

# VARIANCE ESTIMATION FOR THE TWO-PHASE REGRESSION ESTIMATOR - A CALIBRATION APPROACH

Martin Axelson, Statistics Sweden

Research and Development, Statistics Sweden, SE-701 89 Örebro, Sweden

**Key Words:** Two-phase sampling, Regression estimator, Variance estimation, Calibration

## 1 Introduction

In situations where weak or no auxiliary information is available at the population level, two-phase regression estimation constitutes a possible tool for cost-efficient estimation. In this respect, important references are Särndal and Swensson (1987), Dupont (1995), and Hidiroglou and Särndal (1998). Särndal and Swensson presented general results regarding generalized regression estimation under two-phase sampling, while both Dupont and Hidiroglou and Särndal discussed two-phase calibration estimation and its possible relation to generalized regression estimation. Hidiroglou and Särndal showed that under their suggested calibration approach, the generalized regression estimator (*GREG*) under two-phase sampling alternatively may be derived using a two-step calibration approach, a result in line with the findings regarding the relation between regression and calibration estimation under single-phase sampling in Deville and Särndal (1992).

In this paper we focus on variance estimation for the *GREG* under two-phase sampling, and the question we seek to answer is whether or not the available auxiliary information may be used more extensively than is generally the case, in order to obtain more efficient variance estimators. This is by no means a new question (e.g. Rao and Sitter, 1995; Sitter, 1997; Axelson, Breidt, and Carriquiry, 1996), but to our knowledge, the approach presented in this paper is a new development. The main goal of this paper is to extend the calibration technique to allow not only for point estimation purposes under two-phase sampling, but for variance estimation purposes as well. A general framework for calibration estimation of the variance of the two-phase regression estimator will be presented in detail and the method will be evaluated empirically through a small-scale simulation study. Related work under single-phase sampling includes Singh, Horn, and Yu (1998) and Théberge (1999).

The paper is organized as follows. In Section 2 some basic notation is introduced. In Section 3 the *GREG* is defined and a large-sample approximation

to its variance is given. Section 4 deals with variance estimation. The standard approach is presented in Section 4.1, while a detailed account of the proposed calibration approach is given in Section 4.2. In Section 5, finally, the results of a small-scale Monte Carlo study are presented.

## 2 Preliminaries

Let  $U = \{1, \dots, k, \dots, N\}$  denote the finite population of interest. Associated with each population element  $k \in U$  there are fixed values of the auxiliary column vector  $\mathbf{x}_1$  and the variable of interest  $y$ , respectively. The value of  $\mathbf{x}_1$ ,  $\mathbf{x}_{1k}$ , is known for all  $k \in U$ , and the objective of the survey is to estimate  $t_y = \sum_U y_k$ , the finite population total of  $y$ , which is unknown at the outset of the study. (If  $S_1$  is any set of elements such that  $S_1 \subseteq U$  and  $a_k$  is a quantity associated with element  $k$ ,  $\sum_{S_1} a_k$  is our shorthand for  $\sum_{k \in S_1} a_k$ .) To estimate  $t_y$ , two-phase regression estimation will be used.

Let  $s_a$  be a sample drawn from  $U$  according to a sampling design  $p_a(\cdot)$ . For all  $n_a$  elements included in the first-phase sample  $s_a$ , information on the auxiliary vector  $\mathbf{x}$  is recorded. (For a discussion about the possible relationships between the auxiliary variables  $\mathbf{x}_1$  and  $\mathbf{x}$ , see Hidiroglou and Särndal, 1998.) While  $\mathbf{x}_1$  is often primarily of an administrative nature,  $\mathbf{x}$  is typically chosen because it is assumed to be a powerful, yet relatively inexpensive, predictor for  $y$ . However, since both of the auxiliary variables  $\mathbf{x}_1$  and  $\mathbf{x}$  eventually will serve as predictors for the study variable  $y$  in the *GREG*, we assume that each variable in itself has a predictive ability strong enough to motivate the use of regression estimation. Next, a sample  $s$  is drawn from  $s_a$  according to a sampling design  $p(\cdot|s_a)$ , and the value of the study variable,  $y_k$ , is recorded for the  $n$  elements included in the second-phase sample. Throughout the paper, it is assumed that  $s_a$  and  $s$  are organized in such a way that the elements are sorted after increasing size of  $k$ . When necessary, the  $\nu_a$ th element in  $s_a$  will be referred to as  $k_{\nu_a}$  ( $\nu_a = 1, \dots, n_a$ ), while the  $\nu$ th element in  $s$  will be referred to as  $k_\nu^*$  ( $\nu = 1, \dots, n$ ).

The first-phase first- and second-order inclusion probabilities induced by  $p_a(\cdot)$  are denoted  $\pi_{ak}$  and  $\pi_{akl}$ , respectively, while the conditional first-

and second-order inclusion probabilities induced by  $p(\cdot|s_a)$  are denoted  $\pi_{k|s_a}$  and  $\pi_{kl|s_a}$ , respectively. In this paper, we only consider two-phase designs such that (i)  $\pi_{akl} > 0$  for all  $k&l \in U$  and (ii)  $\pi_{kl|s_a} > 0$  for all  $k&l \in s_a$  and every  $s_a$ .

