

WORKING PAPERS

N° TSE-659

April 2016

“Estimation under cross-classified sampling with
application to a childhood survey”

Hélène Juillard, Guillaume Chauvet and Anne Ruiz-Gazen

Estimation under cross-classified sampling with application to a childhood survey

Hélène Juillard*
Guillaume Chauvet†
Anne Ruiz-Gazen‡

April 25, 2016

Estimation under cross-classified sampling with application to a childhood survey

Abstract

The cross-classified sampling design consists in drawing samples from a two-dimension population, independently in each dimension. Such design is commonly used in consumer price index surveys and has been recently applied to draw a sample of babies in the French Longitudinal Survey on Childhood, by crossing a sample of maternity units and a sample of days. We propose to derive a general theory of estimation for this sampling design. We consider the Horvitz-Thompson estimator for a total, and show that the cross-classified design will usually result in a loss of efficiency as compared to the widespread two-stage design. We obtain the asymptotic distribution of the Horvitz-Thompson estimator, and several unbiased variance estimators. Facing the problem of possibly negative values, we propose simplified non-negative variance estimators and study their bias under a super-population model. The proposed estimators are compared for totals and ratios on

*INED, 133 boulevard Davout, 75020 Paris, France

†ENSAI/IRMAR, Campus de Ker Lann, 35170 Bruz, France

‡Toulouse School of Economics, 21 allée de Brienne, 31000 Toulouse, France

simulated data. An application on real data from the French Longitudinal Survey on Childhood is also presented, and we make some recommendations. Supplementary materials are available online.

Some key words: analysis of variance, Horvitz-Thompson estimator, independence, invariance, Sen-Yates-Grundy estimator, two-stage sampling.

Short title: Estimation under cross-classified sampling

1 Introduction

The 2011 French Longitudinal Survey on Childhood ELFE (Etude Longitudinale Française depuis l'Enfance) comprises more than 18,000 children selected on the basis of their place and date of birth. On the one hand, a sample of 320 maternity units has been drawn. On the other hand, a sample of 25 days divided in four time periods and spread across the four seasons of 2011 has been selected. The babies born at the sampled locations and on the sampled days have been approached through midwives. Data were collected on babies whose parents consented to their inclusion during their stay at the maternity unit. ELFE is conducted by the National Institute for Demographic Studies, the National Institute for Health and Medical Research and the French Blood Agency. The objective of observing children born within the same year is to analyze their physical and psychological health together with their living and environmental conditions. This large-scale study of children's development and socialization is the first of its kind in France. The collected data are now available to public and private research teams and many projects are underway in areas such as health, health environment and social sciences. In order to derive reliable confidence intervals for finite population parameters such as totals or ratios, the ELFE sampling design has to be taken into account.

The ELFE sample is drawn according to a non-standard sampling design, called Cross-Classified Sampling (CCS), following Ohlsson (1996). It consists in drawing

independently two samples from each component of a two-dimensional population. In the ELFE survey, a sample of maternity units and a sample of days are independently selected. This sampling design appears in other contexts than the ELFE survey. Some examples include consumer price index surveys, as detailed in Dalén & Ohlsson (1995) for the Swedish survey, where outlets and items are sampled, and business surveys (Skinner, 2015), where businesses and products are sampled. Due to its particular properties, CCS deserves a specific attention. However, as noted by Skinner (2015), "the literature on the theory of cross-classified sampling is very limited". In particular, no general theory is derived under the finite population framework. While the papers by Vos (1964) and Ohlsson (1996) focus on simple random sampling without replacement, Skinner (2015) gives some results under stratified without replacement simple random sampling and under with replacement unequal probability sampling. Dalén & Ohlsson (1995) provide some results under probability proportional to size without-replacement sampling.

In the present paper, we develop a general theory for estimation and variance estimation under CCS. The asymptotic normality of the Horvitz-Thompson estimator is derived under some mild conditions. A comparison with a two-stage sampling design is carried out in a general framework. We also raise an issue, not reported before, of possible negative values for Horvitz-Thompson and Yates-Grundy variance estimates. This problem occurs even in the simplest case of simple random sampling without replacement. Non-negative simplified variance estimators are therefore introduced. Conditions for their approximate unbiasedness are given under a design-based and a model-based approach. The properties of our variance estimators are evaluated through a small but realistic simulation study when estimating totals and ratios. Finally, an application to the ELFE data is detailed.

2 Cross-classified sampling design

2.1 Notations and Horvitz-Thompson estimation

Keeping in mind the ELFE survey, we consider a population U_M of N_M maternity units and a population U_D of N_D days. However, the developments below are completely general and may be applied to any populations U_M and U_D . We will use the indexes i and j for the maternity units, and the indexes k and l for the days. We consider a sampling design $p_M(\cdot)$ on the population U_M , leading to a sample S_M of (average) size n_M , and a sampling design $p_D(\cdot)$ on the population U_D leading to a sample S_D of (average) size n_D . We assume that the two samples are selected independently. The cross-classified sampling design $p(\cdot)$ on the product population $U = U_M \times U_D$ is therefore defined as

$$p(s) = p_M(s_M) \times p_D(s_D) \quad \text{for any } s = s_M \times s_D \subset U_M \times U_D.$$

Let π_i^M denote the probability that i is selected in S_M , π_{ij}^M denote the probability that units i and j are selected jointly in S_M , and let $\Delta_{ij}^M = \pi_{ij}^M - \pi_i^M \pi_j^M$. The quantities π_k^D , π_{kl}^D and Δ_{kl}^D are similarly defined. We assume that the first and second-order inclusion probabilities are non-negative in each population. The probability for the pairs (i, k) to be selected in the product sample $S_M \times S_D$ is $\pi_i^M \pi_k^D$, and the probability for the pairs (i, k) and (j, l) to be selected jointly in the product sample $S_M \times S_D$ is $\pi_{ij}^M \pi_{kl}^D$.

