi (2001, 2002a, 2002b) but here they are unified. In fact, we also show that the estimation of these parameters is equivalent to the one of a mean, at least approximatively. Therefore, a central role is played by the mean case. Finally, in Section 5 we develop the mean case for independent samples. As for the dependent samples case, the estimation of the ratio or product and of the variance follows immediately from the mean estimation. Some general conclusions and remarks are given in Section 6.

## 2 A general class of estimators

Let  $Y$  be a  $m \times 1$  vector of study variables and  $T_Y$  be a scalar parameter of interest. Assuming that two auxiliary variables,  $X$  and  $Z$ , are related with  $Y$  but no information about the population parameters of  $X$  are available, then estimation of  $T_Y$  can be based on double sampling. Thus, we assume that a preliminary large sample of  $n'$  ( $n' < N$ ) units is drawn by some sampling scheme. At this phase only  $X$  and  $Z$  are measured. As a second step, a second sample of size  $n$  ( $n < n'$ ) is drawn. Without loss of generality, at this

phase all the variables  $Y$ ,  $X$  and  $Z$  are observed. From now on ' will denote quantities related to the first phase sample only.

If  $\mathbf{T}_X$  and  $\mathbf{T}_Z$  are  $k \times 1$  and  $l \times 1$  vectors of population parameters of  $X$  and  $Z$  and  $t_X$ ,  $t'_X$  and  $t_Z$ ,  $t'_Z$  are the corresponding first and second phase unbiased estimates, then the proposed class of estimators is defined by

$$\hat{T}_Y = g(t_Y, \mathbf{t}^T), \quad (1)$$

where  $t_Y$  is, at least approximatively, an unbiased sample estimate of  $T_Y$ ,  $\mathbf{t} = (t_X^T, t_Z^T, t_X'^T, t_Z'^T)^T$  and function  $g$  is such that

- a.  $g : \mathcal{S} \rightarrow \mathbb{R}$  where  $\mathcal{S} \in \mathbb{R}^{2(k+l)+1}$  is a convex and bounded set which contains the point  $(T_Y, \mathbf{T}^T)^T$ , where  $\mathbf{T} = E(\mathbf{t}) = (\mathbf{T}_X^T, \mathbf{T}_Z^T, \mathbf{T}_X'^T, \mathbf{T}_Z'^T)^T$ .
- b. It is a continuous and bounded function in  $\mathcal{S}$ .
- c. Its first and second partial derivatives are continuous and bounded in  $\mathcal{S}$ .
- d.  $g(t_Y, \mathbf{T}^T) = t_Y$ .

Let

$$g_0 = \left. \frac{\partial g(t_Y, \mathbf{t}^T)}{\partial t_Y} \right|_{(t_Y, \mathbf{t}^T) = (T_Y, \mathbf{T}^T)}$$

be the partial derivative of  $g$  with respect to  $t_Y$  and

$$\mathbf{g}_t = (\mathbf{g}_X^T, \mathbf{g}_Z^T, \mathbf{g}_X'^T, \mathbf{g}_Z'^T)^T = \left. \frac{\partial g(t_Y, \mathbf{t}^T)}{\partial \mathbf{t}} \right|_{(t_Y, \mathbf{t}^T) = (T_Y, \mathbf{T}^T)}$$

be the  $2(k+l)$  vector of the partial derivatives of  $g$  with respect to each component of  $\mathbf{t}$ . Specifically,  $\mathbf{g}_X$  and  $\mathbf{g}_X'$  are the  $k \times 1$  vectors of the partial derivatives of  $g$  with respect to the elements of  $t_X$  and  $t_X'$ , while  $\mathbf{g}_Z$  and  $\mathbf{g}_Z'$  are  $l \times 1$  vectors defined in the same way but referring to  $t_Z$  and  $t_Z'$ .