For any quantity associated with element  $k$ , we let  $\tilde{\cdot}$  symbolize division by  $\pi_{ak}$  and let  $\check{\cdot}$  symbolize division by  $\pi_{ak}\pi_{k|s_a}$ . Thus, for example,  $\check{y}_k = y_k/\pi_{ak}$ , which is defined for all  $k \in U$ , and  $\check{\check{y}}_k = \check{y}_k/\pi_{k|s_a} = y_k/(\pi_{ak}\pi_{k|s_a})$ , which is defined for all  $k \in s_a$ .

### 3 Generalized regression estimation under two-phase sampling

Let  $\hat{t}_{y_s} = \sum_s \check{y}_k$  and let  $\hat{t}_{x_s}$  be analogous to  $\hat{t}_{y_s}$ . Moreover, let  $\hat{t}_{x_{1s_a}} = \sum_{s_a} \check{x}_{1k}$ , let  $\hat{t}_{x_{s_a}}$  be analogous to  $\hat{t}_{x_{1s_a}}$ , and let  $\mathbf{t}_{x_1} = \sum_U \mathbf{x}_{1k}$ . In line with Särndal and Swensson (1987) and Särndal et al. (1992, section 9.7), we define the *GREG* under two-phase sampling as

$$\hat{t}_{y_r} = \hat{t}_{y_s} + (\hat{t}_{x_{s_a}} - \hat{t}_{x_s})' \hat{\mathbf{B}}_s + (\mathbf{t}_{x_1} - \hat{t}_{x_{1s_a}})' \hat{\mathbf{B}}_{1s}, \quad (1)$$

where

$$\hat{\mathbf{B}}_s = (\sum_s \mathbf{x}_k \check{x}'_k / c_k)^{-1} \sum_s \mathbf{x}_k \check{y}_k / c_k \quad (2)$$

and

$$\hat{\mathbf{B}}_{1s} = (\sum_{s_a} \mathbf{x}_{1k} \check{x}'_{1k} / c_{1k})^{-1} \sum_{s_a} \mathbf{x}_{1k} \check{\check{y}}_k / c_{1k}, \quad (3)$$

with

$$\check{\check{y}}_k = \begin{cases} \check{x}'_k \hat{\mathbf{B}}_s & \text{if } k \in s_a - s. \\ \check{y}_k + (\check{x}_k - \check{\check{x}}_k)' \hat{\mathbf{B}}_s & \text{if } k \in s. \end{cases}$$

In (2),  $c$  is a weight which is assigned after the first phase of sampling but prior to the second phase of sampling. The weight serves to reflect the relative importance the statistician is willing to assign to element  $k$  on the basis of the auxiliary information available for  $k \in s_a$ . Similarly, in (3),  $c_1$  is a weight which serves to reflect the relative importance the statistician is willing to assign to element  $k$  on the basis of the auxiliary information available for  $k \in U$ .

Define

$$g_{1ks_a} = 1 + (\mathbf{t}_{x_1} - \hat{t}_{x_{1s_a}})' \times (\sum_{s_a} \mathbf{x}_{1k} \check{x}'_{1k} / c_{1k})^{-1} \mathbf{x}_{1k} / c_{1k}$$

for  $k \in s_a$ , and

$$g_{2ks} = 1 + (\sum_{s_a} g_{1ks_a} \check{x}_k - \sum_s g_{1ks_a} \check{\check{x}}_k)' \times (\sum_s \mathbf{x}_k \check{x}'_k / c_k)^{-1} \mathbf{x}_k / c_k$$

for  $k \in s$ , and let  $g_{ks} = g_{1ks_a} + g_{2ks} - 1$ , which thus is defined  $k \in s$ . It is matter of algebra to show that an alternative expression for (1) is given by

$$\hat{t}_{y_r} = \sum_s g_{ks} \check{y}_k.$$

This expression is due to Hidiroglou and Särndal (1998), who derived  $\hat{t}_{y_r}$  using a two-phase calibration approach.

The variance of  $\hat{t}_{y_r}$  may be written as  $V(\hat{t}_{y_r}) = V_1 + V_2$ , where  $V_1 = V_{pa}[E(\hat{t}_{y_r}|s_a)]$  and  $V_2 = E_{pa}[V(\hat{t}_{y_r}|s_a)]$  sometimes are referred to as the first- and second-phase variance component, respectively. Let

$$\mathbf{B}_{1U} = (\sum_U \mathbf{x}_{1k} \check{x}'_{1k} / c_{1k})^{-1} \sum_U \mathbf{x}_{1k} \check{y}_k / c_{1k},$$

let

$$\hat{\mathbf{B}}_{s_a} = (\sum_{s_a} \mathbf{x}_k \check{x}'_k / c_k)^{-1} \sum_{s_a} \mathbf{x}_k \check{y}_k / c_k,$$