We are interested in some variable of interest with value Y_{ik} for the maternity unit i and the day k . The total $t_Y = \sum_{i \in U_M} \sum_{k \in U_D} Y_{ik}$ is then unbiasedly estimated by the Horvitz-Thompson (HT) estimator

$$\hat{t}_Y = \sum_{i \in S_M} \sum_{k \in S_D} \frac{Y_{ik}}{\pi_i^M \pi_k^D} = \sum_{i \in S_M} \sum_{k \in S_D} \check{Y}_{ik} \quad \text{where } \check{Y}_{ik} = \frac{Y_{ik}}{\pi_i^M \pi_k^D}. \quad (2.1)$$

Making use of the independence between S_M and S_D , the variance of the HT-

estimator is

$$V_{CCS}(\hat{t}_Y) = \sum_{i,j \in U_M} \sum_{k,l \in U_D} \Gamma_{ijkl} \check{Y}_{ik} \check{Y}_{jl} \quad (2.2)$$

where $\Gamma_{ijkl} = \pi_{ij}^M \pi_{kl}^D - \pi_i^M \pi_j^M \pi_k^D \pi_l^D$. The Sen(1953)-Yates-Grundy(1953) form

$$V_{CCS}(\hat{t}_Y) = -\frac{1}{2} \sum_{(i,k) \neq (j,l) \in U_M \times U_D} \Gamma_{ijkl} (\check{Y}_{ik} - \check{Y}_{jl})^2 \quad (2.3)$$

can be used alternatively when both sampling designs are of fixed size.

Our set-up can be compared to the usual two-stage framework, by considering U_M as a population of Primary Sampling Units (PSUs) and U_D as a population of Secondary Sampling Units (SSUs), each maternity unit i being associated to the same population of days. In case of two-stage sampling, denoted by MD , a first-stage sample S_M is selected in U_M , and some second-stage samples S_i are selected independently using $p_D(S_i)$ for each $i \in S_M$ (see Särndal et al., 1992). The variance of the HT-estimator is then

$$V_{MD}(\hat{t}_Y) = V_{MD}^{PSU}(\hat{t}_Y) + V_{MD}^{SSU}(\hat{t}_Y) \quad (2.4)$$

where

$$V_{MD}^{PSU}(\hat{t}_Y) = \sum_{i,j \in U_M} \sum_{k,l \in U_D} \Delta_{ij}^M \pi_k^D \pi_l^D \check{Y}_{ik} \check{Y}_{jl}, \quad (2.5)$$

$$V_{MD}^{SSU}(\hat{t}_Y) = \sum_{i \in U_M} \sum_{k,l \in U_D} \pi_i^M \Delta_{kl}^D \check{Y}_{ik} \check{Y}_{il}. \quad (2.6)$$

Alternatively, we could consider U_D as a population of PSUs and U_M as a population of SSUs, each day k being associated to the same population of maternity units. In this case, the variance of the HT-estimator under two-stage sampling is

$$V_{DM}(\hat{t}_Y) = V_{DM}^{PSU}(\hat{t}_Y) + V_{DM}^{SSU}(\hat{t}_Y) \quad (2.7)$$

where

$$V_{DM}^{PSU}(\hat{t}_Y) = \sum_{k,l \in U_D} \sum_{i,j \in U_M} \Delta_{kl}^D \pi_i^M \pi_j^M \check{Y}_{ik} \check{Y}_{jl}, \quad (2.8)$$

$$V_{DM}^{SSU}(\hat{t}_Y) = \sum_{k \in U_D} \sum_{i,j \in U_M} \pi_k^D \Delta_{ij}^M \check{Y}_{ik} \check{Y}_{il}. \quad (2.9)$$

The different features of CCS and two-stage sampling on a two-dimension population are illustrated on Figure 1.

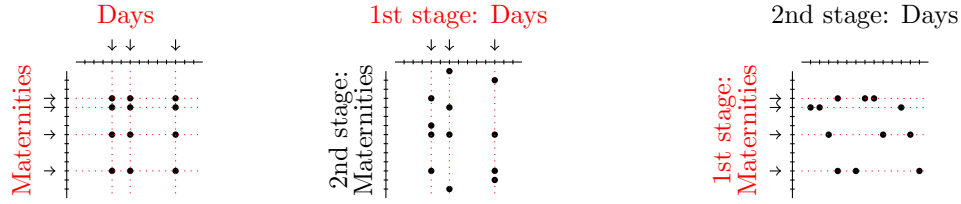


Figure 1: Cross-classified sampling (left panel), two-stage sampling DM with primary units in U_D (central panel), two-stage sampling MD with primary units in U_M (right panel)

2.2 Variance decomposition for cross-classified sampling

The covariance Γ_{ijkl} may be written in several ways, leading to alternative variance decompositions. Plugging $\Gamma_{ijkl} = \pi_{kl}^D \Delta_{ij}^M + \pi_{ij}^M \Delta_{kl}^D - \Delta_{ij}^M \Delta_{kl}^D$ into (2.2) gives

$$V_{CCS}(\hat{t}_Y) = V_1(\hat{t}_Y) + V_2(\hat{t}_Y) - V_3(\hat{t}_Y) \quad (2.10)$$

where

$$V_1(\hat{t}_Y) = \sum_{k,l \in U_D} \sum_{i,j \in U_M} \pi_{kl}^D \Delta_{ij}^M \check{Y}_{ik} \check{Y}_{jl}, \quad (2.11)$$

$$V_2(\hat{t}_Y) = \sum_{i,j \in U_M} \sum_{k,l \in U_D} \pi_{ij}^M \Delta_{kl}^D \check{Y}_{ik} \check{Y}_{jl}, \quad (2.12)$$

$$V_3(\hat{t}_Y) = \sum_{i,j \in U_M} \sum_{k,l \in U_D} \Delta_{ij}^M \Delta_{kl}^D \check{Y}_{ik} \check{Y}_{jl}. \quad (2.13)$$

Plugging $\Gamma_{ijkl} = \Delta_{ij}^M \pi_k^D \pi_l^D + \Delta_{kl}^D \pi_i^M \pi_j^M + \Delta_{ij}^M \Delta_{kl}^D$ into (2.2) gives

$$V_{CCS}(\hat{t}_Y) = V_{MD}^{PSU}(\hat{t}_Y) + V_{DM}^{PSU}(\hat{t}_Y) + V_3(\hat{t}_Y) \quad (2.14)$$

and we have $V_1(\hat{t}_Y) = V_{MD}^{PSU}(\hat{t}_Y) + V_3(\hat{t}_Y)$ and $V_2(\hat{t}_Y) = V_{DM}^{PSU}(\hat{t}_Y) + V_3(\hat{t}_Y)$. This second decomposition was originally derived by Dalén & Ohlsson (1995). It is also given in Ohlsson (1996), and in equation (3) of Theorem 2.2 of Skinner (2015). Other decompositions are possible, e.g. through an analysis of variance decomposition as for two-stage sampling.