From point (d), we have that  $g(T_Y, \mathbf{T}^T) = T_Y$  and  $g_0 = 1$ . Thus, expanding  $\hat{T}_Y$  at the point  $(T_Y, \mathbf{T}^T)$  in a second order Taylor's series we have

$$\begin{aligned} \hat{T}_Y &\cong t_Y + (\mathbf{t} - \mathbf{T})^T \mathbf{g}_t = t_Y + (\mathbf{g}_X - \mathbf{T}_X)^T \mathbf{g}_X + (\mathbf{g}_Z - \mathbf{T}_Z)^T \mathbf{g}_Z \\ &\quad + (t'_X - \mathbf{T}_X)^T \mathbf{g}_X' + (t'_Z - \mathbf{T}_Z)^T \mathbf{g}_Z'. \end{aligned} \quad (2)$$

Since the population parameters of  $X$  are unknown, we have to impose the following constraint,  $\mathbf{g}_X = -\mathbf{g}_X'$ . For computational convenience, it is useful to change parameterization and to rewrite equation (2) as

$$\hat{T}_Y \cong t_Y + (t_V - t'_V)^T \mathbf{g}_V + (t'_Z - \mathbf{T}_Z)^T \mathbf{g}_Z \quad (3)$$

where  $t_V = (t_X^T, t_Z^T)^T$ ,  $t'_V = (t_X'^T, t_Z'^T)^T$ ,  $g_V = (g_X^T, g_Z^T)^T$  and  $g_Z = (g_Z + g'_Z)$ .  
The first order approximation for  $MSE(\hat{T}_Y)$  is given by

$$MSE(\hat{T}_Y) \cong S_{t_Y}^2 + g_V^T S_{t_V - t'_V, t_V - t'_V} g_V + \ddot{g}_Z^T S_{t'_Z, t'_Z} \ddot{g}_Z + 2 \left( g_V^T S_{t_V - t'_V, t_Y} + \ddot{g}_Z^T S_{t'_Z, t_Y} + g_V^T S_{t_V - t'_V, t'_Z} \ddot{g}_Z \right). \quad (4)$$

where

$$S_{t_Y}^2 = E[(t_Y - T_Y)^2], \quad S_{t, t..} = E[(t. - T.)(t.. - T..)^T].$$

This general expression for  $MSE(\hat{T}_Y)$  will be specialized for dependent and independent samples in the next sub-sections.

## 2.1 Dependent Samples

When the second phase sample is drawn from the first phase one, then the first order approximation for  $MSE(\hat{T}_Y)$  is

$$MSE(\hat{T}_Y) \cong S_{t_Y}^2 + g_V^T (S_{t_V, t_V} - S_{t'_V, t'_V}) g_V + \ddot{g}_Z^T S_{t'_Z, t'_Z} \ddot{g}_Z + 2 \left( g_V^T S_{t_V - t'_V, t_Y} + \ddot{g}_Z^T S_{t'_Z, t_Y} \right). \quad (5)$$

Minimizing (5) with respect to  $g_V$  and  $\ddot{g}_Z$  we obtain the optimum values

$$g_V^* = - (S_{t_V, t_V} - S_{t'_V, t'_V})^{-1} S_{t_V - t'_V, t_Y} \quad \text{and} \quad \ddot{g}_Z^* = - S_{t'_Z, t'_Z}^{-1} S_{t'_Z, t_Y}, \quad (6)$$

and replacing them in (5) we get the minimum first order approximation of  $MSE(\cdot)$ , denoted by  $MSE^*(\cdot)$ . It is the minimum variance bound for all the estimators based on the auxiliary vector  $t$  and has a simple form,

$$MSE^*(\hat{T}_Y) = S_{t_Y}^2 \left( 1 - \rho_{t_Y, t_V - t'_V}^2 - \rho_{t_Y, t'_Z}^2 \right), \quad (7)$$

where for any one random vector  $U$ ,

$$\rho_{t_Y, U}^2 = \frac{S_{U, t_Y}^T S_{U, U}^{-1} S_{U, t_Y}}{S_{t_Y}^2}.$$

Finally, replacing  $g_V^*$  and  $\ddot{g}_Z^*$  in (3) we get an optimum estimator. For all the estimation problems that will be considered, it is a regression type estimator. Thus, we will call it "the best" estimator in the class  $g$ .