and define the prediction errors  $E_{1kU} = y_k - \mathbf{x}'_{1k} \mathbf{B}_{1U}$  and  $E_{ks_a} = y_k - \mathbf{x}'_k \hat{\mathbf{B}}_{s_a}$ , which are defined for  $k \in U$  and  $k \in s_a$ , respectively. Moreover, let

$$\check{\mathbf{E}}_{1U} = (\check{E}_{11U}, \dots, \check{E}_{1kU}, \dots, \check{E}_{1NU})'$$

and

$$\check{\mathbf{E}}_{s_a} = (\check{E}_{k_1s_a}, \dots, \check{E}_{k_{\nu_a}s_a}, \dots, \check{E}_{k_{n_a}s_a})',$$

and define the matrices  $\Delta_{1U} = [\Delta_{akl}]$  and  $\Delta_{s_a} = [\Delta_{kl|s_a}]$ , where  $\Delta_{akl} = \pi_{akl} - \pi_{ak}\pi_{al}$  ( $k&l \in U$ ) and  $\Delta_{kl|s_a} = \pi_{kl|s_a} - \pi_{k|s_a}\pi_{l|s_a}$  ( $k&l \in s_a$ ). Large sample approximations for  $V_1$  and  $V_2$  are given by

$$V_1 \doteq AV_1 = \check{\mathbf{E}}'_{1U} \Delta_{1U} \check{\mathbf{E}}_{1U} \quad (4)$$

and

$$V_2 \doteq AV_2 = E_{pa}[AV(\hat{t}_{y_r}|s_a)], \quad (5)$$

where  $AV(\hat{t}_{y_r}|s_a) = \check{\mathbf{E}}'_{s_a} \Delta_{s_a} \check{\mathbf{E}}_{s_a}$ , respectively.

## 4 Variance estimation

### 4.1 The standard approach

For  $k \in s$ , let  $e_{1ks} = y_k - \mathbf{x}'_{1k} \hat{\mathbf{B}}_{1s}$  and  $e_{ks} = y_k - \mathbf{x}'_k \hat{\mathbf{B}}_s$ , and define the vectors

$$\check{\mathbf{e}}_{1s} = (\check{e}_{1k_1^*s}, \dots, \check{e}_{1k_{\nu_s}^*s}, \dots, \check{e}_{1k_{n_s}^*s})'$$

and

$$\check{\mathbf{e}}_s = (\check{e}_{k_1^*s}, \dots, \check{e}_{k_{\nu_s}^*s}, \dots, \check{e}_{k_{n_s}^*s})'.$$

Moreover, define the matrices  $\check{\check{\Delta}}_{1s} = [\check{\check{\Delta}}_{akl}]$  and  $\check{\check{\Delta}}_s = [\check{\check{\Delta}}_{kl|s_a}]$ , where  $\check{\check{\Delta}}_{akl} = \Delta_{akl}/(\pi_{akl}\pi_{kl|s_a})$

( $k \& l \in s$ ) and  $\check{\Delta}_{kl|s_a} = \Delta_{kl|s_a} / \pi_{kl|s_a}$  ( $k \& l \in s$ ). From (4) and (5) it follows that a possible estimator for  $V(\hat{t}_{yr})$  is given by

$$\hat{V}_{REF} = \hat{V}_{REF,1} + \hat{V}_{REF,2}, \quad (6)$$

where

$$\hat{V}_{REF,1} = \check{e}'_{1s} \check{\Delta}_{1s} \check{e}_{1s}$$

and

$$\hat{V}_{REF,2} = \check{e}'_s \check{\Delta}_s \check{e}_s$$

When the sample size is large in each of the two phases, (6) is approximately unbiased for  $V(\hat{t}_{yr})$  (e.g., Särndal and Swensson, 1987).

**Remark 1.** Let

$$\mathbf{G}_{1s} = \text{diag}(g_{1k_1^* s_a}, \dots, g_{1k_\nu^* s_a}, \dots, g_{1k_n^* s_a})$$

and

$$\mathbf{G}_s = \text{diag}(g_{k_1^* s}, \dots, g_{k_\nu^* s}, \dots, g_{k_n^* s}).$$

As an alternative to  $\hat{V}_{REF}$ , Axelson (2000a) proposed

$$\hat{V}_{REF}^{g1,g} = \hat{V}_{REF,1}^{g1} + \hat{V}_{REF,2}^g,$$

where

$$\hat{V}_{REF,1}^{g1} = \check{e}'_{1s} \mathbf{G}_{1s} \check{\Delta}_{1s} \mathbf{G}_{1s} \check{e}_{1s}$$

and

$$\hat{V}_{REF,2}^g = \check{e}'_s \mathbf{G}_s \check{\Delta}_s \mathbf{G}_s \check{e}_s.$$

Other alternatives to  $\hat{V}_{REF}$  based on the so-called  $g$ -weighted residual technique have been proposed by Särndal et al. (1992, section 9.7) and Hidiroglou and Särndal (1998).