2.3 Comparison with two-stage sampling

From expressions (2.7) and (2.14), we obtain after some algebra that

$$V_{CCS}(\hat{t}_Y) - V_{DM}(\hat{t}_Y) = \sum_{i,j \in U_M} \Delta_{ij}^M \sum_{k \neq l \in U_D} \pi_{kl}^D \check{Y}_{ik} \check{Y}_{jl}. \quad (2.15)$$

In case of Poisson sampling (PO) inside U_M and when Y is assumed to be non-negative, the right-hand side in (2.15) is non-negative and CCS is thus less efficient than two-stage sampling. In case of fixed-size sampling inside U_M , equation (2.15) may be alternatively written as

$$V_{CCS}(\hat{t}_Y) - V_{DM}(\hat{t}_Y) = \sum_{i \neq j \in U_M} \frac{(-\Delta_{ij}^M)}{2} \sum_{k \neq l \in U_D} \frac{\pi_{kl}^D}{\pi_k^D \pi_l^D} \left(\frac{Y_{ik}}{\pi_i^M} - \frac{Y_{jk}}{\pi_j^M} \right) \left(\frac{Y_{il}}{\pi_i^M} - \frac{Y_{jl}}{\pi_j^M} \right) \quad (2.16)$$

If the so-called Sen-Yates-Grundy conditions are respected for p_M , the quantities $(-\Delta_{ij}^M)$ are non-negative. If Y_{ik} is roughly proportional to the size of the maternity unit i , as can be expected for count variables, the quantities

$$\left(\frac{Y_{ik}}{\pi_i^M} - \frac{Y_{jk}}{\pi_j^M} \right) \left(\frac{Y_{il}}{\pi_i^M} - \frac{Y_{jl}}{\pi_j^M} \right)$$

will tend to be positive unless the inclusion probabilities π_i^M are defined proportionally to some measure of size. CCS sampling would then be less efficient than two-stage sampling. This result is illustrated in section 4.1 on some simulated populations when both p_M and p_D are simple random sampling without replacement (SI) designs, and for different sample sizes.

3 Variance estimation

3.1 Design-unbiased variance estimation

The HT variance estimator for $V_{CCS}(\hat{t}_Y)$ is

$$\hat{V}_{HT}(\hat{t}_Y) = \sum_{i,j \in S_M} \sum_{k,l \in S_D} \frac{\Gamma_{ijkl}}{\pi_{ij}^M \pi_{kl}^D} \check{Y}_{ik} \check{Y}_{jl}. \quad (3.1)$$

It may be also derived from (2.10), leading to the alternative writing

$$\hat{V}_{HT}(\hat{t}_Y) = \hat{V}_{1,HT}(\hat{t}_Y) + \hat{V}_{2,HT}(\hat{t}_Y) - \hat{V}_{3,HT}(\hat{t}_Y) \quad (3.2)$$

where

$$\hat{V}_{1,HT}(\hat{t}_Y) = \sum_{i,j \in S_M} \sum_{k,l \in S_D} \frac{\Delta_{ij}^M}{\pi_{ij}^M} \check{Y}_{ik} \check{Y}_{jl}, \quad (3.3)$$

$$\hat{V}_{2,HT}(\hat{t}_Y) = \sum_{i,j \in S_M} \sum_{k,l \in S_D} \frac{\Delta_{kl}^D}{\pi_{kl}^D} \check{Y}_{ik} \check{Y}_{jl}, \quad (3.4)$$

$$\hat{V}_{3,HT}(\hat{t}_Y) = \sum_{i,j \in S_M} \sum_{k,l \in S_D} \frac{\Delta_{ij}^M}{\pi_{ij}^M} \frac{\Delta_{kl}^D}{\pi_{kl}^D} \check{Y}_{ik} \check{Y}_{jl}. \quad (3.5)$$

If p_M and p_D are both Poisson sampling designs, this variance estimator is always non-negative. Otherwise, it may take negative values even if p_M and p_D are both SI designs (denoted by SI²) as illustrated in section 4.2. When p_M and p_D are both fixed-size sampling designs, we may alternatively consider the Yates-Grundy like variance estimator:

$$\hat{V}_{YG}(\hat{t}_Y) = \hat{V}_{1,YG}(\hat{t}_Y) + \hat{V}_{2,YG}(\hat{t}_Y) - \hat{V}_{3,YG}(\hat{t}_Y) \quad (3.6)$$

where

$$\hat{V}_{1,YG}(\hat{t}_Y) = -\frac{1}{2} \sum_{i \neq j \in S_M} \frac{\Delta_{ij}^M}{\pi_{ij}^M} \left(\frac{\hat{Y}_{i\bullet}}{\pi_i^M} - \frac{\hat{Y}_{j\bullet}}{\pi_j^M} \right)^2 \quad (3.7)$$

$$\hat{V}_{2,YG}(\hat{t}_Y) = -\frac{1}{2} \sum_{k \neq l \in S_D} \frac{\Delta_{kl}^D}{\pi_{kl}^D} \left(\frac{\hat{Y}_{\bullet k}}{\pi_k^D} - \frac{\hat{Y}_{\bullet l}}{\pi_l^D} \right)^2 \quad (3.8)$$

$$\hat{V}_{3,YG}(\hat{t}_Y) = -\frac{1}{2} \sum_{(i,k) \neq (j,l) \in S_M \times S_D} \frac{\Delta_{ij}^M \Delta_{kl}^D}{\pi_{ij}^M \pi_{kl}^D} (\check{Y}_{ik} - \check{Y}_{jl})^2 \quad (3.9)$$

with $\hat{Y}_{\bullet k} = \sum_{i \in S_M} Y_{ik} / \pi_i^M$ is the estimated sub-total for the day k and $\hat{Y}_{i\bullet} = \sum_{k \in S_D} Y_{ik} / \pi_k^D$ is the estimated sub-total for the maternity **unit** i .

It can be proved that $\hat{V}_{HT}(\hat{t}_Y)$ in (3.2) and $\hat{V}_{YG}(\hat{t}_Y)$ in (3.6) match term by term, when p_M and p_D are stratified simple random sampling (STSI) designs. In the same STSI context, another variance estimator is given in equation (4) of Theorem 2.2 in

Skinner (2015). This variance estimator does not match $\hat{V}_{HT}(\hat{t}_Y)$ or $\hat{V}_{YG}(\hat{t}_Y)$ term by term, since Skinner's variance estimator is based on the variance decomposition in equation (2.14), while our variance estimator is based on the variance decomposition in equation (2.10). Nevertheless, both variance estimators are globally identical in the STSI case.

