

Catalogue no. 12-001-X
ISSN 1492-0921

Survey Methodology

Variance estimation under monotone non-response for a panel survey

by H el ene Juillard and Guillaume Chauvet

Release date: December 20, 2018



Statistics
Canada

Statistique
Canada

Canada

How to obtain more information

For information about this product or the wide range of services and data available from Statistics Canada, visit our website, www.statcan.gc.ca.

You can also contact us by

Email at STATCAN.infostats-infostats.STATCAN@canada.ca

Telephone, from Monday to Friday, 8:30 a.m. to 4:30 p.m., at the following numbers:

- | | |
|---|----------------|
| • Statistical Information Service | 1-800-263-1136 |
| • National telecommunications device for the hearing impaired | 1-800-363-7629 |
| • Fax line | 1-514-283-9350 |

Depository Services Program

- | | |
|------------------|----------------|
| • Inquiries line | 1-800-635-7943 |
| • Fax line | 1-800-565-7757 |

Standards of service to the public

Statistics Canada is committed to serving its clients in a prompt, reliable and courteous manner. To this end, Statistics Canada has developed standards of service that its employees observe. To obtain a copy of these service standards, please contact Statistics Canada toll-free at 1-800-263-1136. The service standards are also published on www.statcan.gc.ca under “Contact us” > “Standards of service to the public.”

Note of appreciation

Canada owes the success of its statistical system to a long-standing partnership between Statistics Canada, the citizens of Canada, its businesses, governments and other institutions. Accurate and timely statistical information could not be produced without their continued co-operation and goodwill.

Standard table symbols

The following symbols are used in Statistics Canada publications:

- . not available for any reference period
- .. not available for a specific reference period
- ... not applicable
- 0 true zero or a value rounded to zero
- 0^s value rounded to 0 (zero) where there is a meaningful distinction between true zero and the value that was rounded
- ^P preliminary
- ^r revised
- X suppressed to meet the confidentiality requirements of the *Statistics Act*
- ^E use with caution
- F too unreliable to be published
- * significantly different from reference category ($p < 0.05$)

Published by authority of the Minister responsible for Statistics Canada

© Her Majesty the Queen in Right of Canada as represented by the Minister of Industry, 2018

All rights reserved. Use of this publication is governed by the Statistics Canada [Open Licence Agreement](#).

An HTML version is also available.

Cette publication est aussi disponible en français.

Variance estimation under monotone non-response for a panel survey

Hélène Juillard and Guillaume Chauvet¹

Abstract

Panel surveys are frequently used to measure the evolution of parameters over time. Panel samples may suffer from different types of unit non-response, which is currently handled by estimating the response probabilities and by reweighting respondents. In this work, we consider estimation and variance estimation under unit non-response for panel surveys. Extending the work by Kim and Kim (2007) for several times, we consider a propensity score adjusted estimator accounting for initial non-response and attrition, and propose a suitable variance estimator. It is then extended to cover most estimators encountered in surveys, including calibrated estimators, complex parameters and longitudinal estimators. The properties of the proposed variance estimator and of a simplified variance estimator are estimated through a simulation study. An illustration of the proposed methods on data from the ELFE survey is also presented.

Key Words: Longitudinal estimation; Non-response model; Product sampling design; Response homogeneity groups; Simplified variance estimation.

1 Introduction

Surveys are not only used to produce estimators for one point in time (cross-sectional estimations), but also to measure the evolution of parameters (longitudinal estimations), and are thus repeated over time. In this paper, we are interested in estimation and variance estimation for panel surveys, in which measures are repeated over time for units in a same sample (Kalton, 2009). Among the panel surveys (also known as longitudinal surveys, see Lynn, 2009), cohort surveys are particular cases where the units in the sample are linked by a common original event, such as being born on the same year for children in the ELFE survey (Enquête longitudinale française depuis l'enfance), which is the motivating example for this work.

ELFE is the first longitudinal study of its kind in France, tracking children from birth to adulthood (Pirus, Bois, Dufourg, Lanoë, Vandentorren, Leridon and the Elfe team, 2010). Covering the whole metropolitan France, it was launched in 2011 and consists of more than 18,000 children whose parents consented to their inclusion. It will examine every aspect of these children's lives from the perspectives of health, social sciences and environmental health. The ELFE survey suffers from unit non-response, which needs to be accounted for by using available auxiliary information, so as to limit the bias of estimators. Though the ELFE survey will be used for illustration in this paper, non-response occurs in virtually any panel survey so that the proposed methods are of general interest; see for example Laurie, Smith and Scott (1999) for the treatment of non-response of the British Household Panel Survey, or Vandecasteele and Debels (2007) for the European Community Household Panel.

Non-response is currently handled by modeling the response probabilities (Kim and Kim, 2007) and by reweighting respondents with the inverse of these estimated probabilities, which leads to the so-called

1. Hélène Juillard, INED, 133 boul. Davout, 75020 Paris, France; Guillaume Chauvet, ENSAI/IRMAR, Campus de Ker Lann, 35170 Bruz, France. E-mail: guillaume.chauvet@ensai.fr.

propensity score adjusted estimator. A panel sample may suffer from three types of unit non-response (Hawkes and Plewis, 2009): initial non-response refers to the original absence of selected units; wave non-response occurs when some units in the panel sample temporarily do not answer at some point in time, while attrition occurs when some units in the panel sample permanently do not answer from some point in time. Wave non-response was fairly uncommon in the first waves of the ELFE survey which were at our disposal. We therefore simplify this set-up by assuming monotone non-response, where only initial non-response and attrition occur.

There is a vast literature on the treatment of unit non-response for surveys over time, see Ekholm and Laaksonen (1991), Fuller, Loughin and Baker (1994), Rizzo, Kalton and Brick (1996), Clarke and Tate (2002), Laaksonen and Chambers (2006), Hawkes and Plewis (2009), Rendtel and Harms (2009), Laaksonen (2007), Slud and Bailey (2010), Zhou and Kim (2012). Variance estimation for longitudinal estimators is considered in Tam (1984), Laniel (1988), Nordberg (2000), Berger (2004), Skinner and Vieira (2005), Qualité and Tillé (2008) and Chauvet and Goga (2018), but with focus on the sampling variance only. Variance estimation in case of non-response weighting adjustments on cross-sectional surveys is considered in Kim and Kim (2007). To the best of our knowledge, and despite the interest for applications, variance estimation accounting for non-response for panel surveys has not been treated in the literature, with the exception of Zhou and Kim (2012).

Zhou and Kim (2012) consider the estimation of a mean for a panel survey, in case of monotone non-response. Instead of using the propensity score adjusted estimator, Zhou and Kim (2012) define an optimal propensity score estimator. It is obtained by noting that for any variable of interest observed before time t , the estimator produced at time t differs from the estimator obtained at the date when the variable was observed, which is based on a larger sample. Adjusting on these differences by means of some form of calibration leads to the estimator proposed by Zhou and Kim (2012). It makes full use of the information collected at previous times, and it is therefore expected to be more efficient than the propensity score adjusted estimator. However, a panel survey may include a large number of variables of interest observed at several times, and calibrating on a too large number of variables may lead to estimators whose performances are worsened (Silva and Skinner, 1997). A careful modeling exercise seems therefore necessary before applying the optimal estimator of Zhou and Kim (2012). In this work, we rather focus on the propensity score adjusted estimator, which is popular in practice.

Zhou and Kim (2012) also consider variance estimation for their optimal estimator, under the so-called reverse framework of Fay (1992). By viewing the sample obtained at time t as the result of a two-phase process, the first phase being associated to the original sampling design and the second phase to the successive non-response steps, it is assumed under the reverse framework that these two phases may be reversed. This requires the two-phase process to be strongly invariant as defined by Beaumont and Haziza (2016). In this paper, we propose a general variance estimator for the propensity score adjusted estimator, for which the strong invariance assumption is not needed. We also extend this variance estimator to account for estimation of complex parameters, possibly with calibrated weights, and to cover longitudinal estimators.

In each case, a simplified conservative variance estimator, which may be easier to compute for secondary users, is also proposed.

The paper is organized as follows. In Section 2, we first define the notation. A parametric model is then postulated, leading to estimated response probabilities and to a reweighted estimator. A variance estimator is then derived by following the approach in Kim and Kim (2007), and a simplified version is also proposed. They are illustrated in the particular case of the logistic regression model. The proposed variance estimator is extended to cover calibrated estimators and complex parameters in Section 3. Longitudinal estimation is discussed in Section 4, and the proposed variance estimator is used to cover such cases. The variance estimators are compared in Section 5 through a simulation study, and an illustration on the ELFE data is proposed in Section 6. We draw some conclusions in Section 7.

2 Correction of non-response and attrition

2.1 Notation and main assumptions

We are interested in a finite population U . A sample s_0 is first selected according to some sampling design $p(\cdot)$, and we assume that the first-order inclusion probabilities π_i are strictly positive for any $i \in U$. This first sampling phase corresponds to the original inclusion of units in the sample.