## 4.2 The calibration approach

In this section we show how the results in Théberge (1999) regarding calibration estimation of bilinear estimators may be extended to allow for potentially efficient estimation of  $V(\hat{t}_{yr})$ . To this end, we need to introduce some new notation. Let  $\mathbf{z}_1$  ( $J_1 \times 1$ ) and  $\mathbf{z}$  ( $J \times 1$ ) denote vectors such that  $\mathbf{z}_{1k}$  is known for (i)  $k \in U$  or (ii)  $k \in s_a$ , and  $\mathbf{z}_k$  is known for  $k \in s_a$ , and define the matrices

$$\mathbf{Z}_{1U} = (\mathbf{z}_{11}, \dots, \mathbf{z}_{1k}, \dots, \mathbf{z}_{1N})',$$

$$\mathbf{Z}_{1s_a} = (\mathbf{z}_{1k_1}, \dots, \mathbf{z}_{1k_\nu}, \dots, \mathbf{z}_{1k_n})',$$

$$\mathbf{Z}_{1s} = (\mathbf{z}_{1k_1^*}, \dots, \mathbf{z}_{1k_\nu^*}, \dots, \mathbf{z}_{1k_n^*})',$$

$$\mathbf{Z}_{s_a} = (\mathbf{z}_{k_1}, \dots, \mathbf{z}_{k_\nu}, \dots, \mathbf{z}_{k_n})',$$

and

$$\mathbf{Z}_s = (\mathbf{z}_{k_1^*}, \dots, \mathbf{z}_{k_\nu^*}, \dots, \mathbf{z}_{k_n^*})'.$$

In deriving a calibration estimator for  $V(\hat{t}_{yr})$ , it is assumed that  $\mathbf{z}_1$  and  $\mathbf{z}$  may be regarded as auxiliary vectors for the residuals  $\check{e}_{1U}$  and  $\check{e}_{s_a}$ , respectively.

Let  $\mathbf{F}_1$  ( $A \times B$ ), let  $\text{vec}(\mathbf{F}_1)$  denote the vector obtained by stacking the successive columns of  $\mathbf{F}_1$  with the first column on top, and let  $\mathbf{F}_2$  ( $AB \times AB$ ) be a positive diagonal matrix. Following Théberge, we define the distance measure  $\|\mathbf{F}_1\|_{\mathbf{F}_2}^2 = \|\text{vec}(\mathbf{F}_1)\|_{\mathbf{F}_2}^2 = \text{vec}(\mathbf{F}_1)' \mathbf{F}_2 \text{vec}(\mathbf{F}_1)$ . Now, a calibration estimator for  $V_1$  based on  $\mathbf{Z}_{1U}$  is given by

$$\hat{V}_{CAL,1} = \check{e}'_{1s} \mathbf{W}_{1s} \check{e}_{1s}, \quad (7)$$

where  $\mathbf{W}_{1s}$  is the matrix which minimizes

$$\|\mathbf{W}_{1s}^* - \check{\Delta}_{1s}\|_{\mathbf{Q}_1}^2, \quad (8)$$

subject to the condition that it minimizes

$$\|\mathbf{Z}'_{1s} \mathbf{W}_{1s}^* \mathbf{Z}_{1s} - \mathbf{Z}'_{1U} \Delta_{1U} \mathbf{Z}_{1U}\|_{\mathbf{T}_1}^2. \quad (9)$$

The matrices  $\mathbf{Q}_1$  and  $\mathbf{T}_1$  in (8) and (9), respectively, are positive diagonal matrices, assigned by the statistician, which serve to reflect the relative importance of the units in the distance measures. When  $\mathbf{z}_{1k}$  is known only for  $k \in s_a$ , we simply minimize (8) with respect to  $\mathbf{W}_{1s}^*$ , subject to the condition that

$$\|\mathbf{Z}'_{1s} \mathbf{W}_{1s}^* \mathbf{Z}_{1s} - \mathbf{Z}'_{1s_a} \check{\Delta}_{1s_a} \mathbf{Z}_{1s_a}\|_{\mathbf{T}_1}^2, \quad (10)$$

where  $\check{\Delta}_{1s_a} = [\check{\Delta}_{akl}]$  with  $\check{\Delta}_{akl} = \Delta_{akl} / \pi_{akl}$  ( $k \& l \in s_a$ ), is minimized. A calibration estimator for  $V_2$  based on  $\mathbf{Z}_{s_a}$  is given by

$$\hat{V}_{CAL,2} = \check{e}'_s \mathbf{W}_s \check{e}_s, \quad (11)$$

where  $\mathbf{W}_s$  is the matrix which minimizes

$$\|\mathbf{W}_s^* - \check{\Delta}_s\|_{\mathbf{Q}}, \quad (12)$$

subject to the condition that it minimizes

$$\|\mathbf{Z}'_s \mathbf{W}_s^* \mathbf{Z}_s - \mathbf{Z}'_{s_a} \Delta_{s_a} \mathbf{Z}_{s_a}\|_{\mathbf{T}}^2. \quad (13)$$

The matrices  $\mathbf{Q}$  and  $\mathbf{T}$  in (12) and (13) are positive diagonal matrices, analogous to  $\mathbf{Q}_1$  and  $\mathbf{T}_1$ , respectively. Combining (7) and (11), we thus have

$$\hat{V}_{CAL} = \hat{V}_{CAL,1} + \hat{V}_{CAL,2}, \quad (14)$$

which may be viewed as a calibration estimator for  $V(\hat{t}_{yT})$ .