Another variance estimator is obtained in Dalén & Ohlsson (1995), in case of a probability proportional to size without-replacement sampling design in both dimensions. Summing the variance component estimators in equations (4.1)-(4.3) of Dalén & Ohlsson (1995) leads to a similar variance estimator than in our equation (3.6), except that $-\hat{V}_{3,YG}(\hat{t}_Y)$ is replaced with $+\hat{V}_{3,YG}(\hat{t}_Y)$ which results in an overestimation of the variance. This overestimation can be neglected in cases when $V_3(\hat{t}_Y)$ is small as compared to the other variance components (see Table 1 in Section 4.2).

If both sampling designs satisfy the Sen-Yates-Grundy conditions (SYG), the terms $\hat{V}_{1,YG}(\hat{t}_Y)$ and $\hat{V}_{2,YG}(\hat{t}_Y)$ are non-negative. However, the term $\hat{V}_{3,YG}(\hat{t}_Y)$ is usually non-negative, which may lead to negative values for $\hat{V}_{YG}(\hat{t}_Y)$ as illustrated in the simulations of section 4.2. It is thus desirable to have access to non-negative variance estimators with limited bias.

3.2 Non-negative variance estimators

We consider the variance decomposition in (2.10), and study the relative order of magnitude of the components. We make the following assumptions:

H1: There exist some constants α_1 and α_2 such that

$$\forall k \in U_D, \quad \frac{1}{N_M} \sum_{i \in U_M} Y_{ik}^2 \leq \alpha_1, \quad \text{and} \quad \forall i \in U_M, \quad \frac{1}{N_D} \sum_{k \in U_D} Y_{ik}^2 \leq \alpha_2.$$

H2: There exist some constants $\lambda_1 > 0$ and $\lambda_2 > 0$ such that

$$\forall k \in U_D, \pi_k^D \geq \lambda_1 \frac{n_D}{N_D}, \quad \text{and} \quad \forall i \in U_M, \pi_i^M \geq \lambda_2 \frac{n_M}{N_M}.$$

H3: There exist some constants γ_1 and γ_2 such that

$$\forall k \neq l \in U_D, \frac{N_D^2}{n_D} \sup_{k \neq l \in U_D} |\Delta_{kl}^D| \leq \gamma_1, \quad \text{and} \quad \forall i \neq j \in U_M, \frac{N_M^2}{n_M} \sup_{i \neq j \in U_M} |\Delta_{ij}^M| \leq \gamma_2.$$

H4: There exists some constant $\delta > 0$ such that

$$V_{CCS}(\hat{t}_Y) \geq \delta N_M^2 N_D^2 \left(\frac{1}{n_M} + \frac{1}{n_D} \right).$$

It is assumed in (H1) that the variable y has bounded moments of order 2 for each maternity unit i and for each day k . Assumptions (H2) and (H3) are classical in survey sampling and are satisfied for many sampling designs, see for example Cardot et al. (2013). It is assumed in (H4) that the variance of the HT-estimator under CCS sampling has the order $N_M^2 N_D^2 (n_M^{-1} + n_D^{-1})$. From assumptions (H1-H4), there exist some constants C_1, C_2 and C_3 such that

$$\frac{V_1(\hat{t}_Y)}{V_{CCS}(\hat{t}_Y)} \leq C_1 \frac{1}{1 + n_M n_D^{-1}}, \quad (3.10)$$

$$\frac{V_2(\hat{t}_Y)}{V_{CCS}(\hat{t}_Y)} \leq C_2 \frac{1}{1 + n_D n_M^{-1}}, \quad (3.11)$$

$$\frac{V_3(\hat{t}_Y)}{V_{CCS}(\hat{t}_Y)} \leq C_3 \frac{1}{n_D n_M^{-1} + n_M n_D^{-1}} \quad (3.12)$$

The proof is given in Appendix 8. It follows from (3.10)-(3.12) that if n_D is large and n_M is bounded, both $V_2(\hat{t}_Y)$ and $V_3(\hat{t}_Y)$ are negligible and a non-negative simplified variance estimator can be derived by focusing on $V_1(\hat{t}_Y)$ only. This leads to

$$\hat{V}_{\text{SIMP1}}(\hat{t}_Y) = \hat{V}_{1,YG}(\hat{t}_Y). \quad (3.13)$$

If the sampling design p_D satisfies the SYG conditions, this simplified estimator is always non-negative. In the particular SI² case, we obtain

$$\hat{V}_{\text{SIMP1}}(\hat{t}_Y) = N_M^2 \left(\frac{1}{n_M} - \frac{1}{N_M} \right) s_{Y\bullet}^2. \quad (3.14)$$

where

$$s_{\hat{Y}_{\bullet\bullet}}^2 = \frac{1}{n_M - 1} \sum_{i \in S_M} \left(\hat{Y}_{i\bullet} - \frac{1}{n_M} \sum_{j \in S_M} \hat{Y}_{j\bullet} \right)^2. \quad (3.15)$$

Symmetrically, both $V_1(\hat{t}_Y)$ and $V_3(\hat{t}_Y)$ may be seen as negligible if n_M is large and n_D is bounded. Another simplified variance estimator is thus

$$\hat{V}_{\text{SIMP2}}(\hat{t}_Y) = \hat{V}_{2,YG}(\hat{t}_Y). \quad (3.16)$$

If the sampling design p_M satisfies the SYG conditions, this estimator is non-negative. In the particular SI^2 case, we have

$$\hat{V}_{\text{SIMP2}}(\hat{t}_Y) = N_D^2 \left(\frac{1}{n_D} - \frac{1}{N_D} \right) s_{\hat{Y}_{\bullet\bullet}}^2 \quad (3.17)$$

where

$$s_{\hat{Y}_{\bullet\bullet}}^2 = \frac{1}{n_D - 1} \sum_{k \in S_D} \left(\hat{Y}_{\bullet k} - \frac{1}{n_D} \sum_{l \in S_D} \hat{Y}_{\bullet l} \right)^2. \quad (3.18)$$

A third possible simplified variance estimator is

$$\begin{aligned} \hat{V}_{\text{SIMP3}}(\hat{t}_Y) &= \hat{V}_{\text{SIMP1}} + \hat{V}_{\text{SIMP2}} \\ &= \hat{V}_{1,YG}(\hat{t}_Y) + \hat{V}_{2,YG}(\hat{t}_Y). \end{aligned} \quad (3.19)$$

This estimator is non-negative if both p_D and p_M satisfy the SYG conditions. It is approximately unbiased for $V_{\text{CCS}}(\hat{t}_Y)$ if n_D is large and n_M is bounded, or if n_M is large and n_D is bounded. In the particular SI^2 case

$$\hat{V}_{\text{SIMP3}}(\hat{t}_Y) = N_M^2 \left(\frac{1}{n_M} - \frac{1}{N_M} \right) s_{\hat{Y}_{\bullet\bullet}}^2 + N_D^2 \left(\frac{1}{n_D} - \frac{1}{N_D} \right) s_{\hat{Y}_{\bullet\bullet}}^2. \quad (3.20)$$

Similar formula can be easily derived in the case of stratified simple random sampling without replacement and will be used in Section 5.