## 2.2 Independent Samples

When both the first and the second phase samples are drawn independently from the population then the first order approximation for  $MSE(\hat{T}_Y)$ , given in (4), becomes

$$\begin{aligned} MSE(\hat{T}_Y) &\cong S_{t_Y}^2 + \mathbf{g}_V^T (\mathbf{S}_{t_V, t_V} + \mathbf{S}_{t'_V, t'_V}) \mathbf{g}_V + \mathbf{g}_Z^T \mathbf{S}_{t'_Z, t'_Z} \mathbf{g}_Z \\ &+ 2 (\mathbf{g}_V^T \mathbf{S}_{t_V, t_Y} - \mathbf{g}_V^T \mathbf{S}_{t'_V, t'_Z} \mathbf{g}_Z). \end{aligned} \quad (8)$$

Minimizing (8) with respect to  $\mathbf{g}_V$  and  $\mathbf{g}_Z$  we obtain the optimum values

$$\begin{aligned} \mathbf{g}_V^* &= -(\mathbf{S}_{t_V, t_V} + \mathbf{S}_{t'_V, t'_V|t'_Z})^{-1} \mathbf{S}_{t_V, t_Y} = -\mathbf{S}_{t_V+t'_V, t_V+t'_V|t'_Z}^{-1} \mathbf{S}_{t_V+t'_V, t_Y|t'_Z} \\ \mathbf{g}_Z^* &= \mathbf{S}_{t'_Z, t'_Z}^{-1} \mathbf{S}_{t'_V, t'_Z}^T \mathbf{g}_V^*, \end{aligned} \quad (9)$$

where, if  $U_i$ ,  $i = 1, 2, 3$ , are any three random vectors, then

$$\mathbf{S}_{U_1, U_2|U_3} = \mathbf{S}_{U_1, U_2} - \mathbf{S}_{U_1, U_3} \mathbf{S}_{U_3, U_3}^{-1} \mathbf{S}_{U_3, U_2}.$$

Replacing  $\mathbf{g}_V^*$  and  $\mathbf{g}_Z^*$  in (8) we have the minimum first order approximation of  $MSE(\hat{T}_Y)$ ,

$$MSE^*(\hat{T}_Y) = S_{t_Y}^2 (1 - \rho_{t_Y, t_V+t'_V|t'_Z}^2), \quad (10)$$

where we have set

$$\rho_{t_Y, t_V+t'_V|t'_Z}^2 = \frac{\mathbf{S}_{t_V+t'_V, t_Y|t'_Z}^T \mathbf{S}_{t_V+t'_V, t_V+t'_V|t'_Z}^{-1} \mathbf{S}_{t_V+t'_V, t_Y|t'_Z}}{S_{t_Y}^2}.$$

As in the previous sub-section, replacing  $\mathbf{g}_V^*$  and  $\mathbf{g}_Z^*$  in (3) we get an optimum estimator.

## 3 Estimation of a finite population mean

All the results of the previous section are valid whatever the used sampling scheme and the involved auxiliary statistics. If both  $t_Y$  and all the auxiliary statistics are, at least at the first order of approximation, sample means, then the computations of Section 2 can be developed further, when a specific sampling scheme is chosen. For this reason, the estimation of a finite population mean plays a central role. This case is treated in this section for dependent samples and then the obtained results are used in Section 4, where the estimation of ratio and product of two population means and the estimation of the population variance are dealt with. It will be seen that the

sample estimates of such parameters may be approximated by sample means, at least at the first order of approximation.

Henceforth we assume that both the first and the second phase samples are drawn by a simple random sample without replacement (SRSWOR). Let  $\mathcal{U} = \{1, \dots, i, \dots, N\}$  be a finite population.  $U_i$  denotes the value of any one observable  $h \times 1$  random vector  $\mathbf{U}$ , for the  $i$ -th population unit. On the other hand, the  $i$ -th sampled value will be denoted by the lowercase letter  $u_i$ . The population mean and covariance matrix of  $\mathbf{U}$  are  $\bar{\mathbf{U}} = \sum_{i=1}^N \mathbf{U}_i / N$  and  $\mathbf{S}_{\mathbf{U}\mathbf{U}} = \{S_{U_r, U_s}\}$ , respectively, where  $S_{U_r, U_s} = \sum_{i=1}^N (U_{ri} - \bar{U}_r)(U_{si} - \bar{U}_s) / (N - 1)$ . While  $\bar{\mathbf{u}}'$ ,  $\mathbf{s}'_{\mathbf{U}\mathbf{U}} = \{s'_{U_r, U_s}\}$  and  $\bar{\mathbf{u}}$ ,  $\mathbf{s}_{\mathbf{U}\mathbf{U}} = \{s_{U_r, U_s}\}$  are the first and second phase sample means and covariance matrices of  $\mathbf{U}$ . Specifically,  $s_{U_r, U_s} = \sum_{i=1}^n (u_{ri} - \bar{u}_r)(u_{si} - \bar{u}_s) / (n - 1)$  and  $s'_{U_r, U_s}$  is defined in the same way with  $n'$  instead of  $n$ . When  $U_r = U_s$  we will write  $S_{U_r, U_r}^2$ ,  $s_{U_r, U_r}^2$  and  $s_{U_r, U_r}^{\prime 2}$  for  $S_{U_r, U_r}$ ,  $s_{U_r, U_r}$  and  $s'_{U_r, U_r}$ , respectively.