To find closed-form expressions for  $\mathbf{W}_{1s}$  and  $\mathbf{W}_s$ , respectively, let  $\mathbf{V}_{1U} = \mathbf{Z}_{1U} \otimes \mathbf{Z}_{1U}$ ,  $\mathbf{V}_{1s_a} = \mathbf{Z}_{1s_a} \otimes \mathbf{Z}_{1s_a}$ ,  $\mathbf{V}_{1s} = \mathbf{Z}_{1s} \otimes \mathbf{Z}_{1s}$ ,  $\mathbf{V}_{s_a} = \mathbf{Z}_{s_a} \otimes \mathbf{Z}_{s_a}$ , and  $\mathbf{V}_s = \mathbf{Z}_s \otimes \mathbf{Z}_s$ , where  $\otimes$  denotes the Kronecker product (e.g., Rao, 1973, p. 29). Moreover, for a matrix  $\mathbf{F}_1$ , let  $\mathbf{F}_1^\dagger$  denote the Moore-Penrose generalized inverse (e.g., Rao, 1973, p. 26) of  $\mathbf{F}_1$ . Using the results in Théberge (1999, section 5), it may be shown that the minimization of (8) with respect to  $\mathbf{W}_{1s}^*$ , subject to (9), yields

$$\begin{aligned} \text{vec}(\mathbf{W}_{1s}) &= \text{vec}(\check{\check{\Delta}}_{1s}) + \mathbf{Q}_1^{-1} \mathbf{V}_{1s} \mathbf{T}_1^{1/2} \\ &\quad \times [\mathbf{T}_1^{1/2} \mathbf{V}'_{1s} \mathbf{Q}_1^{-1} \mathbf{V}_{1s} \mathbf{T}_1^{1/2}]^\dagger \\ &\quad \times \mathbf{T}_1^{1/2} [\mathbf{V}'_{1U} \text{vec}(\Delta_{1U}) \\ &\quad - \mathbf{V}'_{1s} \text{vec}(\check{\check{\Delta}}_{1s})]. \end{aligned}$$

When (10) is used as the calibration condition, we get

$$\begin{aligned} \text{vec}(\mathbf{W}_{1s}) &= \text{vec}(\check{\check{\Delta}}_{1s}) + \mathbf{Q}_1^{-1} \mathbf{V}_{1s} \mathbf{T}_1^{1/2} \\ &\quad \times [\mathbf{T}_1^{1/2} \mathbf{V}'_{1s} \mathbf{Q}_1^{-1} \mathbf{V}_{1s} \mathbf{T}_1^{1/2}]^\dagger \\ &\quad \times \mathbf{T}_1^{1/2} [\mathbf{V}'_{1s_a} \text{vec}(\check{\check{\Delta}}_{1s_a}) \\ &\quad - \mathbf{V}'_{1s} \text{vec}(\check{\check{\Delta}}_{1s})]. \end{aligned}$$

Furthermore, the minimization of (12) with respect to  $\mathbf{W}_s^*$ , subject to (13), yields

$$\begin{aligned} \text{vec}(\mathbf{W}_s) &= \text{vec}(\check{\check{\Delta}}_s) + \mathbf{Q}^{-1} \mathbf{V}_s \mathbf{T}^{1/2} \\ &\quad \times [\mathbf{T}^{1/2} \mathbf{V}'_s \mathbf{Q}^{-1} \mathbf{V}_s \mathbf{T}^{1/2}]^\dagger \\ &\quad \times \mathbf{T}^{1/2} [\mathbf{V}'_{s_a} \text{vec}(\Delta_{s_a}) \\ &\quad - \mathbf{V}'_s \text{vec}(\check{\check{\Delta}}_s)]. \end{aligned}$$

**Remark 2.** An alternative to (14), based on  $g$ -weighted residuals, is given by

$$\hat{V}_{CAL}^{g1,s} = \hat{V}_{CAL,1}^{g1} + \hat{V}_{CAL,2}^g,$$

where

$$\hat{V}_{CAL,1}^{g1} = \check{\mathbf{e}}_{1s}' \mathbf{G}_{1s} \mathbf{W}_{1s} \mathbf{G}_{1s} \check{\mathbf{e}}_{1s}$$

and

$$\hat{V}_{CAL,2}^g = \check{\mathbf{e}}_s' \mathbf{G}_s \mathbf{W}_s \mathbf{G}_s \check{\mathbf{e}}_s.$$

**Remark 3.** Under stratified sampling designs, alternative calibration estimators are obtained if the minimization is carried out separately within each stratum.