We consider the case of a panel survey in which the sole units in the original sample s_0 are followed over time, without reentry or late entry units at subsequent times to represent possible newborns. We are therefore interested in estimating some parameter defined over the population U , for some study variable y_t taking the value y_{it} for the unit i at time t . The units in the sample s_0 are followed at subsequent times $\delta = 1, \dots, t$, and the sample is prone to unit non-response at each time. We note r_i^δ for the response indicator for unit i at time δ , and s_δ for the subset of respondents at time δ .

We assume monotone non-response resulting in the nested sequence $s_0 \supset s_1 \supset \dots \supset s_t$. For $\delta = 1, \dots, t$, we note $p_i^\delta = \Pr(i \in s_\delta | s_{\delta-1})$ for the response probability of some unit i to be a respondent at time δ . We assume that the data are missing at random, i.e. the response probability p_i^δ at time δ can be explained by the variables observed at times $0, \dots, \delta - 1$, including the variables of interest, see for example Zhou and Kim (2012). Also, we assume that at any time δ the units answer independently of one another, and we note $p_{ij}^\delta = p_i^\delta p_j^\delta$ for the probability that two distinct units i and j answer jointly at time δ .

2.2 Reweighted estimator

We are interested in estimating the total $Y(t) = \sum_{i \in U} y_{it}$ at time t . In practice, the response probabilities at each time are unknown and need to be estimated. We assume that at each time δ the probability of response is parametrically modeled as

$$p_i^\delta = f^\delta(z_i^\delta, \alpha^\delta) \quad (2.1)$$

for some known function $f^\delta(\cdot, \cdot)$, where z_i^δ is a vector of variables observed for all the units in $s_{\delta-1}$, and α^δ denotes some unknown parameter. Here and elsewhere, the superscript δ will be used when we account for non-response at time δ , like for the probability p_i^δ of unit i to be a respondent at time δ . Following the approach in Kim and Kim (2007), we assume that the true parameter is estimated by $\hat{\alpha}^\delta$, the solution of the estimating equation

$$\frac{\partial}{\partial \alpha} \sum_{i \in s_{\delta-1}} k_i^\delta \{r_i^\delta \ln(p_i^\delta) + (1 - r_i^\delta) \ln(1 - p_i^\delta)\} = 0, \quad (2.2)$$

with k_i^δ some weight of unit i in the estimating equation. Customary choices for these weights include $k_i^\delta = 1$ and $k_i^\delta = \pi_i^{-1}$, see Fuller and An (1998), Beaumont (2005) and Kim and Kim (2007).

The estimated response probability at time δ is $\hat{p}_i^\delta = f^\delta(z_i^\delta, \hat{\alpha}^\delta)$. The propensity score adjusted estimator at time t , which will be simply called the reweighted estimator in what follows, is defined as

$$\hat{Y}_t(t) = \sum_{i \in s_t} \frac{y_{it}}{\pi_i \hat{p}_i^{1 \rightarrow t}} \quad \text{with} \quad \hat{p}_i^{1 \rightarrow t} = \prod_{\delta=1}^t \hat{p}_i^\delta. \quad (2.3)$$

Here and elsewhere, the subscript t will be used when the sample observed at time t is used for estimation, like for $\hat{Y}_t(\cdot)$ which makes use of the sample s_t . We simplify the notation as $\hat{Y}_t(t) \equiv \hat{Y}_t$ when the total at time t is estimated by using the sample observed at time t .

2.3 Variance computation

Under some regularity assumptions on the response mechanisms and some regularity conditions on the $p^\delta(\cdot, \cdot)$'s, we obtain from Theorem 1 in Kim and Kim (2007) that we can write

$$\hat{Y}_t = \hat{Y}_{\text{lin}, t}(t) + O_p(Nn^{-1}), \quad (2.4)$$

where

$$\hat{Y}_{\text{lin}, t}(t) = \sum_{i \in s_{t-1}} \frac{1}{\pi_i \hat{p}_i^{1 \rightarrow t-1}} \left\{ k_i^t \pi_i \hat{p}_i^{1 \rightarrow t-1} p_i^t (h_i^t)^\top \gamma^t + \frac{r_i^t}{p_i^t} \left(y_{it} - k_i^t \pi_i \hat{p}_i^{1 \rightarrow t-1} p_i^t (h_i^t)^\top \gamma^t \right) \right\}, \quad (2.5)$$

and where for any $\delta = 1, \dots, t$ we denote by h_i^δ the value of $h_i^\delta(\alpha) = \partial \logit(p_i^\delta) / \partial \alpha$ evaluated at $\alpha = \alpha^\delta$, and

$$\gamma^\delta = \left\{ \sum_{i \in s_{\delta-1}} k_i^\delta p_i^\delta (1 - p_i^\delta) h_i^\delta (h_i^\delta)^\top \right\}^{-1} \sum_{i \in s_{\delta-1}} \frac{1 - p_i^\delta}{\hat{p}_i^{1 \rightarrow \delta-1}} h_i^\delta \frac{y_{it}}{\pi_i}. \quad (2.6)$$

From (2.5), we obtain that

$$E\{\hat{Y}_{\text{lin}, t}(t) | s_{t-1}\} = \hat{Y}_{t-1}(t), \quad (2.7)$$

with $\hat{Y}_{t-1}(t)$ the estimator of $Y(t)$ computed on s_{t-1} . Using a proof by induction, it follows from (2.4) and (2.7) that \hat{Y}_t is approximately unbiased for $Y(t)$. Also, the variance of \hat{Y}_t may be asymptotically approximated by

$$V_{\text{app}}(\hat{Y}_t) = V\left(\sum_{i \in s_0} \frac{y_{it}}{\pi_i}\right) + E\left[\sum_{\delta=1}^t V\{\hat{Y}_{\text{lin},\delta}(t) | s_{\delta-1}\}\right]. \tag{2.8}$$

The first term in the right-hand side of (2.8) is the variance due to the sampling design, that we note as $V^p(\hat{Y}_t)$. The second term in the right-hand side of (2.8) is the variance due to non-response, that we note as $V^{\text{nr}}(\hat{Y}_t)$. From (2.5), this asymptotic variance is given by

$$V^{\text{nr}}(\hat{Y}_t) = E\left(\sum_{\delta=1}^t V^{\text{nr}\delta}(\hat{Y}_t)\right), \tag{2.9}$$

where

$$V^{\text{nr}\delta}(\hat{Y}_t) = \sum_{i \in s_{\delta-1}} p_i^\delta (1 - p_i^\delta) \left(\frac{y_{it}}{\pi_i \hat{p}_i^{1 \rightarrow \delta-1} p_i^\delta} - k_i^\delta (h_i^\delta)^\top \gamma^\delta \right)^2. \tag{2.10}$$

We note that for each of its component $\delta = 1, \dots, t$, the term $V^{\text{nr}\delta}(\hat{Y}_t)$ in (2.10) includes a centering term $k_i^\delta (h_i^\delta)^\top \gamma^\delta$, which is essentially a prediction of $(\pi_i \hat{p}_i^{1 \rightarrow \delta-1} p_i^\delta)^{-1} y_i$ by means of regressors h_i^δ . This centering is due to the estimation of the response probabilities. Suppressing these centering terms, equations (2.9) and (2.10) would lead to the variance of the estimator of $Y(t)$ we would obtain by replacing in (2.3) the estimated probabilities by their true values. The variance of this estimator is usually larger than that of the reweighted estimator in (2.3); see also Beaumont (2005), equation (5.7) and Kim and Kim (2007), equation (17), for the case $t = 1$.

2.4 Variance estimation

At time t , an approximately unbiased estimator for the variance due to the sampling design $V^p(\hat{Y}_t)$ is

$$\hat{V}_t^p(\hat{Y}_t) = \sum_{i,j \in s_t} \frac{\Delta_{ij}}{\pi_{ij}} \frac{1}{\hat{p}_{ij}^{1 \rightarrow t}} \frac{y_{it}}{\pi_i} \frac{y_{jt}}{\pi_j}, \tag{2.11}$$

where $\hat{p}_{ij}^{1 \rightarrow t} \equiv \prod_{\delta=1}^t \hat{p}_{ij}^\delta$, and where $\hat{p}_{ij}^\delta = \hat{p}_i^\delta$ if $i = j$, and $\hat{p}_{ij}^\delta = \hat{p}_i^\delta \hat{p}_j^\delta$ otherwise. Following equation (25) in Kim and Kim (2007), $V^{\text{nr}}(\hat{Y}_t)$ may be approximately unbiasedly estimated at time t by

$$\hat{V}_t^{\text{nr}}(\hat{Y}_t) = \sum_{\delta=1}^t \hat{V}_t^{\text{nr}\delta}(\hat{Y}_t) \tag{2.12}$$

where

$$\hat{V}_t^{\text{nr}\delta}(\hat{Y}_t) = \sum_{i \in s_t} \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \left(\frac{y_{it}}{\pi_i \hat{p}_i^{1 \rightarrow \delta}} - k_i^\delta (\hat{h}_i^\delta)^\top \hat{\gamma}_t^\delta \right)^2, \tag{2.13}$$

$$\hat{h}_i^\delta = h(z_i, \hat{\alpha}^\delta), \tag{2.14}$$

$$\hat{\gamma}_t^\delta = \left\{ \sum_{i \in s_t} k_i^\delta \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \hat{h}_i^\delta (\hat{h}_i^\delta)^\top \right\}^{-1} \sum_{i \in s_t} \frac{1 - \hat{p}_i^\delta}{\hat{p}_i^{1 \rightarrow t}} \hat{h}_i^\delta \frac{y_{it}}{\pi_i}. \quad (2.15)$$