Now we are interested in  $T_Y = \bar{Y}$ , thus  $t_Y = \bar{y}$ . Furthermore we consider the following auxiliary vectors,  $\mathbf{t}_X = (\bar{x}, s_X^2)^T$ ,  $\mathbf{t}_Z = (\bar{z}, s_Z^2)^T$ ,  $\mathbf{t}'_X$  and  $\mathbf{t}'_Z$ . With this notation equation (3) becomes

$$\bar{y}_g \cong \bar{y} + (\mathbf{t}_V - \mathbf{t}'_V)^T \mathbf{g}_V + (\mathbf{t}'_Z - \mathbf{T}_Z)^T \ddot{\mathbf{g}}_Z. \quad (11)$$

For a better interpretation of the results, it is convenient to define a random vector  $\mathbf{V} = (\tilde{\mathbf{X}}^T, \tilde{\mathbf{Z}}^T)^T$  where, for any variable  $U$ ,  $\tilde{\mathbf{U}} = (U, \delta_U)^T$  and  $\delta_U = (U - \bar{U})^2$ . In addition, let  $\mathbf{U}_1$  and  $\mathbf{U}_2$  be any two random vector of size  $h \times 1$  and  $k \times 1$ . Then  $\mathbf{S}_{\mathbf{U}_1, \mathbf{U}_2} = \{S_{U_1, U_2}\}$  is the  $h \times k$  matrix of the covariances among the elements of  $\mathbf{U}_1$  and  $\mathbf{U}_2$ .

Setting

$$\theta_1 = \left( \frac{1}{n'} - \frac{1}{N} \right), \quad \theta_2 = \left( \frac{1}{n} - \frac{1}{N} \right) \quad \text{and} \quad \theta = \left( \frac{1}{n} - \frac{1}{n'} \right)$$

we have that, at least up to terms of order  $O(n^{-1})$ ,

$$S_{t_Y}^2 \cong \theta_2 S_Y^2, \quad \mathbf{S}_{\mathbf{t}_V, \mathbf{t}_V} \cong \theta_2 \mathbf{S}_{\mathbf{V}, \mathbf{V}}, \quad \mathbf{S}_{\mathbf{t}'_V, \mathbf{t}'_V} \cong \theta_1 \mathbf{S}_{\mathbf{V}, \mathbf{V}}, \quad \mathbf{S}_{\mathbf{t}'_Z, \mathbf{t}'_Z} \cong \theta_1 \mathbf{S}_{\tilde{\mathbf{Z}}, \tilde{\mathbf{Z}}},$$

$$\mathbf{S}_{\mathbf{t}_V, \mathbf{t}_Y} \cong \theta_2 \mathbf{S}_{\mathbf{V}, Y}, \quad \mathbf{S}_{\mathbf{t}'_V, \mathbf{t}_Y} \cong \theta_1 \mathbf{S}_{\mathbf{V}, Y} \quad \text{and} \quad \mathbf{S}_{\mathbf{t}'_Z, \mathbf{t}_Y} \cong \theta_1 \mathbf{S}_{\tilde{\mathbf{Z}}, Y}.$$