3.3 Relative bias under a superpopulation model

We consider the following superpopulation model

$$Y_{ik} = \mu + \sigma_M U_i + \sigma_D V_k + \sigma_E W_{ik} \quad (3.21)$$

where U_i , V_k and W_{ik} are independently generated according to a standard normal distribution. This is a particular case for a single stratum of the stratified cross-classified population model introduced in equation (8) of Skinner (2015), where the fixed and random effects are allowed to depend on the strata. Model (3.21) is an analysis of variance model with two crossed random factors and without repetition. Let “ E_m ” denote the expectation with respect to the model (3.21) and “ E_p ” denote the expectation with respect to the CCS design. For each simplified variance estimator $\hat{V}_{\text{SIMP}i}$, $i = 1, 2, 3$, the relative bias RB under the model and under the sampling design is defined by

$$\text{RB}_{m,p} \left[\hat{V}_{\text{SIMP}i}(\hat{t}_Y) \right] = \frac{E_m \left\{ E_p \left[\hat{V}_{\text{SIMP}i}(\hat{t}_Y) \right] - V_{\text{CCS}}(\hat{t}_Y) \right\}}{E_m \left[V_{\text{CCS}}(\hat{t}_Y) \right]}. \quad (3.22)$$

In the SI^2 case, these relative biases are of the form

$$\text{RB}_{m,p} \left[\hat{V}_{\text{SIMP}i}(\hat{t}_Y) \right] = -1/(1 + A_i) \quad (3.23)$$

for $i = 1$ and 2 and

$$\text{RB}_{m,p} \left[\hat{V}_{\text{SIMP}3}(\hat{t}_Y) \right] = 1/(1 + A_3) \quad (3.24)$$

for some positive constant A_i , $i = 1, 2, 3$, depending on σ_M , σ_D , σ_E and n_M , N_M , n_D and N_D , see equations (3.25)-(3.27). Equations (3.23) and (3.24) imply that the two first simplified variance estimators are negatively biased while the third one is positively biased. Using the notations $r_M = \sigma_M^2/\sigma_E^2$, $r_D = \sigma_D^2/\sigma_E^2$, $f_M = n_M/N_M$

and $f_D = n_D/N_D$, we have

$$A_1 = \frac{1 - f_M}{1 - f_D} \frac{n_D r_M + 1}{n_M r_D + f_M}, \quad (3.25)$$

$$A_2 = \frac{1 - f_D}{1 - f_M} \frac{n_M r_D + 1}{n_D r_M + f_D}, \quad (3.26)$$

$$A_3 = \frac{n_D r_M + f_D}{1 - f_D} + \frac{n_M r_D + f_M}{1 - f_M}. \quad (3.27)$$

The bias of \hat{V}_{SIMP1} increases from -1 to 0 when A_1 increases, which occurs in particular when the ratio r_M or the sample size n_D increases. In other words, \hat{V}_{SIMP1} will have a small bias under model (3.21) if the variable of interest contains some maternity unit effect or if the number of sampled days is large enough. Symmetrically, \hat{V}_{SIMP2} will have a small bias under model (3.21) if the variable of interest contains some day effect or if the number of sampled maternity units is large enough. The bias of \hat{V}_{SIMP3} decreases from 1 to 0 when A_3 increases, which occurs in particular when r_M or r_D increases, or when n_M or n_D increases. In other words, \hat{V}_{SIMP3} will have a small bias under model (3.21) if the variable of interest contains some maternity unit or some day effect, or if the number of sampled days or the number of sampled maternity units is large enough. The simulation study in section 4 supports these results, and confirms that the variance tends to be underestimated with \hat{V}_{SIMP1} or \hat{V}_{SIMP2} , and overestimated with \hat{V}_{SIMP3} .

3.4 A central-limit theorem

To produce confidence intervals with appropriate asymptotic coverage, it is of interest to state a central-limit theorem (CLT) for CCS. Roughly speaking, Theorem 1 below states that if the HT-estimator follows a CLT under both sampling designs p_D and p_M , then the HT-estimator also follows a CLT under CCS. It is derived almost directly from Theorem 2 in Chen and Rao (2007), and the proof is therefore omitted.

Theorem 1. *Suppose that assumptions (H1)-(H4) hold. Suppose that*

H5: $\sigma_1^{-1}V_1 \rightarrow_{\mathcal{L}} \mathcal{N}(0, 1)$, where $\rightarrow_{\mathcal{L}}$ stands for the convergence in distribution under the sampling-design, with

$$V_1 = \frac{1}{N} \left(\sum_{i \in S_M} \frac{Y_{i\bullet}}{\pi_i^M} - \sum_{i \in U_M} Y_{i\bullet} \right) \quad \text{and} \quad \sigma_1^2 = V(V_1) \quad (3.28)$$

where $Y_{i\bullet} = \sum_{k \in U_D} Y_{ik}$.

H6: $\sup_t |P(\sigma_2^{-1}U_1 \leq t | S_M) - \Phi(t)| = o_p(1)$, where Φ is the cumulative distribution function of the standard normal distribution, and where

$$U_1 = \frac{1}{N} \sum_{i \in S_M} \frac{1}{\pi_i^M} (\hat{Y}_{i\bullet} - Y_{i\bullet}) \quad \text{and} \quad \sigma_2^2 = V(U_1 | S_M). \quad (3.29)$$

H7: $\sigma_1^2/\sigma_2^2 \rightarrow_P \gamma^2$, where \rightarrow_P stands for the convergence in probability under the sampling-design.

Then

$$\frac{N^{-1}(\hat{t}_Y - t_Y)}{\sqrt{\sigma_1^2 + \sigma_2^2}} \rightarrow_{\mathcal{L}} \mathcal{N}(0, 1). \quad (3.30)$$

For illustration, we consider the particular case when p_D and p_M are both SI designs. Suppose that (H2)-(H4) hold, and that (H1) is strengthened to

H1b: There exists $\delta > 0$ and some constants α_1 and α_2 such that

$$\forall k \in U_D \quad \frac{1}{N_M} \sum_{i \in U_M} Y_{ik}^{2+\delta} \leq \alpha_1, \quad \text{and} \quad \forall i \in U_M \quad \frac{1}{N_D} \sum_{k \in U_D} Y_{ik}^{2+\delta} \leq \alpha_2.$$

Then by using the CLT in Hajek (1961), the assumption (H5) can be shown to hold. By mimicking the proof of Lemma 2 in Chen and Rao (1997), the assumption (H6) can be shown to hold as well.