This leads to the global variance estimator at time t

$$\hat{V}_t(\hat{Y}_t) = \hat{V}_t^p(\hat{Y}_t) + \hat{V}_t^{\text{nr}}(\hat{Y}_t). \quad (2.16)$$

A simplified estimator of the variance due to non-response is obtained by ignoring the prediction terms $k_i^\delta (\hat{h}_i^\delta)^\top \hat{\gamma}_t^\delta$ for each of the $\delta = 1, \dots, t$ variance components. After some algebra, this leads to the simplified variance estimator

$$\hat{V}_{t, \text{simp}}^{\text{nr}}\{\hat{Y}_t(t)\} = \sum_{i \in s_t} \frac{1 - \hat{p}_i^{1 \rightarrow t}}{(\hat{p}_i^{1 \rightarrow t})^2} \left(\frac{y_{it}}{\pi_i} \right)^2. \quad (2.17)$$

The main advantage of this simplified variance estimator is that it only requires the knowledge of the estimated response probabilities. On the other hand, the computation of the variance estimator in (2.12) requires the knowledge of the response models used at all times. The simplified variance estimator is therefore of particular interest for secondary users of the survey data, for which the estimated response probabilities may be the only available information related to the response modeling. This simplified variance estimator will tend to overestimate the variance due to non-response of (\hat{Y}_t) if the prediction term $k_i^\delta (\hat{h}_i^\delta)^\top \hat{\gamma}_t^\delta$ partly explains $(\pi_i \hat{p}_i^{1 \rightarrow \delta-1} p_i^\delta)^{-1} y_{it}$.

2.5 Application to the logistic regression model

In the particular case when a logistic regression model is used at each time δ , the model (2.1) may be rewritten as

$$\text{logit}(p_i^\delta) = (z_i^\delta)^\top \alpha^\delta. \quad (2.18)$$

We obtain $\hat{h}_i^\delta = z_i^\delta$, and the estimator for the variance due to non-response is given by (2.12), with

$$\hat{V}_t^{\text{nr}\delta}(\hat{Y}_t) = \sum_{i \in s_t} \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \left(\frac{y_{it}}{\pi_i \hat{p}_i^{1 \rightarrow \delta}} - k_i^\delta (z_i^\delta)^\top \hat{\gamma}_t^\delta \right)^2, \quad (2.19)$$

$$\hat{\gamma}_t^\delta = \left\{ \sum_{i \in s_t} k_i^\delta \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} z_i^\delta (z_i^\delta)^\top \right\}^{-1} \sum_{i \in s_t} \frac{1 - \hat{p}_i^\delta}{\hat{p}_i^{1 \rightarrow t}} z_i^\delta \frac{y_{it}}{\pi_i}. \quad (2.20)$$

If the reweighted estimator is computed at time $t = 1$, the estimator in (2.12) for the variance due to non-response may be rewritten as

$$\hat{V}_1^{\text{nr}}(\hat{Y}_1) = \sum_{i \in s_1} (1 - \hat{p}_i^1) \left(\frac{y_{i1}}{\pi_i \hat{p}_i^1} - k_i^1 (z_i^1)^\top \hat{\gamma}_1^1 \right)^2. \quad (2.21)$$

If the reweighted estimator is computed at time $t = 2$, the estimator in (2.12) for the variance due to non-response may be rewritten as

$$\begin{aligned} \hat{V}_2^{nr}(\hat{Y}_2) &= \sum_{i \in s_2} \frac{(1 - \hat{p}_i^1)}{\hat{p}_i^2} \left(\frac{y_{i2}}{\pi_i \hat{p}_i^1} - k_i^1 (z_i^1)^\top \hat{\gamma}_2^1 \right)^2 \\ &+ \sum_{i \in s_2} (1 - \hat{p}_i^2) \left(\frac{y_{i2}}{\pi_i \hat{p}_i^1 \hat{p}_i^2} - k_i^2 (z_i^2)^\top \hat{\gamma}_2^2 \right)^2. \end{aligned} \tag{2.22}$$

In practice, the model of Response Homogeneity Groups (RHG) is often assumed when correcting for unit non-response. Under this model, it is assumed that at each time $\delta = 1, \dots, t$, the sub-sample $s_{\delta-1}$ may be partitioned into $C(\delta - 1)$ groups $s_{\delta-1}^c$, $c = 1, \dots, C(\delta - 1)$, such that the response probability p_i^δ is constant inside a group. This model is a particular case of the logistic regression model in (2.18), obtained with

$$z_i^\delta = [1\{i \in s_{\delta-1}^1\}, \dots, 1\{i \in s_{\delta-1}^{C(\delta-1)}\}]^\top, \tag{2.23}$$

and the variance due to non-response is estimated accordingly. Explicit formulas are given in Appendix.

3 Calibration and complex parameters

In most surveys, a calibration step is used to obtain adjusted weights which enable to improve the accuracy of total estimates. Such calibrated estimators are considered in Section 3.1. Also, more complex parameters than totals are frequently of interest, and a linearization step can be used for variance estimation. This is the purpose of Section 3.2. The estimation of complex parameters with calibrated weights is treated in Section 3.3. In each case, explicit formulas for variance estimation and simplified variance estimation are derived, and the bias of the simplified variance estimator is discussed.

3.1 Variance estimation for calibrated total estimators

Assume that a vector x_i of auxiliary variables is available for any unit $i \in s_t$, and that the vector of totals X on the population U is known. Then an additional calibration step (Deville and Särndal, 1992) is usually applied to \hat{Y}_t . It consists in modifying the weights $d_{ii} = \pi_i^{-1} (\hat{p}_i^{1 \rightarrow t})^{-1}$ to obtain calibrated weights w_{ii} which enable to match the real total X , in the sense that

$$\sum_{i \in s_t} w_{ii} x_i = X. \tag{3.1}$$

The new calibrated weights are chosen to minimize a distance function with the original weights, while satisfying (3.1). This leads to the calibrated estimator

$$\hat{Y}_{wt} = \sum_{i \in s_t} w_{ii} y_{it}. \tag{3.2}$$

The estimated residual for the weighted regression of y_{it} on x_i is denoted by

$$e_{it} = y_{it} - \hat{b}_t x_i \quad (3.3)$$

with

$$\hat{b}_t = \left(\sum_{i \in s_t} \frac{1}{\pi_i \hat{p}_i^{1 \rightarrow t}} x_i x_i^\top \right)^{-1} \sum_{i \in s_t} \frac{1}{\pi_i \hat{p}_i^{1 \rightarrow t}} x_i y_{it}. \quad (3.4)$$

Replacing in (2.11) the variable y_{it} with e_{it} yields the estimator of the variance due to the sampling design

$$\hat{V}_t^p(\hat{Y}_{wt}) = \sum_{i, j \in s_t} \frac{\Delta_{ij}}{\pi_{ij}} \frac{1}{\hat{p}_{ij}^{1 \rightarrow t}} \frac{e_{it}}{\pi_i} \frac{e_{jt}}{\pi_j}. \quad (3.5)$$

Similarly, replacing in (2.12) the variable y_{it} with e_{it} yields the estimator of the variance due to the non-response

$$\hat{V}_t^{nr}(\hat{Y}_{wt}) = \sum_{\delta=1}^t \sum_{i \in s_t} \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \left(\frac{e_{it}}{\pi_i \hat{p}_i^{1 \rightarrow \delta}} - k_i^\delta (\hat{h}_i^\delta)^\top \hat{\gamma}_{te}^\delta \right)^2 \quad (3.6)$$

$$\hat{\gamma}_{te}^\delta = \left\{ \sum_{i \in s_t} k_i^\delta \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \hat{h}_i^\delta (\hat{h}_i^\delta)^\top \right\}^{-1} \sum_{i \in s_t} \frac{1 - \hat{p}_i^\delta}{\hat{p}_i^{1 \rightarrow t}} \hat{h}_i^\delta \frac{e_{it}}{\pi_i}. \quad (3.7)$$

The global variance estimator for \hat{Y}_{wt} is

$$\hat{V}_t(\hat{Y}_{wt}) = \hat{V}_t^p(\hat{Y}_{wt}) + \hat{V}_t^{nr}(\hat{Y}_{wt}). \quad (3.8)$$

The simplified estimator of the variance due to non-response is

$$\hat{V}_{t, \text{simp}}^{nr}(\hat{Y}_{wt}) = \sum_{i \in s_t} \frac{1 - \hat{p}_i^{1 \rightarrow t}}{(\hat{p}_i^{1 \rightarrow t})^2} \left(\frac{e_{it}}{\pi_i} \right)^2. \quad (3.9)$$

Here again, this simplified variance estimator ignores the prediction terms $k_i^\delta (\hat{h}_i^\delta)^\top \hat{\gamma}_{te}^\delta$. If the underlying calibration model is appropriate, then the explanatory power of \hat{h}_i^δ for e_{it} is expected to be small, as well as the bias of the simplified variance estimator. On the other hand, if there remains in e_{it} some significant part of y_{it} that may not be explained by x_i , the bias of the simplified variance estimator may be non-negligible. This may occur in case of domain estimation, when the calibration variables do not include any auxiliary information specific of the domain.