## 5 A Simulation Study

To study the design-based properties of  $\hat{V}_{CAL}$  relative to  $\hat{V}_{REF}$ , a small-scale Monte Carlo study was performed. For the study, the population generated by Axelson (2000b) was used. The population  $U$  was generated according to a slightly modified version of the method suggested by Vale and Maurelli (1983), which allows for the generation a multivariate non-normal distribution with specified correlation structure and given marginal means, variances, and coefficients of skewness and kurtosis. In generating  $U$ , Axelson used the correlation structure and the univariate moments for the variables in the real-world population  $MU281$  (Särndal et al., 1992, Appendix B, pp. 652–659) as input. Hence, although artificial,  $U$  is similar to  $MU281$  in terms of the univariate moments as well as the correlation structure.

As an initial choice, the study variable and the auxiliary information was chosen in accordance to the choices made by Särndal et al. (1992, p. 278). That is,  $y$  corresponded to  $RMT85$ ,  $x_1$  corresponded to  $CS82$ , and  $x_2$  corresponded to  $SS82$ . The population correlation matrix of  $(y, x_1, x_2)'$  is given by

$$\mathbf{r} = (r_{ij}) = \begin{bmatrix} 1.00 & 0.65 & 0.66 \\ 0.65 & 1.00 & 0.14 \\ 0.66 & 0.14 & 1.00 \end{bmatrix}$$

and  $R^2 = 0.75$ , where  $R^2$  is the multiple coefficient of determination associated with regression of  $y$  on  $(1, x_1, x_2)$ . To get an indication of the extent to which the behavior of  $\hat{V}_{CAL}$  relative to  $\hat{V}_{REF}$  depends on the predictive ability of the vector  $(1, x_1, x_2)$ , the choice of  $x_2$  corresponding to  $REV84$  was also considered. For this choice, the population correlation matrix is given by

$$\mathbf{r} = (r_{ij}) = \begin{bmatrix} 1.00 & 0.65 & 0.91 \\ 0.65 & 1.00 & 0.59 \\ 0.91 & 0.59 & 1.00 \end{bmatrix}$$

and  $R^2 = 0.85$ . These two choices of auxiliary information will henceforth be referred to as  $X1$  and  $X2$ .

While the choice of the first-phase design often is governed by administrative arguments, the choice of the second-phase design is typically more directly related to efficiency arguments. It is not uncommon that the auxiliary information collected for the elements included in the first-phase sample is used in order to obtain an efficient stratified design for the second phase. To study the extent to which the behavior of  $\hat{V}_{CAL}$  relative to  $\hat{V}_{REF}$  depends on the choice of the second-phase design, the following two-phase sampling designs were considered:

D1  $p_a(\cdot)$ : simple random sampling ( $n_a/N = 0.1$ )  
 $p(\cdot|s_a)$ : simple random sampling ( $n/n_a = 0.1$ )

D2  $p_a(\cdot)$ : simple random sampling ( $n_a/N = 0.1$ )  
 $p(\cdot|s_a)$ :  $s_a$  divided into  $H = 2$  equally sized strata by increasing size of  $x_2$ , stratified simple random sampling ( $n_h/n_{ah} = 0.1$ )

D3  $p_a(\cdot)$ : simple random sampling,  $n_a/N = 0.1$   
 $p(\cdot|s_a)$ :  $s_a$  divided into  $H = 4$  equally sized strata by increasing size of  $x_2$ , stratified simple random sampling ( $n_h/n_{ah} = 0.1$ )

Now, let  $\delta_{hk} = 1$  if  $k$  belongs to second-phase strata  $h$  and define  $\delta_k = (\delta_{1k}, \dots, \delta_{hk}, \dots, \delta_{Hk})'$ . Throughout the simulation study,  $x_1$  was assumed available for all  $k \in U$  while  $x_2$  was known only for all  $k \in s_a$ , and the *GREG* was defined according to (1) with  $\mathbf{x}_{1k} = (1, x_{1k})'$ ,  $\mathbf{x}_k = \delta_k \otimes (1, x_{1k}, x_{2k})'$  and  $c_{1k} = c_k = 1$ .

It is likely that the performance of  $\hat{V}_{CAL}$  to a large extent is governed by the choice of the auxiliary vectors  $\mathbf{z}_1$  and  $\mathbf{z}$ . In the Monte Carlo study, the following two choices of  $\mathbf{z}_1$  and  $\mathbf{z}$  was considered:

Z1  $\mathbf{z}_{1k} = \check{\mathbf{x}}'_k \hat{\mathbf{B}}_s - \check{\mathbf{x}}'_{1k} \hat{\mathbf{B}}_{1s} = \check{d}_{ks}$  ( $k \in s_a$ ) and  $\mathbf{z}_k = -\check{d}_{ks}$  ( $k \in s_a$ )

Z2  $\mathbf{z}_{1k} = \check{\mathbf{x}}'_{1k} \hat{\mathbf{B}}_{1s}$  ( $k \in U$ ) and  $\mathbf{z}_k = \check{\mathbf{x}}'_k \hat{\mathbf{B}}_s$  ( $k \in s_a$ )

The motivation for Z1 is that  $\check{d}_{ks} = \check{e}_{1ks} - \check{e}_{ks}$  and  $-\check{d}_{ks} = \check{e}_{ks} - \check{e}_{1ks}$  may be regarded as proxies for  $\check{E}_{1ks_a}$  and  $\check{E}_{ks_a}$ , respectively. Z2 may be motivated through an extension of the framework presented by Singh et al. (1998), to allow for generalized regression estimation under two-phase sampling based on multivariate auxiliary information.