4 Simulations

In this Section, two artificial populations are first generated using the superpopulation model (3.21). In Section 4.1, CCS is compared with two stage sampling

in terms of variance, which illustrates the results in Section 2.3. A Monte Carlo experiment is then presented in Section 4.2, and the variance estimators introduced in Section 3 are compared for the estimation of a total. Some attention is paid to the issue of negative values for the unbiased variance estimator. In Section 4.3, two other populations with two variables of interest for each are generated. We focus on variance estimation for a ratio, making use of the variance estimators introduced in Section 3 with estimated linearized variables instead of the variable of interest. The results from Tables 1 and 2 are readily reproducible using the R code provided in the supplementary materials of the present paper.

4.1 Comparison with two-stage sampling

Two populations are generated according to model (3.21), with $N_M = 1000$ maternity units and $N_D = 1000$ days for each population, and with $\mu = 200$ and $\sigma_E = 5$. Equal random effects standard deviations $\sigma_M = \sigma_D = 5$ are used for population 1, while we use $\sigma_M = 0.5$ and $\sigma_D = 5$ for population 2. For each population, the SI² sampling design is used, with sample sizes, n_M and n_D , equal to 5, 10, 100 and 500. The ratios V_{MD}/V_{CCS} between the variance under two-stage sampling and the variance under CCS are computed, and plotted as a percentage on Figure 2. A ratio smaller than 100 indicates that two-stage sampling is more accurate than CCS, which holds true in all cases considered in our experiment.

The ratio increases with n_D and decreases when n_M increases. Also, it can be observed that the ratio decreases with σ_M . This impact of the maternity unit effect is noticeable, and illustrates the substantial loss in accuracy induced by using a CCS instead of a two-stage sampling design if the maternity unit effect is small. Similar conclusions could be derived when computing the ratio V_{DM}/V_{CCS} .

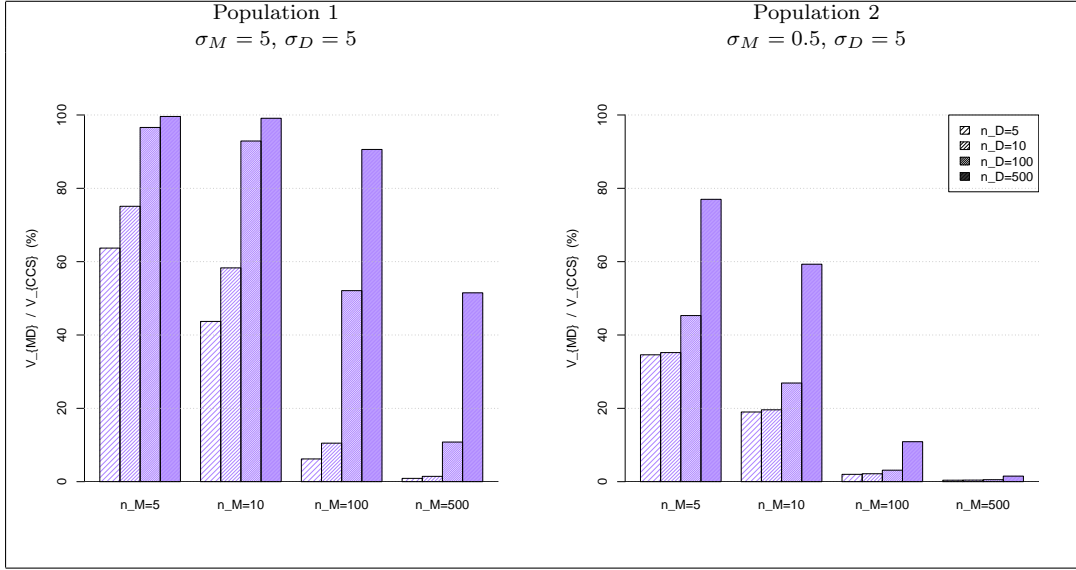


Figure 2: V_{MD}/V_{CCS} (%) for population 1 (left panel) and population 2 (right panel)

4.2 Variance estimation for a total

We consider the two artificial populations generated as described in Section 4.1. For each population, the SI^2 sampling design is used, with sample sizes equal to 5, 10, 100 and 500, and the sample selection is repeated $B = 10,000$ times. For each sample $b = 1, \dots, B$, we compute the estimate $\hat{t}_Y^{(b)}$ of the total t_Y . The unbiased variance estimator $\hat{V}^{(b)}$ and the simplified variance estimators $\hat{V}_{SIMP1}^{(b)}$, $\hat{V}_{SIMP2}^{(b)}$, $\hat{V}_{SIMP3}^{(b)}$ are also computed for $\hat{t}_Y^{(b)}$.

For each variance estimator \hat{V} , we compute 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 true variance V was approximated through an independent set of 50,000 simulations. The number ($\#NEG$) of negative variance estimators $\hat{V}^{(b)}$ is also computed.

The results are reported in Table 1. The variance estimator \hat{V} is almost unbiased in all situations, as expected. However, this variance estimator is prone to negative values with small sample sizes when the value of σ_M and/or the value of σ_D is small

as compared to σ_E . The problem vanishes when the sample sizes increase. We now turn to the simplified variance estimators. The relative bias of \hat{V}_{SIMP1} decreases when n_D increases or when n_M decreases, and when σ_M increases or when σ_D decreases. This supports the findings in Section 3.3. Symmetrical conclusions are drawn for the relative bias of \hat{V}_{SIMP2} . Turning to \hat{V}_{SIMP3} , we note that the relative bias decreases when either σ_M or σ_D increases. This variance estimator is therefore advisable in all cases but those where there is no maternity unit nor day effect.