3.2 Variance estimation for complex parameters

We may be interested in estimating more complex parameters than totals. Suppose that the variable of interest y_{it} is q -multivariate, and that the parameter of interest is $\theta(t) = f\{Y(t)\}$ with $f(\cdot)$ a known function. At time t , substituting \hat{Y}_t into $\theta(t)$ yields the plug-in estimator $\hat{\theta}_t = f(\hat{Y}_t)$.

The estimated linearized variable of $\theta(t)$ is

$$u_{it} = \left\{ f'(\hat{Y}_t) \right\}^\top y_{it}, \quad (3.10)$$

with $f'(\hat{Y}_t)$ the q -vector of first derivatives of f at point \hat{Y}_t . Replacing in (2.11) the variable y_{it} with u_{it} yields the estimator of the variance due to the sampling design

$$\hat{V}_t^p(\hat{\theta}_t) = \sum_{i,j \in s_t} \frac{\Delta_{ij}}{\pi_{ij}} \frac{1}{\hat{p}_{ij}^{1 \rightarrow t}} \frac{u_{it}}{\pi_i} \frac{u_{jt}}{\pi_j}. \tag{3.11}$$

Similarly, replacing in (2.12) the variable y_{it} with u_{it} yields the estimator of the variance due to the non-response

$$\hat{V}_t^{nr}(\hat{\theta}_t) = \sum_{\delta=1}^t \sum_{i \in s_t} \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \left(\frac{u_{it}}{\pi_i \hat{p}_i^{1 \rightarrow \delta}} - k_i^\delta (\hat{h}_i^\delta)^\top \hat{\gamma}_{t\theta}^\delta \right)^2 \tag{3.12}$$

$$\hat{\gamma}_{t\theta}^\delta = \left\{ \sum_{i \in s_t} k_i^\delta \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \hat{h}_i^\delta (\hat{h}_i^\delta)^\top \right\}^{-1} \sum_{i \in s_t} \frac{1 - \hat{p}_i^\delta}{\hat{p}_i^{1 \rightarrow t}} \hat{h}_i^\delta \frac{u_{it}}{\pi_i}. \tag{3.13}$$

The global variance estimator for $\hat{\theta}_t$ is

$$\hat{V}_t(\hat{\theta}_t) = \hat{V}_t^p(\hat{\theta}_t) + \hat{V}_t^{nr}(\hat{\theta}_t). \tag{3.14}$$

The simplified estimator of the variance due to non-response is

$$\hat{V}_{t, \text{simpl}}^{nr}(\hat{\theta}_t) = \sum_{i \in s_t} \frac{1 - \hat{p}_i^{1 \rightarrow t}}{(\hat{p}_i^{1 \rightarrow t})^2} \left(\frac{u_{it}}{\pi_i} \right)^2. \tag{3.15}$$

The bias of this simplified variance estimator will depend on the explanatory power for \hat{h}_i^δ on the linearized variable u_{it} .

3.3 Variance estimation for complex parameters under calibration

The calibrated weights w_{it} may be used to obtain an estimator of the parameter $\theta(t)$. Substituting \hat{Y}_{wt} into $\theta(t) = f\{Y(t)\}$ yields the calibrated plug-in estimator $\hat{\theta}_{wt} = f(\hat{Y}_{wt})$. To obtain a variance estimator for $\hat{\theta}_{wt}$, we first compute the estimated linearized variable $u_{it} = \{f'(\hat{Y}_t)\}^\top y_{it}$ and take

$$e_{\theta it} = u_{it} - \hat{b}_{\theta t} x_i \tag{3.16}$$

with

$$\hat{b}_{\theta t} = \left(\sum_{i \in s_t} \frac{1}{\pi_i \hat{p}_i^{1 \rightarrow t}} x_i x_i^\top \right)^{-1} \sum_{i \in s_t} \frac{1}{\pi_i \hat{p}_i^{1 \rightarrow t}} x_i u_{it}. \tag{3.17}$$

Replacing in (2.11) the variable y_{it} with $e_{\theta it}$ yields the estimator of the variance due to the sampling design

$$\hat{V}_t^p(\hat{\theta}_{wt}) = \sum_{i,j \in s_t} \frac{\Delta_{ij}}{\pi_{ij}} \frac{1}{\hat{p}_{ij}^{1 \rightarrow t}} \frac{e_{\theta it}}{\pi_i} \frac{e_{\theta jt}}{\pi_j}. \tag{3.18}$$

Similarly, replacing in (2.12) the variable y_{it} with $e_{\theta it}$ yields the estimator of the variance due to the non-response

$$\hat{V}_t^{\text{nr}}(\hat{\theta}_{wt}) = \sum_{\delta=1}^t \sum_{i \in s_t} \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \left(\frac{e_{\theta it}}{\pi_i \hat{p}_i^{1 \rightarrow \delta}} - k_i^\delta (\hat{h}_i^\delta)^\top \hat{\gamma}_{te\theta}^\delta \right)^2 \quad (3.19)$$

$$\hat{\gamma}_{te\theta}^\delta = \left\{ \sum_{i \in s_t} k_i^\delta \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \hat{h}_i^\delta (\hat{h}_i^\delta)^\top \right\}^{-1} \sum_{i \in s_t} \frac{1 - \hat{p}_i^\delta}{\hat{p}_i^{1 \rightarrow t}} \hat{h}_i^\delta \frac{e_{\theta it}}{\pi_i}. \quad (3.20)$$

The global variance estimator for $\hat{\theta}_{wt}$ is

$$\hat{V}_t(\hat{\theta}_{wt}) = \hat{V}_t^p(\hat{\theta}_{wt}) + \hat{V}_t^{\text{nr}}(\hat{\theta}_{wt}). \quad (3.21)$$

The simplified estimator of the variance due to non-response is

$$\hat{V}_{t, \text{simp}}^{\text{nr}}(\hat{\theta}_{wt}) = \sum_{i \in s_t} \frac{1 - \hat{p}_i^{1 \rightarrow t}}{(\hat{p}_i^{1 \rightarrow t})^2} \left(\frac{e_{\theta it}}{\pi_i} \right)^2. \quad (3.22)$$

Since the variable $e_{\theta it}$ is obtained as the residual in the regression of the linearized variable u_{it} on the calibration variables x_i , the explanatory power for \hat{h}_i^δ on $e_{\theta it}$ is expected to be small in practice, and the bias of the simplified variance estimator is expected to be small as well.

4 Longitudinal estimators

We may be interested in a change in parameters, such as

$$\Delta(u \rightarrow t) = Y(t) - Y(u), \quad (4.1)$$

the difference between the totals of a variable of interest measured at two different times $u < t$. Since the variable y_{iu} is measured on all sub-samples $s_{u'}$ for $u' = u, \dots, t$, there are several possible estimators for $\Delta(u \rightarrow t)$. For $u' = u, \dots, t$, we denote by

$$\hat{\Delta}_{u't}(u \rightarrow t) = \sum_{i \in s_t} \frac{y_{it}}{\pi_i \hat{p}_i^{1 \rightarrow t}} - \sum_{i \in s_{u'}} \frac{y_{iu}}{\pi_i \hat{p}_i^{1 \rightarrow u'}} \quad (4.2)$$

the estimator which makes use of s_t for the estimation of $Y(t)$, and of $s_{u'}$ for the estimation of $Y(u)$. The case $u' = u$ corresponds to the estimation of $Y(u)$ on the largest available sub-sample, s_u . The case $u' = t$ corresponds to the estimation of $Y(u)$ and $Y(t)$ on the common sub-sample s_t .