With these mean values, from (6) we have that the "best" values  $\mathbf{g}_V^*$  and  $\ddot{\mathbf{g}}_Z^*$  are the partial regression coefficients of  $Y$  on  $\mathbf{V}$  and  $\tilde{\mathbf{Z}}$ , respectively, that is

$$\mathbf{g}_V^* = -\mathbf{S}_{\mathbf{V}, \mathbf{V}}^{-1} \mathbf{S}_{\mathbf{V}, Y} = - \begin{bmatrix} \beta_{Y, X | \delta_X, \tilde{\mathbf{Z}}} \\ \beta_{Y, \delta_X | X, \tilde{\mathbf{Z}}} \\ \beta_{Y, Z | \tilde{X}, \delta_Z} \\ \beta_{Y, \delta_Z | \tilde{X}, Z} \end{bmatrix} \quad \text{and} \quad \ddot{\mathbf{g}}_Z^* = -\mathbf{S}_{\tilde{\mathbf{Z}}, \tilde{\mathbf{Z}}}^{-1} \mathbf{S}_{\tilde{\mathbf{Z}}, Y} = - \begin{bmatrix} \beta_{Y, Z | \delta_Z} \\ \beta_{Y, \delta_Z | Z} \end{bmatrix}$$



Replacing such optimal values for  $g_V$  and  $g_Z$  in (5) we get

$$\text{MSE}^*(\bar{y}_g) \cong \theta_2 S_Y^2 \left( 1 - \frac{\theta}{\theta_2} \rho_{Y,V}^2 - \frac{\theta_1}{\theta_2} \rho_{Y,Z}^2 \right), \quad (12)$$

where  $\rho_{\cdot}$  denotes the multiple correlation coefficient between the specified variables. Up to terms of  $O(n^{-1})$ ,  $\text{MSE}^*(\bar{y}_g)$  is the minimum variance bound for all the estimators based on the auxiliary vector  $t$ , which summarizes now the first and second phase sample means and variances of  $X$  and  $Z$ .

Replacing  $g_V^*$  and  $g_Z^*$  in (11) we get the following chain regression type estimator

$$\begin{aligned} \bar{y}_{reg} = & \bar{y} - \beta_{Y,X|\delta_X,\bar{Z}}(\bar{x} - \bar{x}') - \beta_{Y,\delta_X|X,\bar{Z}}(s_X^2 - s_X'^2) - \beta_{Y,Z|\bar{x},\delta_Z}(\bar{z} - \bar{z}') \\ & - \beta_{Y,\delta_Z|\bar{x},Z}(s_Z^2 - s_Z'^2) - \beta_{Y,Z|\delta_Z}(\bar{z}' - \bar{Z}) - \beta_{Y,\delta_Z|Z}(s_Z'^2 - S_Z^2), \quad (13) \end{aligned}$$

whose bias is of order  $n^{-1}$ . Expression (13) shows that the construction of  $\bar{y}_{reg}$  is based on two steps. At the first step a regression type estimator for the unknown first phase sample mean of  $Y$  is provided. Of course, this estimator is based on  $X$  and  $Z$  conditionally to the first phase sample. Then such estimator is improved through another regression type estimator which uses the first phase sample information about  $Z$ . Actually, a similar interpretation was provided by Sahoo and Sahoo (1999) in a less general context.

**Remark 1.** Usually the partial regression coefficients of  $Y$  on  $V$  or  $\bar{Z}$  are unknown but replacing suitable estimates of such coefficients in (13) we get a new estimator which is equivalent to  $\bar{y}_{reg}$ , at the first order of approximation.

## 4 Other parameter estimation

The estimation of ratio or product of two population means, i.e.  $R = \bar{Y}_0/\bar{Y}_1$  or  $P = \bar{Y}_0\bar{Y}_1$ , and the population variance,  $S_Y^2 = \sum_{i=1}^N (Y_i - \bar{Y})^2/(N-1)$ , follow immediately from the mean case. In fact for all these parameters  $t_Y$  is approximatively a sample mean.

Let us first consider the ratio estimation, thus  $T_Y = R$  and  $t_Y = \hat{R} = \bar{y}_0/\bar{y}_1$ . Expanding  $\hat{R}$  at the point  $(\bar{Y}_0, \bar{Y}_1)$  in a second order Taylor's series, we have

$$\hat{R} \cong R \left( 1 + \frac{\bar{y}_0}{\bar{Y}_0} - \frac{\bar{y}_1}{\bar{Y}_1} \right) = R(1 + \bar{d}),$$

where  $\bar{d}$  is the second phase sample mean of  $D = Y_0/\bar{Y}_0 - Y_1/\bar{Y}_1$ . Since  $\hat{R}$  can be considered as a convenient sample mean (at the first order of approximation), we can use all the results of the previous section replacing  $Y$