n_M	5	10	10	100	500	5	10	10	100	500
n_D	5	10	100	100	500	5	10	100	100	500
σ_M	5					50				
σ_D	5					5				
RB _{MC} (\hat{V})	1	-1	2	0	-0	1	-1	1	0	0
#NEG	6	0	0	0	0	0	0	0	0	0
RB _{MC} (\hat{V}_{SIMP1})	-43	-47	-6	-49	-49	-	-2	1	-1	-1
RB _{MC} (\hat{V}_{SIMP2})	-46	-50	-91	-51	-51	-99	-99	-100	-99	-99
RB _{MC} (\hat{V}_{SIMP3})	11	3	2	1	-0	1	-0	1	0	0
σ_M	0.5					0.5				
σ_D	5					0.5				
RB _{MC} (\hat{V})	1	-1	0	1	-1	1	-1	2	-0	-0
#NEG	91	0	0	0	0	1393	298	0	0	0
RB _{MC} (\hat{V}_{SIMP1})	-82	-90	-81	-98	-99	-4	-9	-3	-34	-47
RB _{MC} (\hat{V}_{SIMP2})	-1	-2	-10	-0	-2	-5	-10	-52	-36	-49
RB _{MC} (\hat{V}_{SIMP3})	18	8	9	2	-0	90	81	45	29	4

Table 1: Comparison between variance estimators for a total

4.3 Variance estimation for a ratio

We now consider variance estimation for a ratio. Two populations are generated with $N_M = 1000$ maternity units and $N_D = 1000$ days. In each population, two count variables are generated so as to mimic the data encountered in the ELFE survey. More precisely, we first generate an auxiliary variable Z_{ik} according to model (3.21) with $\mu = 200$, $\sigma_E = \sigma_D = 5$, and $\sigma_M = 5$ or 50. The first variable of interest X_{ik} is generated according to a Poisson distribution with parameter Z_{ik} . The second variable of interest Y_{ik} is generated according to a binomial distribution with parameters X_{ik} and p_{ik} . We consider two cases: (i) equal probabilities with

$p_{ik} = 0.3$; (ii) unequal probabilities with $\text{logit}(p_{ik}) = \beta Z_{ik}$, where β was chosen so that the average probability is approximately 0.3. Note that Y_{ik} follows a Poisson distribution with parameter $p_{ik}Z_{ik}$.

The reason for this generating process is that some variable of interest X_{ik} , like the number of births in the ELFE survey, may contain some maternity unit and/or day effect which is reflected in the way Z_{ik} is generated. On the other hand, some maternity unit and/or day effect may also be contained in some other variable of interest Y_{ik} , like the number of births per caesarean. Such effects may be either similar to those for X_{ik} like with pattern (i), or may occur differently like with pattern (ii).

For each population, the SI^2 sampling design is used, with sample sizes equal to 5, 10, 100 and 500, and the sample selection is repeated $B = 10,000$ times. For each sample $b = 1, \dots, B$, we compute the substitution estimator $\hat{R}^{(b)} = \hat{t}_Y^{(b)}/\hat{t}_X^{(b)}$ of the ratio $R = t_Y/t_X$. The variance estimator $\hat{V}^{(b)}$ and the simplified variance estimators $\hat{V}_{\text{SIMP1}}^{(b)}$, $\hat{V}_{\text{SIMP2}}^{(b)}$, $\hat{V}_{\text{SIMP3}}^{(b)}$ are also computed for $\hat{t}_Y^{(b)}$, where the variable of interest Y_{ik} is replaced with the estimated linearized variable of the ratio.

The results are reported in Table 2. The variance estimator \hat{V} is almost unbiased in all situations, as expected, but is prone to negative values even when the maternity unit or day effect is small. We now turn to the relative bias for the simplified variance estimators. With pattern (i), the situation is much different from that when a total is estimated, since the relative bias of \hat{V}_{SIMP3} is much larger than for the other two simplified estimators. This can be explained as follows: when the probabilities p_{ik} are uniform, both Y_{ik} and X_{ik} contain the same maternity unit and day effect, but these effects wear off in the linearized variable. Whatever the values of σ_M and σ_D are, the situation is therefore comparable to that observed in the bottom right cell of Table 1. With pattern (ii), the probabilities p_{ik} depend on i and k , leading potentially to some remaining maternity unit and/or day effect in the linearized

variable. In such situation, which seems more realistic in practice, the relative bias of \hat{V}_{SIMP1} and \hat{V}_{SIMP2} increase when σ_M or σ_D increase, while the relative bias of \hat{V}_{SIMP3} decreases.

	n_M	5	10	10	100	500	5	10	10	100	500
	n_D	5	10	100	100	500	5	10	100	100	500
	σ_M	5					50				
	σ_D	5					5				
Case (i)	$RB_{MC}(\hat{V})$	-0	-1	-1	0	-0	-2	-1	-1	0	1
	#NEG	1645	484	14	0	0	1656	499	12	0	0
$p_{ik} = 0.3$	$RB_{MC}(\hat{V}_{SIMP1})$	-1	-1	-2	-10	-37	-1	-1	-1	-8	-32
	$RB_{MC}(\hat{V}_{SIMP2})$	0	-2	-10	-8	-30	-2	-1	-9	-8	-31
	$RB_{MC}(\hat{V}_{SIMP3})$	99	96	89	82	33	97	98	90	84	37
Case (ii)	$RB_{MC}(\hat{V})$	0	-1	2	0	-0	-4	-3	-1	-0	0
	#NEG	1351	235	0	0	0	67	0	0	0	0
$p_{ik} = \frac{e^{\beta Z_{ik}}}{1+e^{\beta Z_{ik}}}$	$RB_{MC}(\hat{V}_{SIMP1})$	-7	-13	-4	-39	-48	-5	-4	-1	-1	-1
	$RB_{MC}(\hat{V}_{SIMP2})$	-6	-14	-61	-40	-49	-87	-93	-99	-98	-99
	$RB_{MC}(\hat{V}_{SIMP3})$	87	73	35	22	3	8	3	-0	0	0

Table 2: Comparison between variance estimators for a ratio

5 Application to the ELFE survey

ELFE is the first longitudinal study of its kind in France, tracking children from birth to adulthood (Pirus et al., 2010). This cohort comprises more than 18,000 children whose parents consented to their inclusion. The population of inference consists of babies born during 2011 in French maternity units, excluding very premature infants. This is a two-dimensional population with 544 maternity units as spatial units and 365 days as time units. The crossing of one day and one maternity unit represents a cluster of infants.

The sample is obtained by CCS, where days and maternity units are selected independently with selected families surveyed shortly after birth in 320 metropolitan maternity units and during 25 days for one year. The population of maternity units is divided into five strata of equal size. The allocation per stratum is proportional to the number of deliveries recorded in 2008. The sample selection for maternity units is stratified systematic sampling, which can be approximated by stratified simple

random sampling (STSI). The sample selection of days is not actually random, due to logistic constraints. A number of $n_D = 25$ days is selected during 4 waves, each wave covering a season. It may be approximated by STSI, with four strata associated to the four seasons of 2011. The sample sizes inside strata are provided in Tables 3 and 4.