In the context of full response, several authors have recommended the estimator $\hat{\Delta}_{u't}(u \rightarrow t)$ which makes use of the common sample only, if the variables y_{ui} and y_{it} are strongly positively correlated; see Caron and Ravalet (2000), Qualité and Tillé (2008), Goga, Deville and Ruiz-Gazen (2009), Chauvet and Goga (2018). In our context, this choice may be heuristically justified as follows. For $u' < t$, and by conditioning on the sub-sample $s_{u'}$, we obtain

$$V\left\{\hat{\Delta}_{u't}(u \rightarrow t)\right\} \simeq V\left\{\sum_{i \in s_{u'}} \frac{y_{it} - y_{iu}}{\pi_i \hat{p}_i^{1 \rightarrow u'}}\right\} + EV\left\{\sum_{i \in s_t} \frac{y_{it}}{\pi_i \hat{p}_i^{1 \rightarrow t}} \middle| s_{u'}\right\}, \quad (4.3)$$

$$V \left\{ \hat{\Delta}_u(u \rightarrow t) \right\} \simeq V \left\{ \sum_{i \in s_u} \frac{y_{it} - y_{iu}}{\pi_i \hat{p}_i^{1 \rightarrow u'}} \right\} + EV \left\{ \sum_{i \in s_t} \frac{y_{it} - y_{iu}}{\pi_i \hat{p}_i^{1 \rightarrow t}} \middle| s_u \right\}. \tag{4.4}$$

In equations (4.3) and (4.4), the first term in the right-hand side is identical. Since the variables y_{iu} and y_{it} are expected to be positively correlated, the difference $y_{it} - y_{iu}$ is expected to be smaller than y_{it} . Therefore, the estimator $\hat{\Delta}_u(u \rightarrow t)$ based on the common sample is expected to be more efficient in terms of variance. The results of a small simulation study in Section 5.2 support this heuristic reasoning. Therefore, we focus only in this Section on the estimator $\hat{\Delta}_u(u \rightarrow t)$ for the estimation of $\Delta(u \rightarrow t)$. As pointed out by a Referee, and following the approach in Zhou and Kim (2012), we may obtain a gain in efficiency by using the full information on s_u , namely by calibrating the weights $(\pi_i \hat{p}_i^{1 \rightarrow t})^{-1}$ on the estimator \hat{Y}_u .

Replacing in (2.11) the variable y_{it} with $y_{it} - y_{iu}$ yields the estimator of the variance due to the sampling design

$$\hat{V}_t^p \left\{ \hat{\Delta}_u(u \rightarrow t) \right\} = \sum_{i, j \in s_t} \frac{\Delta_{ij}}{\pi_{ij} \hat{p}_{ij}^{1 \rightarrow t}} \frac{1}{\pi_i} \frac{(y_{it} - y_{iu})(y_{jt} - y_{ju})}{\pi_j}. \tag{4.5}$$

Similarly, replacing in (2.12) the variable y_{it} with $y_{it} - y_{iu}$ yields the estimator of the variance due to the non-response

$$\hat{V}_t^{nr} \left\{ \hat{\Delta}_u(u \rightarrow t) \right\} = \sum_{\delta=1}^t \sum_{i \in s_t} \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \left(\frac{y_{it} - y_{iu}}{\pi_i \hat{p}_i^{1 \rightarrow \delta}} - k_i^\delta (\hat{h}_i^\delta)^\top \hat{\gamma}_{t\Delta}^\delta \right)^2 \tag{4.6}$$

with

$$\hat{\gamma}_{t\Delta}^\delta = \left\{ \sum_{i \in s_t} k_i^\delta \frac{\hat{p}_i^\delta (1 - \hat{p}_i^\delta)}{\hat{p}_i^{\delta \rightarrow t}} \hat{h}_i^\delta (\hat{h}_i^\delta)^\top \right\}^{-1} \sum_{i \in s_t} \frac{1 - \hat{p}_i^\delta}{\hat{p}_i^{1 \rightarrow t}} \hat{h}_i^\delta \frac{y_{it} - y_{iu}}{\pi_i}. \tag{4.7}$$

The global variance estimator for $\hat{\Delta}_u(u \rightarrow t)$ is

$$\hat{V}_t \left\{ \hat{\Delta}_u(u \rightarrow t) \right\} = \hat{V}_t^p \left\{ \hat{\Delta}_u(u \rightarrow t) \right\} + \hat{V}_t^{nr} \left\{ \hat{\Delta}_u(u \rightarrow t) \right\}. \tag{4.8}$$

Variance estimation for measures of change is also considered in Berger (2004), Qualité and Tillé (2008), Goga et al. (2009), Chauvet and Goga (2018), among others.

The simplified estimator of the variance due to non-response is

$$\hat{V}_{t, \text{simp}}^{nr} \left\{ \hat{\Delta}_u(u \rightarrow t) \right\} = \sum_{i \in s_t} \frac{1 - \hat{p}_i^{1 \rightarrow t}}{(\hat{p}_i^{1 \rightarrow t})^2} \left(\frac{y_{it} - y_{iu}}{\pi_i} \right)^2. \tag{4.9}$$

If the variables y_{it} and y_{iu} are strongly positively correlated, the bias of the simplified variance estimator is expected to be small.

5 A simulation study

In this section, several artificial populations are generated according to the model described in Section 5.1. In Section 5.2, we consider several estimators for a change between totals, which illustrates the heuristic reasoning in Section 4. A Monte Carlo experiment is presented in Section 5.3, and several variance estimators for estimating a total, a ratio or a parameter change are compared. The results from Tables 5.1 and 5.2 are readily reproducible using the R code provided in the Supplementary Material.

5.1 Simulation set-up

We consider seven populations of size 10,000, each containing three variables of interest y_{i1} , y_{i2} and y_{i3} observed at times $t = 1, 2$ and 3 , respectively. The variables of interest are generated according to the superpopulation model

$$y_{i1} = \alpha^0 + \alpha^a x_{ai} + \alpha^b x_{bi} + \sigma u_{i1}, \quad (5.1)$$

$$y_{i2} = \rho y_{i1} + \sigma u_{i2}, \quad (5.2)$$

$$y_{i3} = \rho y_{i2} + \sigma u_{i3}. \quad (5.3)$$

The auxiliary variables x_{ai} and x_{bi} are independently generated from a Gamma distribution with shape and scale parameters 2 and 1. Two auxiliary variables x_{ci} and x_{di} , not related to the variables of interest, are generated similarly. The variables u_{i1} , u_{i2} and u_{i3} are independently generated according to a standard normal distribution. We use $\alpha^0 = 10$, $\alpha^a = \alpha^b = 5$ and $\sigma = 10$, which leads to a coefficient of determination (R^2) in model (5.1) approximately equal to 0.50. The parameter ρ is set to 0, 0.2, 0.4, 0.6, 0.8, 1.0 and 1.2 for populations 1 to 7, respectively.

For each population, a simple random sample s_0 of size $n = 1,000$ is selected. Three non-response phases are then successively simulated. At each phase $\delta = 1, 2, 3$, the sub-sample of respondents s_δ is obtained by Poisson sampling with a response probability p_i^δ for unit i , defined as

$$\text{logit}(p_i^\delta) = \beta^{\delta 0} + \beta^{\delta a} x_{ai} + \beta^{\delta b} x_{bi}. \quad (5.4)$$

We use $\beta^{\delta 0} = -1$ at each phase $\delta = 1, 2, 3$. For $\delta = 1$, we use $\beta^{1a} = \beta^{1b} = 0.60$, which corresponds to an average response rate of 0.75. For $\delta = 2, 3$, we use $\beta^{\delta a} = \beta^{\delta b} = 0.75$, which corresponds to an average response rate of 0.81. Inside each sub-sample s_δ , the estimated response probabilities \hat{p}_i^δ are obtained by means of an unweighted logistic regression.

5.2 Comparison of estimators for a difference of totals

In this section, we are interested in comparing the accuracy of two estimators for a difference of totals $\Delta(u \rightarrow t)$ for $u = 1$ and $t = 2$, for $u = 1$ and $t = 3$, and for $u = 2$ and $t = 3$. We consider the estimator $\hat{\Delta}_{ut}(u \rightarrow t)$, which makes use of the whole appropriate sub-samples for variables y_{iu} and y_{it} ,

and the estimator $\hat{\Delta}_u(u \rightarrow t)$, which makes use of the common sub-sample only. These two estimators are compared through the relative difference (RD) of their variances, which are defined as follows:

$$RD(u \rightarrow t) = 100 \times \frac{V\{\hat{\Delta}_{ut}(u \rightarrow t)\} - V\{\hat{\Delta}_u(u \rightarrow t)\}}{V\{\hat{\Delta}_u(u \rightarrow t)\}}. \tag{5.5}$$

The true variances are replaced by their Monte Carlo approximation, obtained by repeating $B = 100,000$ times the sample selection and the non-response phases.

The results are presented in Table 5.1. A positive RD indicates that the use of the common sample only leads to a more accurate estimator. As could be expected, the RD increases in all cases with ρ , that is, when the correlation between y_{it} and y_{iu} increases. For $u = 1$ and $t = 2$, and for $u = 2$ and $t = 3$, the estimator $\hat{\Delta}_u(u \rightarrow t)$ is more accurate for ρ greater than 0.6. For $u = 1$ and $t = 3$, $\hat{\Delta}_u(u \rightarrow t)$ is more accurate for ρ greater than 0.8.

Table 5.1
Relative Difference (RD) between two estimators for a difference of totals

ρ	RD(1 → 2)	RD(1 → 3)	RD(2 → 3)
0.0	-12	-27	-13
0.2	-09	-25	-11
0.4	-04	-20	-03
0.6	05	-09	11
0.8	17	11	39
1.0	30	33	83
1.2	40	46	127

5.3 Performances of the variance estimators

In this section, we consider the artificial population 5 ($\rho = 0.8$) generated as described in Section 5.1. The sample selection by means of simple random sampling of size $n = 1,000$ and the three non-response phases are applied $B = 5,000$ times. We are interested in evaluating the variance estimators and the simplified variance estimators, in case of estimating a total, a ratio or a change in totals.