with  $R(1 + D)$ . Specifically, denoting now the general class (1) by  $\hat{R}_g$ , the minimum first order approximation of  $\text{MSE}(\hat{R}_g)$  is

$$\begin{aligned} \text{MSE}^*(\hat{R}_g) &= \theta_2 R^2 S_D^2 \left( 1 - \frac{\theta}{\theta_2} \rho_{D,V}^2 - \frac{\theta_1}{\theta_2} \rho_{D,\bar{Z}}^2 \right) + \\ &= \text{MSE}'(\hat{R}) - \theta_2 R^2 S_D^2 \left( \frac{\theta}{\theta_2} \rho_{D,V}^2 + \frac{\theta_1}{\theta_2} \rho_{D,\bar{Z}}^2 \right), \end{aligned}$$

where  $\text{MSE}'(\hat{R}) = \theta_2 R^2 S_D^2$  is the first order approximation of the MSE of the standard estimator  $\hat{R}$ . For giving the chain regression type estimator,  $\hat{R}_{reg}$ , which reaches  $\text{MSE}^*(\hat{R}_g)$  we need to introduce more notation. Let  $U_1$  and  $U_2$  be any two auxiliary variables (where  $U_2$  can be a random vector) and let us denote the partial regression coefficient of  $Y_j$  on  $U_1$  given  $U_2$  as  $\beta_{j,U_1|U_2}$ ,  $j = 0, 1$ . If  $\hat{\beta}_{D,U_1|U_2}$  is a suitable estimate of  $\beta_{D,U_1|U_2} = \beta_{0,U_1|U_2}/\bar{Y}_0 - \beta_{1,U_1|U_2}/\bar{Y}_1$ , then the best estimator is given by

$$\begin{aligned} \hat{R}_{reg} &= \hat{R} \left[ 1 - \hat{\beta}_{D,X|\delta_X,\bar{Z}}(\bar{x} - \bar{x}') - \hat{\beta}_{D,\delta_X|X,\bar{Z}}(s_X^2 - s_X'^2) - \hat{\beta}_{D,Z|\bar{X},\delta_Z}(\bar{z} - \bar{z}') \right. \\ &\quad \left. - \hat{\beta}_{D,\delta_Z|\bar{X},Z}(s_Z^2 - s_Z'^2) - \hat{\beta}_{D,Z|\delta_Z}(\bar{z}' - \bar{Z}) - \hat{\beta}_{D,\delta_Z|Z}(s_Z'^2 - S_Z^2) \right]. \end{aligned}$$

If the interest is in the population product  $P = \bar{Y}_0 \bar{Y}_1$ , all the results for the ratio estimation can be used. It is enough to replace  $R$  and  $\hat{R}$  with  $P$  and  $\hat{P}$  respectively, where now  $\hat{P} = \bar{y}_0 \bar{y}_1$  and  $D = Y_0/\bar{Y}_0 + Y_1/\bar{Y}_1$ .

Finally when we want to estimate a population variance, i.e.  $T_Y = S_Y^2$ , we have that  $t_Y = s_Y^2$  is the sample mean of  $\delta_Y = (Y - \bar{Y})^2$ , except for terms of order  $n^{-1}$ . In fact

$$s_Y^2 = \frac{n}{n-1} \left[ \bar{s}_Y^2 - (\bar{y} - \bar{Y})^2 \right] \cong \bar{s}_Y^2,$$

where  $\bar{s}_Y^2 = \sum_{i=1}^n \delta_{Y_i}/n$ . Therefore replacing  $Y$  with  $\delta_Y$ , we can use again all the results of Section 4, getting

$$\text{MSE}^*(\hat{S}_Y^2) = \theta_2 S_{\delta_Y}^2 \left( 1 - \frac{\theta}{\theta_2} \rho_{\delta_Y,V}^2 - \frac{\theta_1}{\theta_2} \rho_{\delta_Y,\bar{Z}}^2 \right)$$