Let

$$\mathbf{\Pi}_{as} = \text{diag}(\pi_{ak_1^*}, \dots, \pi_{ak_v^*}, \dots, \pi_{ak_n^*})$$

and

$$\mathbf{\Pi}_{s|s_a} = \text{diag}(\pi_{k_1^*|s_a}, \dots, \pi_{k_v^*|s_a}, \dots, \pi_{k_n^*|s_a}).$$

The following two choices of  $\mathbf{Q}_1$  and  $\mathbf{Q}$  were considered:

Q1  $\mathbf{Q}_1 = (\mathbf{\Pi}_{s|s_a} \mathbf{\Pi}_{as}^{-1}) \otimes (\mathbf{\Pi}_{as}^{-1} \mathbf{\Pi}_{s|s_a})$  and  $\mathbf{Q} = (\mathbf{\Pi}_{s|s_a} \mathbf{\Pi}_{as})^{-1} \otimes (\mathbf{\Pi}_{as} \mathbf{\Pi}_{s|s_a})^{-1}$

Q2  $\mathbf{Q}_1 = \text{diag}(\mathbf{V}_{1s}) \text{diag}[\text{vec}(\check{\mathbf{\Delta}}_{1s})]^{-1}$  and  $\mathbf{Q} = \text{diag}(\mathbf{V}_s) \text{diag}[\text{vec}(\check{\mathbf{\Delta}}_s)]^{-1}$

Throughout the study, the matrices  $\mathbf{T}_1 = \mathbf{T} = \mathbf{1}$  were used. Hence, in addition to  $\hat{V}_{REF}$ , the following four choices of  $\hat{V}_{CAL}$  were included in the study:

$\hat{V}_{C_1} \hat{V}_{CAL}$  according to (14) using Z1 and Q1

$\hat{V}_{C_2} \hat{V}_{CAL}$  according to (14) using Z1 and Q2

$\hat{V}_{C_3} \hat{V}_{CAL}$  according to (14) using Z2 and Q1

$\hat{V}_{C_4} \hat{V}_{CAL}$  according to (14) using Z2 and Q2

For each of combination of sampling design and auxiliary information,  $M = 10000$  two-phase samples were realized, and for each simulation run, the variance estimates corresponding to  $\hat{V}_{REF}$ ,  $\hat{V}_{C_1}$ ,  $\hat{V}_{C_2}$ ,  $\hat{V}_{C_3}$ , and  $\hat{V}_{C_4}$  were computed. Let  $\hat{V}_m$  and  $\hat{t}_{yr,m}$  denote the estimates of  $V(\hat{t}_{yr})$  and  $t_y$ , respectively, corresponding to the  $m$ th simulation run. The Monte Carlo mean squared error of  $\hat{V}$  was calculated as

$$MSE_{MC}(\hat{V}) = \sum_{m=1}^M (\hat{V}_m - S_{\hat{t}_{yr}}^2)^2 / (M - 1),$$

where

$$S_{\hat{t}_{yr}}^2 = \sum_{m=1}^M (\hat{t}_{yr,m} - \bar{\hat{t}}_{yr})^2 / (M - 1),$$

with  $\bar{\hat{t}}_{yr} = \sum_{m=1}^M \hat{t}_{yr,m} / M$ . Also computed for each variance estimator were 95% confidence intervals of the form

$$CI(\hat{t}_{yr,m}, \hat{V}_m) = \hat{t}_{yr,m} \pm z_{0.975} \sqrt{\hat{V}_m},$$

where  $z_{0.975}$  is the 97.5 percentile of the standard normal distribution, and the empirical coverage rate was calculated as

$$EC_{MC}(\hat{t}_{yr}, \hat{V}) = 100 \sum_{m=1}^M I_{[CI(\hat{t}_{yr,m}, \hat{V}_m) \ni t_y]} / M,$$

where  $I_{[\cdot]}$  is the indicator function.

In Table 1, some of the results from the Monte Carlo study are presented.