As for the total $Y(t)$, we consider at each time $t = 1, 2, 3$, three estimators. The estimator \hat{Y}_t makes use of the weights $d_{ii} = \pi_i^{-1}(\hat{p}_i^{1 \rightarrow t})^{-1}$. The estimator \hat{Y}_{wt} makes use of the weights w_i , obtained by calibrating the weights d_{ii} on the population size and on the totals of the auxiliary variables x_{ai} and x_{bi} . The estimator $\hat{Y}_{\tilde{w}t}$ makes use of the weights \tilde{w}_i , obtained by calibrating the weights d_{ii} on the population size and on the totals of the auxiliary variables x_{ci} and x_{di} . The working model is therefore well-specified for \hat{Y}_{wt} , but not for $\hat{Y}_{\tilde{w}t}$. The proposed variance estimator for \hat{Y}_t is obtained from equation (2.16), and the simplified variance estimator is obtained by plugging in (2.16) the simplified variance estimator for non-response given in (2.17). The proposed variance estimators for \hat{Y}_{wt} and $\hat{Y}_{\tilde{w}t}$ are obtained from equation (3.8), and the simplified variance estimators are obtained by plugging in (3.8) the simplified variance estimator for non-response given in (3.9).

We are also interested in estimating the ratio $R(t) = Y(t)/Y(1)$ for $t = 2, 3$. At each time t , we consider three estimators. The estimator \hat{R}_t makes use of the weights d_i . The proposed variance estimator is obtained from equation (3.14), by using the estimated linearized variable $u_{it} = (\hat{Y}_1)^{-1} (y_{it} - \hat{R}_t y_{1i})$. The simplified variance estimator is obtained by plugging in (3.14) the simplified variance estimator for non-response given in (3.15). The estimators \hat{R}_{w_t} and $\hat{R}_{\tilde{w}_t}$ make use of the calibrated weights w_i and \tilde{w}_i . The proposed variance estimators are obtained from equation (3.21). The simplified variance estimators are obtained by plugging in (3.21) the simplified variance estimator for non-response given in (3.22).

Finally, we are interested in estimating the change in totals $\Delta(1 \rightarrow t)$ for $t = 2, 3$. At each time t , we consider three estimators. The estimator $\hat{\Delta}_t(1 \rightarrow t)$ makes use of the weights d_i . The proposed variance estimator is obtained from equation (4.8), and the simplified variance estimator is obtained by plugging in (4.8) the simplified variance estimator for non-response given in (4.9). The estimators $\hat{\Delta}_{t,w}(1 \rightarrow t)$ and $\hat{\Delta}_{t,\tilde{w}}(1 \rightarrow t)$ make use of the calibrated weights w_i and \tilde{w}_i . The proposed variance estimators are obtained from equation (4.8), by replacing $y_{it} - y_{1i}$ by the estimated residual for the weighted regression of $y_{it} - y_{1i}$ on the calibration variables. The simplified variance estimators are obtained by plugging in (4.8) the simplified variance estimator for non-response given in (4.9).

For a proposed variance estimator \hat{V} , we computed the Monte Carlo Percent Relative Bias

$$RB_{mc}(\hat{V}) = 100 \times \frac{B^{-1} \sum_{b=1}^B \hat{V}^{(b)} - V}{V}$$

where the global variance V was approximated through an independent set of 100,000 simulations. To evaluate the contribution of some component \hat{V}_a into the variance estimator \hat{V} , we computed the contribution (in percent)

$$CONTR_{mc}(\hat{V}_a) = 100 \times \frac{\frac{1}{B} \sum_{b=1}^B \hat{V}_a^{(b)}}{\frac{1}{B} \sum_{b=1}^B \hat{V}^{(b)}}$$

To evaluate the simplified variance estimator for the non-response \hat{V}_{simp}^{nr} , we computed the Monte Carlo Percent Relative Bias

$$RB_{mc}(\hat{V}_{simp}^{nr}) = 100 \times \frac{B^{-1} \sum_{b=1}^B \hat{V}_{simp}^{(b)} - V^{nr}}{V^{nr}},$$

where the variance V^{nr} due to non-response was approximated through an independent set of 100,000 simulations.

The simulation results are presented in Table 5.2. The proposed variance estimator is almost unbiased in all cases. As could be expected, the contribution of the variance due to the sampling design decreases with time, as the number of respondents decreases and as the variance due to non-response becomes larger. The simplified variance estimator is highly biased for the variance due to non-response in case of \hat{Y}_t . The bias decreases quickly with time, but remains large at time $t = 3$. The simplified variance estimator is almost unbiased for a calibrated estimator when the working model is adequately specified, but is severely biased

otherwise. This is consistent with our reasoning in Section 3.1. The simplified variance estimator is almost unbiased for the three estimators of the ratio, and for the calibrated estimators of the change in totals. In case of the non-calibrated estimator for the change in totals, the bias can be as high as 30%.

Table 5.2

Relative bias of a global variance estimator, relative contribution to the estimators of variance components and relative bias of a simplified variance estimator for the variance due to non-response for the estimation of a total, a ratio or a change in totals with three sets of weights

	$t = 1$	$t = 2$	$t = 3$	$t = 1$	$t = 2$	$t = 3$	$t = 1$	$t = 2$	$t = 3$
	\hat{Y}_t			\hat{Y}_{wt}			$\hat{Y}_{\bar{w}t}$		
$RB_{mc}(\hat{V})$	0	-1	-2	-1	-1	-2	-1	-1	-3
$CONTR_{mc}(\hat{V}_t^P)$	81	57	35	69	49	32	80	56	35
$CONTR_{mc}(\hat{V}_t^{nr1})$	19	19	13	31	22	15	20	18	13
$CONTR_{mc}(\hat{V}_t^{nr2})$	-	25	18	-	28	19	-	25	17
$CONTR_{mc}(\hat{V}_t^{nr3})$	-	-	34	-	-	34	-	-	34
$RB_{mc}(\hat{V}_{t,simp}^{nr})$	559	188	80	0	-1	-2	83	34	15
	\hat{R}_t			\hat{R}_{wt}			$\hat{R}_{\bar{w}t}$		
$RB_{mc}(\hat{V})$	-	0	-2	-	-1	-2	-	-1	-2
$CONTR_{mc}(\hat{V}_t^P)$	-	49	32	-	49	32	-	50	33
$CONTR_{mc}(\hat{V}_t^{nr1})$	-	22	15	-	22	15	-	22	15
$CONTR_{mc}(\hat{V}_t^{nr2})$	-	28	19	-	28	19	-	28	19
$CONTR_{mc}(\hat{V}_t^{nr3})$	-	-	34	-	-	34	-	-	34
$RB_{mc}(\hat{V}_{t,simp}^{nr})$	-	0	0	-	-1	-2	-	-1	-1
	$\hat{\Delta}_{tt}(1 \rightarrow t)$			$\hat{\Delta}_{tt,w}(1 \rightarrow t)$			$\hat{\Delta}_{tt,\bar{w}}(1 \rightarrow t)$		
$RB_{mc}(\hat{V})$	-	0	-2	-	0	-2	-	-1	-3
$CONTR_{mc}(\hat{V}_t^P)$	-	50	33	-	49	32	-	50	33
$CONTR_{mc}(\hat{V}_t^{nr1})$	-	22	14	-	22	15	-	22	14
$CONTR_{mc}(\hat{V}_t^{nr2})$	-	28	18	-	28	19	-	28	18
$CONTR_{mc}(\hat{V}_t^{nr3})$	-	-	34	-	-	34	-	-	34
$RB_{mc}(\hat{V}_{t,simp}^{nr})$	-	19	30	-	-1	-2	-	3	5

6 Illustration

In this section, we aim at illustrating our results on a real data set from the ELFE survey. The population of inference consists of infants born in one of the 544 French maternity units during 2011, except very

premature infants. Our illustration is meant to mimic as closely as possible the methodology of the ELFE survey. In particular, the modeling of attrition at each time is performed with variables available at baseline as explanatory variables only. As pointed out by the Associate Editor, under the MAR assumption, the variables of interest measured at any times $\delta < t$ may also have been used to model attrition between times $t - 1$ and t .

An original sample s_0 of about 35,600 infants was originally selected when the babies were just a few days old and were still at the maternity unit. The sample was selected using a cross-classified sampling design (Skinner, 2015; Juillard, Chauvet and Ruiz-Gazen, 2016). A sample of days and a sample of maternity units were independently selected, and both sample selections may be approximated by stratified simple random sampling (STSI). The sample consisted in all the infants born during one of the 25 selected days in one of the 320 selected maternity units.

Among the 35,600 infants originally selected, a total of 18,329 face-to-face interviews were completed with their families, which represents a response rate of 51%. This led to the subsample s_1 after accounting for non-response. The weights at time $t = 1$ were computed on the basis of the original sampling weights, adjusted in two steps. First, response probabilities were estimated by means of a model of Response Homogeneity Groups (RHGs), with 20 RHGs defined by using a logistic regression model with explanatory variables *Age of the mother*, *Gemellary identity* and *Season of birth*. Then, a calibration by means of the raking ratio method was performed on the binary variables *Born within marriage*, *Immigrant mother* and *Gemellary identity*.