and

$$\begin{aligned} \hat{S}_{reg,Y}^2 &= s_Y^2 - \beta_{\delta_Y,X|\delta_X,\bar{Z}}(\bar{x} - \bar{x}') - \beta_{\delta_Y,\delta_X|X,\bar{Z}}(s_X^2 - s_X'^2) - \beta_{\delta_Y,Z|\bar{X},\delta_Z}(\bar{z} - \bar{z}') \\ &\quad - \beta_{\delta_Y,\delta_Z|\bar{X},Z}(s_Z^2 - s_Z'^2) - \beta_{\delta_Y,Z|\delta_Z}(\bar{z}' - \bar{Z}) - \beta_{\delta_Y,\delta_Z|Z}(s_Z'^2 - S_Z^2), \end{aligned}$$

where  $\hat{S}_Y^2$  denotes the general class (1) and  $\hat{S}_{reg,Y}^2$  is the best estimator in the class.

## 5 Estimating the population mean with independent samples

In this section we consider the estimation of the population mean when the first and second phase samples are independent.

Using the same notation given in Section 3 and after some algebra, from the general results of Section 2.2 follow that

$$\mathbf{g}_v^* = \left[ \begin{array}{c} -\frac{\theta_2}{\theta_{12}} S_{\bar{X}\bar{X}|\bar{Z}}^{-1} S_{\bar{X}Y|\bar{Z}} \\ -\frac{\theta_1}{\theta_{12}} S_{\bar{Z},\bar{Z}}^{-1} S_{\bar{Z},Y} - \frac{\theta_2}{\theta_{12}} S_{\bar{Z},\bar{Z}|\bar{X}}^{-1} S_{\bar{Z},Y|\bar{X}} \end{array} \right] \quad \text{and} \quad \mathbf{g}_z^* = -S_{\bar{Z},\bar{Z}}^{-1} S_{\bar{Z},Y},$$

where  $\theta_{12} = \theta_1 + \theta_2$ . In this case too, these optimum values depend on some coefficient regressions. For instance, the quantity  $S_{\bar{X}\bar{X}|\bar{Z}}^{-1} S_{\bar{X}Y|\bar{Z}}$  is the regression coefficient of  $Y$  on the residuals from the regression of  $\bar{X}$  on  $\bar{Z}$ . The minimum first order approximation for  $\text{MSE}(\bar{y}_g)$  is now

$$\text{MSE}^*(\bar{y}_g) = \theta_2 S_Y^2 \left( 1 - \frac{\theta_2}{\theta_{12}} \rho_{YV}^2 - \frac{\theta_1}{\theta_{12}} \rho_{Y\bar{Z}}^2 \right).$$

From the comparison between double sampling with the dependent and the independent samples, i.e. comparing the last expression for  $\text{MSE}^*(\bar{y}_g)$  with (12), we have that the sampling scheme with independent samples is always more efficient. It is more expensive, however, and thus it is less used in practice.

As for the case of dependent samples, the estimation of a ratio or product of two means and of the variance follow immediately from the mean case.

## 6 Conclusions

The main result of this paper is that, given an auxiliary vector  $\mathbf{t}$ , any estimator for a parameter  $T_Y$ , belongs to general class  $g(t_Y, \mathbf{t})$ . Furthermore we provide the minimum MSE for the class and a regression type estimator which reaches that minimum. No estimator based on the same  $\mathbf{t}$  can be more efficient than the best one. If an estimator is, however, equivalent to the best one, at the first order of approximation, it is optimum as well.

Among the parameters of interest a special role is played by the mean. The general results for this case let us to treat the estimation of a ratio or product of two means and the variance estimation. This is, actually, possible whenever the used estimator  $t_Y$  is a linear combination of samples means, at least at the first order of approximation.

Finally, we stress that in this paper some general results about the double sampling with independent samples are provided. This sampling scheme was not widely treated in literature, even if it gives more efficient estimators. From a practical point of view it is more expensive, however.