Design & auxiliary information	$MSE_{MC}(\hat{V}) / MSE_{MC}(\hat{V}_{REF})$ ( $EC_{MC}(\hat{t}_{yr}, \hat{V})$ )				
	$\hat{V}_{REF}$	$\hat{V}_{C_1}$	$\hat{V}_{C_2}$	$\hat{V}_{C_3}$	$\hat{V}_{C_4}$
D1, X1	1.00 (93.7)	0.95 (93.7)	1.04 (94.0)	1.00 (93.7)	0.79 (94.0)
D2, X1	1.00 (93.8)	0.90 (93.8)	1.32 (94.1)	1.00 (93.8)	1.00 (94.1)
D3, X1	1.00 (93.1)	0.91 (93.2)	1.45 (93.8)	1.00 (93.2)	1.13 (93.6)
D1, X2	1.00 (94.0)	0.88 (94.1)	1.02 (94.2)	1.00 (94.0)	0.83 (94.4)
D2, X2	1.00 (93.6)	0.87 (93.7)	1.17 (93.8)	1.00 (93.6)	1.20 (94.2)
D3, X2	1.00 (93.3)	0.90 (93.4)	1.16 (94.0)	1.00 (93.3)	1.56 (94.0)

Table 1. Relative Monte Carlo MSEs and empirical coverage rates of 95% confidence intervals

Clearly, all of the studied variance estimators display acceptable properties in terms of the empirical coverage rates. Comparing the estimators on the basis of the difference between the the empirical and the nominal coverage rate, the data implies that either  $\hat{V}_{C_2}$  or  $\hat{V}_{C_4}$  is to be preferred. However, since the differences between the studied estimators are quite small, we refrain from drawing any far-reaching conclusions solely on the basis of the empirical coverage rates.

If the relative MSEs are included in the comparison a slightly different picture emerges, in which  $\hat{V}_{C_1}$  appears as the best choice. The major reasons for this are:

- i Apart for the combinations  $D1, X1$  and  $D1, X2$ ,  $\hat{V}_{C_1}$  is the most efficient estimator.
- ii Even under  $D1, X1$  and  $D1, X2$ ,  $\hat{V}_{C_1}$  is more efficient than  $\hat{V}_{REF}$ .
- iii Throughout the simulation study,  $\hat{V}_{C_2}$  is less efficient than  $\hat{V}_{REF}$ .
- iv  $\hat{V}_{C_3}$  displays the same efficiency as  $\hat{V}_{REF}$  throughout the simulation study. This may seem surprising, but it is a matter of algebra to show that  $\hat{V}_{C_3}$  is almost identical to  $\hat{V}_{REF}$  under the used definition of the *GREG*.
- v Although  $\hat{V}_{C_4}$ , which is most efficient under  $D1, X1$  and  $D1, X2$ , performs reasonably well under  $D2, X1$  and  $D3, X1$ , it is the most inefficient choice under  $D2, X2$  and  $D3, X2$ .
- vi The degree to which the performance depends on the choice of the auxiliary vector  $\mathbf{x}$  seems to be smaller for  $\hat{V}_{C_1}$  than for any of the other three calibration-type variance estimators.

In conclusion, the results from the Monte Carlo study imply that more efficient variance estimators may be obtained through the use of the suggested calibration approach, but they also indicate that the performance depends on the choice of auxiliary information to be used. However, before any general recommendations can be made, both the theoretical and practical properties of the suggested approach to variance estimation need to be further investigated.

## References

Axelson, M. (2000a). A modified approach to variance estimation for the generalized regression estimator under two-phase sampling. In *On variance estimation for the two-phase regression estimator*. Ph. D. thesis, Department of Statistics, Uppsala University.

Axelson, M. (2000b). A note on variance estimation for the regression estimator under two-phase sampling. In *On variance estimation for the two-phase regression estimator*. Ph. D. thesis, Department of Statistics, Uppsala University.

Axelson, M., Breidt, F. J., and Carriquiry, A. L. (1996). Two-phase regression estimation for policy analysis using computer simulation experiments. In *Proceedings of the Section on Survey Research Methods*, pp. 320–325. American Statistical Association.

Deville, J. C. and Särndal, C. E. (1992). Calibration estimators in survey sampling. *Journal of the American Statistical Association* **87**, 376–382.

Dupont, F. (1995). Alternative adjustments when there are several levels of auxiliary information. *Survey Methodology* **21**, 125–135.

Hidiroglou, M. A. and Särndal, C. E. (1998). Use of auxiliary information for two-phase sampling. *Survey Methodology* **24**, 11–20.

Rao, C. R. (1973). *Linear statistical inference and its applications* (2 ed.). New York: Wiley.

Rao, J. N. K. and Sitter, R. R. (1995). Variance estimation under two-phase sampling with application to imputation for missing data. *Biometrika* **82**, 453–460.

Särndal, C. E. and Swensson, B. (1987). A general view of estimation for two phases of selection with applications to two-phase sampling and non-response. *International Statistical Review* **55**, 279–294.

Särndal, C. E., Swensson, B., and Wretman, J. (1992). *Model Assisted Survey Sampling*. New York: Springer-Verlag.

Singh, S., Horn, S., and Yu, F. (1998). Estimation of variance of general regression estimator: Higher level calibration approach. *Survey Methodology* **24**, 41–50.

Sitter, R. R. (1997). Variance estimation for the regression estimator in two-phase sampling. *Journal of the American Statistical Association* **92**, 780–787.

Théberge, A. (1999). Extensions of calibration estimators in survey sampling. *Journal of the American Statistical Association* **94**, 635–644.

Vale, C. D. and Maurelli, V. A. (1983). Simulating multivariate nonnormal distributions. *Psychometrika* **48**, 465–471.