When the children reached the age of two months, the parents had the first phone interview with a response rate of 87%. This leads to the subsample s_2 . The weights at time $t = 2$ were computed on the basis on the weight obtained at time $t = 1$, with a two-step adjustment. First, response probabilities were estimated by means of 20 RHGs, defined by using a logistic regression with explanatory variables *Age of the mother*, *Mother nationality* and *Father present at childbirth*. Then, a calibration by the raking ratio method was performed on the same calibration variables as at time $t = 1$.

When the children were one year old, the parents were contacted by phone with a response rate of 77%. This led to the subsample s_3 . The weights at time $t = 3$ were computed on the basis on the weights obtained at time $t = 2$, with a two-step adjustment similar to that realized at time $t = 2$.

We considered three variables of interest: *Breastfeeding exclusivity at the childbirth*, *at two month*, *at one year*. For each of these variables, we computed the estimator \hat{R}_t and the calibrated estimator \hat{R}_{wt} for the percentage $R(t)$ of breastfeeding among all the children at time t , and the associated variance estimators. We also computed the estimated coefficient of variation (in percent), defined as

$$\widehat{CV}_t(\hat{Y}_t) = 100 \times \frac{\sqrt{\hat{V}_t(\hat{Y}_t)}}{\hat{Y}_t}. \quad (6.1)$$

For each component \hat{V}_{ia} in the estimated variance \hat{V}_t , we computed its contribution (in percent) defined as

$$\text{CONTR}(\hat{V}_{ta}) = 100 \times \frac{\hat{V}_{ta} - \hat{V}_t}{\hat{V}_t}. \tag{6.2}$$

We also computed the simplified variance estimator for non-response $\hat{V}_{t, \text{simp}}^{\text{nr}}$, and the relative difference (in percent) with the approximately unbiased variance estimator \hat{V}_t^{nr} defined as

$$\text{RD}(\hat{V}_{t, \text{simp}}^{\text{nr}}) = 100 \times \frac{\hat{V}_{t, \text{simp}}^{\text{nr}} - \hat{V}_t^{\text{nr}}}{\hat{V}_t^{\text{nr}}}. \tag{6.3}$$

The results are given in Table 6.1. As observed in the simulation study, the RD of the simplified variance estimator for non-response is negligible in all cases.

Table 6.1
Estimates for a ratio, variance estimates, coefficient of variation, relative contributions of variance components and relative difference of a simplified variance estimator for a variable in the ELFE survey

Breastfeeding exclusivity	<i>t</i> = 1	<i>t</i> = 2	<i>t</i> = 3	<i>t</i> = 1	<i>t</i> = 2	<i>t</i> = 3
	maternity	2 months	1 year	maternity	2 months	1 year
	without calibration			with calibration		
\hat{R}_t (%)	59.0	30.6	3.3	59.4	31.0	3.4
$\hat{V}(\hat{R}_t)$	1.34E-05	1.50E-05	2.58E-06	1.28E-05	1.48E-05	2.60E-06
$\hat{C}V(\hat{Y}_t)$ (%)	0.6	1.3	4.8	0.6	1.2	4.7
$\text{CONTR}(\hat{V}_t^p)$	31	34	24	28	34	25
$\text{CONTR}(\hat{V}_t^{\text{nr}1})$	69	51	42	72	51	41
$\text{CONTR}(\hat{V}_t^{\text{nr}2})$	-	15	13	-	15	13
$\text{CONTR}(\hat{V}_t^{\text{nr}3})$	-	-	21	-	-	21
$\text{RD}(\hat{V}_{t, \text{simp}}^{\text{nr}})$	2	2	0	1	2	0

7 Conclusion

In this paper, we considered variance estimation accounting for weighting adjustments in panel surveys. We proposed both an approximately unbiased variance estimator and a simplified variance estimator for estimators of totals, complex parameters and measures of change, which covers most cases that may be encountered in practice. Our simulation results indicate that the proposed variance estimator performs well in all cases considered. The simplified variance estimator tends to overestimate the variance of the expansion estimator for totals, and to overestimate the variance for calibrated estimators of totals when the calibration variables lack of explanatory power for the variable of interest. However, the simplified variance estimator performs well for the estimation of ratios and change in totals with calibrated weights, even if the calibration model is not appropriate for the study variable.

The assumption of independent response behaviour is usually not tenable for multi-stage surveys, since units within clusters tend to be correlated with respect to the response behaviour. In this context, estimation of response probabilities based upon conditional logistic regression in the context of correlated responses has been studied by Skinner and D'Arrigo (2011), see also Kim, Kwon and Park (2016). Extending the present work in the context of correlated response behaviour is a challenging problem for further research.

Acknowledgements

We thank the Editors, an Associate Editor and the referees for useful comments and suggestions which led to an improvement of the paper.

Appendix

Estimation of the variance due to non-response for Response Homogeneity Groups

We consider the model of Response Homogeneity Groups introduced in Section 2.5. Recall that this model may be summarized as follows: at each time $\delta = 1, \dots, t$, the sub-sample $s_{\delta-1}$ is partitioned into $C(\delta - 1)$ groups $s_{\delta-1}^c$, $c = 1, \dots, C(\delta - 1)$. The response probabilities are assumed to be constant within the groups.

This model is equivalent to the logistic regression model in (2.18), with

$$z_i^\delta = [1\{i \in s_{\delta-1}^1\}, \dots, 1\{i \in s_{\delta-1}^{C(\delta-1)}\}]^\top. \quad (\text{A.1})$$

The equation (2.2) leads to the estimated response probabilities

$$\hat{p}_i^\delta = \frac{\sum_{i \in s_{\delta-1}^c} k_i^\delta r_i^\delta}{\sum_{i \in s_{\delta-1}^c} k_i^\delta} \quad \text{for} \quad i \in s_{\delta-1}^c. \quad (\text{A.2})$$

We first consider the case when the reweighted estimator is computed at time $t = 1$. In the estimator of the variance due to non-response given in (2.21), the vector $\hat{\gamma}_1^1$ simplifies as

$$\hat{\gamma}_1^1 = \left(\frac{\sum_{i \in s_1 \cap s_0^1} \frac{y_{i1}}{\pi_i}}{\hat{p}_1^1 \sum_{i \in s_1 \cap s_0^1} k_i^1}, \dots, \frac{\sum_{i \in s_1 \cap s_0^{C(0)}} \frac{y_{i1}}{\pi_i}}{\hat{p}_{C(0)}^1 \sum_{i \in s_1 \cap s_0^{C(0)}} k_i^1} \right)^\top. \quad (\text{A.3})$$

After some algebra, the variance estimator in (2.21) may be rewritten as

$$\hat{V}_1^{\text{nr}}(\hat{Y}_1) = \sum_{c=1}^{C(0)} \frac{(1 - \hat{p}_c^1)}{(\hat{p}_c^1)^2} \sum_{i \in s_1 \cap s_0^c} \left(\frac{y_{i1}}{\pi_i} - k_i^1 \frac{\sum_{j \in s_1 \cap s_0^c} \frac{y_{j1}}{\pi_j}}{\sum_{j \in s_1 \cap s_0^c} k_j^1} \right)^2. \quad (\text{A.4})$$

We now consider the case when the reweighted estimator is computed at time $t = 2$. We focus on the simpler case when the same system of RHGs is kept over time. In the estimator of the variance due to non-response given in (2.22), the vectors $\hat{\gamma}_2^1$ and $\hat{\gamma}_2^2$ simplify as

$$\hat{\gamma}_2^1 = \left(\frac{\sum_{i \in s_2 \cap s_1^1} \frac{y_{i2}}{\pi_i}}{\hat{p}_1^1 \sum_{i \in s_2 \cap s_1^1} k_i^1}, \dots, \frac{\sum_{i \in s_2 \cap s_1^{C(0)}} \frac{y_{i2}}{\pi_i}}{\hat{p}_{C(0)}^1 \sum_{i \in s_2 \cap s_1^{C(0)}} k_i^1} \right)^\top, \tag{A.5}$$

$$\hat{\gamma}_2^2 = \left(\frac{\sum_{i \in s_2 \cap s_1^1} \frac{y_{i2}}{\pi_i}}{\hat{p}_1^1 \hat{p}_1^2 \sum_{i \in s_2 \cap s_1^1} k_i^2}, \dots, \frac{\sum_{i \in s_2 \cap s_1^{C(0)}} \frac{y_{i2}}{\pi_i}}{\hat{p}_{C(0)}^1 \hat{p}_{C(0)}^2 \sum_{i \in s_2 \cap s_1^{C(0)}} k_i^2} \right)^\top. \tag{A.6}$$