## 7 Appendix

In this appendix we provide another expression for  $g_V^*$  and  $g_Z^*$  given in Section 2.2. With these new expression it is easier to derive  $g_V^*$  and  $g_Z^*$  for the mean case.

$$g_V^* = \begin{bmatrix} g_X^* \\ g_Z^* \end{bmatrix} = \begin{bmatrix} -\left(S_{t_X, t_X | t_Z} - S_{t'_X, t'_X | t'_Z}\right)^{-1} S_{t_Y, t_X | t_Z} \\ -S_{t_Z, t_Z}^{-1} S_{t_Y, t_Z} - S_{t_Z, t_Z}^{-1} S_{t_X, t_Z}^T g_X^* \end{bmatrix}$$

$$g_Z^* = -S_{t_Z, t_Z}^{-1} S_{t_Y, t_Z} - S_{t_Z, t_Z}^{-1} S_{t_X, t_Z}^T g_X^* + S_{t'_Z, t'_Z}^{-1} S_{t'_X, t'_Z}^T g_X^*.$$

Specifically,  $g_X^*$  and  $g_Z^*$  for the mean case follow immediately from the above expression. While some algebra leads to the equivalence of the two expressions for  $g_Z^*$ .

### References

- Chand L. (1975), *Some Ratio-Type Estimators based on Two or More Auxiliary Variates*, (Unpublished Ph.D. Dissertation, Iowa State University, Ames, IOWA, USA).
- Diana G., Tommasi C. (2001), *Estimation for Finite Population Variance in Double Sampling*, Working Paper #2001.22, University of Padova, Padova, Italy.
- Diana G., Tommasi C. (2002a), *Optimal Estimation for Finite Population Mean in Two Phase Sampling*, Working Paper #2002.3, University of Padova, Padova, Italy.
- Diana G., Tommasi C. (2002b), *Estimation for Ratio of Two Population Means in Double Sampling*, Working Paper #2002.5, University of Padova, Padova, Italy.
- Gupta R.K., Singh S., Mangat N.S. (1992). Some chain ratio type estimators for estimating finite population variance, *Aligarh J. Statist.*, 12-13, 65-69.
- Khare B.B. (1991), Determination of Sample Sizes for a Class of Two-Phase Sampling Estimators for Ratio and Product of two Population Means Using Auxiliary Character, *Metron*, 49, 185-197.

- Prasad B., Singh R.S., Singh H.P. (1996), Some Chain Ratio-Type Estimators for Ratio of two Population Means Using Two Auxiliary Characters in Two Phase Sampling, *Metron*, 54, 95-113.
- Sahoo J., Sahoo L.N. (1999), An Alternative Class of Estimators in Double Sampling Procedures, *Cal. Statist. Assoc. Bull.*, 49, 193-194.
- Singh S. (1991), Estimation of finite population variance using double sampling, *Aligarh J. Statist.*, 11, 53-56.
- Singh H.P., Espejo Ruiz M. (2000), An Improved Class of Chain Regression Estimator in Two-Phase Sampling, *Statistics & Decisions*, 18, 205-218.
- Singh V.K., Singh G.N. (1991), Chain Type Regression Estimators with Two Auxiliary Variables under Double Sampling Scheme, *Metron*, 49, 279-289.
- Singh V.K., Singh H. P. (1994), Estimation of Ratio and Product of two Finite Population Means in Two Phase Sampling, *Journ. Statist. Plann. Infer.*, 41, 163-171.
- Singh V.K., Singh Hari P., Singh Housila P., Shukla D. (1994), A General Class of Chain Estimators for Ratio and Product of two Means in Finite Population, *Commun. Statist. Theory & Meth.*, 23(5), 1341-1355.
- Srivastava R., Khare B.B., Srivastava S.R. (1990), A Generalized Chain Ratio Estimator for Mean of Finite Population, *J. Indian Soc. Agric. Statist.*, 42, 108-117.
- Srivastava S. Rani, Khare B.B., Srivastava S.R. (1988), On Generalized Chain Estimators Ratio and Product of two Population Means Using Auxiliary Characters, *Assam. Stat. Review*, 2(1), 21-29.
- Srivastava S. Rani, Srivastava S.R., Khare B.B. (1989), Chain Ratio-Type Estimators for Ratio of two Population Means Using Auxiliary Characters, *Commun. Statist. Theory & Meth.*, 18(10), 3917-3926.
- Tracy D.S., Singh H.P. (1999), A General Class of Chain Regression Estimators in Two-Phase Sampling, *J. Appl. Statist. Sci.*, 8, 205-216.
- Upadhyaya L.N., Dubey S.P., Singh H.P. (1992), A Class of Ratio-in-Regression Estimators Using Two Auxiliary Variables in Double Sampling, *Journ. Scient. Res.*, 42, 127-134.