After some algebra, the variance estimator in (2.22) may be rewritten as

$$\begin{aligned} \hat{V}_2^{nr}(\hat{Y}_2) &= \sum_{c=1}^{C(0)} \frac{(1 - \hat{p}_c^1)}{\hat{p}_c^2} \sum_{i \in s_2 \cap s_1^c} \left(\frac{y_{i2}}{\pi_i \hat{p}_c^1} - k_i^1 \frac{\sum_{j \in s_2 \cap s_1^c} \frac{y_{j2}}{\pi_j}}{\sum_{j \in s_2 \cap s_1^c} k_j^1} \right)^2 \\ &+ \sum_{c=1}^{C(0)} (1 - \hat{p}_c^2) \sum_{i \in s_2 \cap s_1^c} \left(\frac{y_{i2}}{\pi_i \hat{p}_c^1 \hat{p}_c^2} - k_i^2 \frac{\sum_{j \in s_2 \cap s_1^c} \frac{y_{j2}}{\pi_j}}{\sum_{j \in s_2 \cap s_1^c} k_j^2} \right)^2. \end{aligned} \tag{A.7}$$

If we further assume that k_i^δ is constant over times $\delta = 1, 2$, and may thus be rewritten as k_i , the expression in (A.7) simplifies as

$$\hat{V}_2^{nr}(\hat{Y}_2) = \sum_{c=1}^{C(0)} \frac{(1 - \hat{p}_c^{1 \rightarrow 2})}{(\hat{p}_c^{1 \rightarrow 2})^2} \sum_{i \in s_2 \cap s_1^c} \left(\frac{y_{i2}}{\pi_i} - k_i \frac{\sum_{j \in s_2 \cap s_1^c} \frac{y_{j2}}{\pi_j}}{\sum_{j \in s_2 \cap s_1^c} k_j} \right)^2. \tag{A.8}$$

with $\hat{p}_c^{1 \rightarrow 2} = \prod_{\delta=1}^2 \hat{p}_c^\delta$ for $c = 1, \dots, C(0)$. This simplification of the variance estimator can be extended to the reweighted estimator at time t . Assuming that the RHGs are kept over time, and that $k_i^\delta = k_i$ for any $\delta = 1, \dots, t$, the variance estimator in (2.12) may be written as

$$\hat{V}_t^{nr}(\hat{Y}_t) = \sum_{c=1}^{C(0)} \frac{(1 - \hat{p}_c^{1 \rightarrow t})}{(\hat{p}_c^{1 \rightarrow t})^2} \sum_{i \in s_t \cap s_{t-1}^c} \left(\frac{y_{it}}{\pi_i} - k_i \frac{\sum_{j \in s_t \cap s_{t-1}^c} \frac{y_{jt}}{\pi_j}}{\sum_{j \in s_t \cap s_{t-1}^c} k_j} \right)^2 \tag{A.9}$$

with $\hat{p}_c^{1 \rightarrow t} = \prod_{\delta=1}^t \hat{p}_c^\delta$ for $c = 1, \dots, C(0)$.

References

Beaumont, J.-F. (2005). Calibrated imputation in surveys under a quasimodel-assisted approach. *Journal of the Royal Statistical Society, Series B*, 67, 445-458.

Beaumont, J.-F., and Haziza, D. (2016). A note on the concept of invariance in two-phase sampling designs. *Survey Methodology*, 42, 2, 319-323. Paper available at <https://www150.statcan.gc.ca/n1/pub/12-001-x/2016002/article/14662-eng.pdf>.

- Berger, Y. (2004). Variance estimation for measures of change in probability sampling. *Canadian Journal of Statistics*, 32, 4, 451-467.
- Caron, N., and Ravalet, P. (2000). Estimation dans les enquêtes répétées : application à l'enquête emploi en continu. Technical report INSEE, Paris.
- Chauvet, G., and Goga, C. (2018). Linearization versus bootstrap for variance estimation of the change between Gini indexes. *Survey Methodology*, 44, 1, 17-42. Paper available at <https://www150.statcan.gc.ca/n1/pub/12-001-x/2018001/article/54926-eng.pdf>.
- Clarke, P., and Tate, P. (2002). An application of non-ignorable non-response models for gross flows estimation in the British labour force survey. *Australian & New Zealand Journal of Statistics*, 4, 413-425.
- Deville, J.-C., and Särndal, C.-E. (1992). Calibration estimators in survey sampling. *Journal of the American Statistical Association*, 87, 376-382.
- Ekholm, A., and Laaksonen, S. (1991). Weighting via response modeling in the finnish household budget survey. *Journal of Official Statistics*, 7, 325-327.
- Fay, R. (1992). When are inferences from multiple imputation valid? *Proceedings of the Survey Research Methods Section*, American Statistical Association, 81, 1, 227-232.
- Fuller, W., and An, A. (1998). Regression adjustment for non-response. *Journal of the Indian Society of Agricultural Statistics*, 51, 331-342.
- Fuller, W.A., Loughin, M.M. and Baker, H.D. (1994). Regression weighting in the presence of nonresponse with application to the 1987-1988 Nationwide Food Consumption Survey. *Survey Methodology*, 20, 1, 75-85. Paper available at <https://www150.statcan.gc.ca/n1/pub/12-001-x/1994001/article/14429-eng.pdf>.
- Goga, C., Deville, J.-C. and Ruiz-Gazen, A. (2009). Composite estimation and linearization method for two-sample survey data. *Biometrika*, 96, 691-709.
- Hawkes, D., and Plewis, I. (2009). Modelling nonresponse in the national child development study. *Journal of the Royal Statistical Society, Series A*, 169, 479-491.
- Juillard, H., Chauvet, G. and Ruiz-Gazen, A. (2017). Estimation under cross-classified sampling with application to a childhood survey. *Journal of the American Statistical Association*, 112, 850-858.
- Kalton, G. (2009). Design for surveys over time. *Handbook of Statistics*, 29, 89-108.
- Kim, J.K., and Kim, J.J. (2007). Nonresponse weighting adjustment using estimated response probability. *Canadian Journal of Statistics*, 35, 501-514.
- Kim, J.K., Kwon, Y. and Park, M. (2016). Calibrated propensity score method for survey nonresponse in cluster sampling. *Biometrika*, 103, 461-473.
- Laaksonen, S. (2007). Weighting for two-phase surveyed data. *Survey Methodology*, 33, 2, 121-130. Paper available at <https://www150.statcan.gc.ca/n1/pub/12-001-x/2007002/article/10489-eng.pdf>.

- Laaksonen, S., and Chambers, R.L. (2006). Survey estimation under informative nonresponse with follow-up. *Journal of Official Statistics*, 22, 81-95.
- Laniel, N. (1988). Variances for a rotating sample from a changing population. *Proceedings of the Business and Economics Statistics Section*, American Statistical Association, 246-250.
- Laurie, H., Smith, R. and Scott, L. (1999). Strategies for reducing nonresponse in a longitudinal panel survey. *Journal of Official Statistics*, 15, 269-282.
- Lynn, P. (2009). Methods for longitudinal surveys. *Methodology of Longitudinal Surveys*, 1-19.
- Nordberg, L. (2000). On variance estimation for measures of change when samples are coordinated by the use of permanent random numbers. *Journal of Official Statistics*, 16, 363-378.
- Pirus, C., Bois, C., Dufourg, M., Lanoë, J., Vandentorren, S., Leridon, H. and the Elfe team (2010). Constructing a cohort: Experience with the French Elfe project. *Population*, 65, 637-670.
- Qualité, L., and Tillé, Y. (2008). Variance estimation of changes in repeated surveys and its application to the Swiss survey of value added. *Survey Methodology*, 34, 2, 173-181. Paper available at <https://www150.statcan.gc.ca/n1/pub/12-001-x/2008002/article/10758-eng.pdf>.
- Rendtel, U., and Harms, T. (2009). Weighting and calibration for household panels. *Methodology of Longitudinal Surveys*, 265-286.
- Rizzo, L., Kalton, G. and Brick, J.M. (1996). A comparison of some weighting adjustment methods for panel nonresponse. *Survey Methodology*, 22, 1, 43-53. Paper available at <https://www150.statcan.gc.ca/n1/pub/12-001-x/1996001/article/14386-eng.pdf>.
- Silva, P., and Skinner, C. (1997). Cross-classified sampling: Some estimation theory. *Variable Selection for Regression Estimation in Finite Populations*, 23, 23-32.
- Skinner, C. (2015). Cross-classified sampling: Some estimation theory. *Statistics & Probability Letters*, 104, 163-168.
- Skinner, C., and D'Arrigo, J. (2011). Inverse probability weighting for clustered non-response. *Biometrika*, 98, 953-966.
- Skinner, C., and Vieira, M. (2005). Design effects in the analysis of longitudinal survey data. S3RI Methodology Working Papers, M05/13. Southampton, UK: Southampton Statistical Sciences Research Institute.
- Slud, E.V., and Bailey, L. (2010). Evaluation and selection of models for attrition nonresponse adjustment. *Journal of Official Statistics*, 26, 1-18.
- Tam, S. (1984). On covariance from overlapping samples. *The American Statistician*, 38, 1-18.
- Vandecasteele, L., and Debels, A. (2007). Attrition in panel data: The effectiveness of weighting. *European Sociological Review*, 23, 1, 81-97.
- Zhou, M., and Kim, J. (2012). An efficient method of estimation for longitudinal surveys with monotone missing data. *Biometrika*, 99, 631-648.