Strata g	Strata size N_{Mg}	Sample size n_{Mg}
1	108	21
2	108	41
3	109	55
4	108	80
5	111	90
Total	544	287

Table 3: Population and sample strata sizes for the maternity units design p_M .

Strata h	Strata size N_{Dh}	Sample size n_{Dh}
1	91	4
2	91	6
3	91	7
4	92	8
Total	365	25

Table 4: Population and sample strata sizes for the days design p_D .

In this Section, we aim at illustrating the results previously obtained on a real data set. Some aspects of the ELFE survey, like the non-response issue or the calibration step, deserve a specific attention but are beyond the scope of the present paper and are therefore not considered. In particular, the ELFE survey is prone to several levels of non-response, since some sampled maternity units and some families refused to participate either for some specific days or for the whole period. In the present study, the sample of respondents is viewed as the original sample and in particular, we consider only the 287 maternity units that participate during the 25 days of survey. The calibration step is not taken into account. The results below are meant

to illustrate our theoretical results, but are not intended for use in other contexts. We consider seven count variables from the ELFE survey. Some of them depend on the characteristics of the maternity units (e.g., the spatial location), like the variable indicating whether the mother is followed by a midwife. Others are related to the days of the survey, like the variable indicating whether the birth occurred by caesarean. For each variable, the estimated total \hat{t}_Y from equation (2.1), the estimated variance $\hat{V}(\hat{t}_Y)$ from equation (3.2) and the three simplified estimators are given in the upper part of Table 5. Similar indicators are given in the bottom part of Table 5 for ratios, when the totals of the variables of interest are divided by the total number of births.

	Birth	Born by Caesarean	Twins	Born within marriage	Mother followed by a midwife	Mother aged between 18 and 25 years	Primiparous mother	Immigrant mother
\hat{t}_Y	362924	33873	10187	160283	42337	43238	162316	44169
$\hat{V}(\hat{t}_Y)$	7.6E+07	1.5E+07	5.3E+05	2.0E+07	3.9E+06	2.6E+06	1.5E+07	3.6E+06
RD (\hat{V}_{SIMP1})	-63.7 %	-95.5 %	-63.5 %	-64.6 %	-13.2 %	-49.7 %	-46.5 %	-58.2 %
RD (\hat{V}_{SIMP2})	-31.1 %	-1.9 %	-13.3 %	-29.7 %	-76.3 %	-35.2 %	-41.4 %	-33.4 %
RD (\hat{V}_{SIMP3})	5.2 %	2.6 %	23.2 %	5.7 %	10.5 %	15.1 %	12.2 %	8.4 %
\hat{R}	1.00	0.09	0.03	0.44	0.12	0.12	0.45	0.12
$\hat{V}(\hat{R})$		7.9E-05	2.8E-06	2.4E-05	2.5E-05	1.2E-05	3.0E-05	1.6E-05
RD (\hat{V}_{SIMP1})		-96.2 %	-51.0 %	-31.0 %	-7.9 %	-40.2 %	-69.3 %	-49.2 %
RD (\hat{V}_{SIMP2})		-0.4 %	-17.0 %	-44.7 %	-80.5 %	-35.5 %	-5.0 %	-37.5 %
RD (\hat{V}_{SIMP3})		3.4 %	31.9 %	24.3 %	11.5 %	24.3 %	25.7 %	13.3 %

Table 5: Variance estimates of estimated total and ratio on some ELFE variables

The relative difference RD between \hat{V}_{SIMP} and the unbiased estimator \hat{V} is

$$RD = \frac{\hat{V}_{\text{SIMP}}(\hat{t}_{Y^*}) - \hat{V}(\hat{t}_{Y^*})}{\hat{V}(\hat{t}_{Y^*})}.$$

Different behaviours may be observed for the variables of interest, depending on the maternity unit/day effect. For instance, the variable indicating whether the birth occurred by caesarean contains an important day effect, and the RD of \hat{V}_{SIMP2} is therefore small while that of \hat{V}_{SIMP1} is large. Symmetrically, the variable indicating whether the mother is followed by a midwife contains a small day effect as compared to the maternity unit effect, and the RD of \hat{V}_{SIMP2} is therefore large while that of \hat{V}_{SIMP1} is small. Also, we note that the RD of \hat{V}_{SIMP3} is relatively stable for all variables when estimating a total, which is an important feature in favour of this

third simplified estimator. We note however that the absolute RD of \hat{V}_{SIMP_3} can be large when estimating a ratio, which confirms the simulation results.

6 Conclusion

The present paper derives some general estimation theory for the cross-classified sampling design which was used in the recent ELFE survey on childhood. The issue of possibly negative variance estimates may arise even in case of simple random sampling without replacement. Alternative estimators to the usual Horvitz-Thompson and Yates-Grundy variance estimators are thus proposed, and proved to be non-negative under the usual Sen-Yates-Grundy conditions. The relative bias of the proposed variance estimators is derived for a superpopulation model. The behavior of these estimators is also investigated for totals and ratios on simulated data and on data extracted from the ELFE survey. Among the proposals, one variance estimator that leads to a slight overestimation of the variance in many cases, appears to be advisable.

Despite the present results and the recent paper by Skinner (2015), the cross-classified sampling design still deserves some attention. In particular, the treatment of non-response and the calibration problem should also be taken into account, and is currently under investigation.

7 Bibliography

Cardot, H., and Goga, C. and Lardin, P. (2013). Uniform convergence and asymptotic confidence bands for model-assisted estimators of the mean of sampled functional data. *Electronic Journal of Statistics*, 7, 562-596.

Chen, J. and Rao, J.N.K. (2007). Asymptotic normality under two-phase sampling

designs. *Statistica Sinica*, 17, 1047-1064.

Dalén, J. and Ohlsson, E. (1995). Variance Estimation in the Swedish Consumer Price Index. *Journal of Business & Economic Statistics*, 13, No.3, 347-356.

Hajek, J. (1961). Some extensions of the Wald-Wolfowitz-Noether theorem. *The Annals of Mathematical Statistics*, 32, 506-523.

Ohlsson, E. (1996). Cross-Classified Sampling. *Journal of Official Statistics*, 12, No.3, 241-251.

Pirus, C., Bois, C., Dufourg, M.N., Lanoë, J.L., Vandentorren, S., Leridon, H. and the Elfe team (2010). Constructing a Cohort: Experience with the French Elfe Project. *Population*, 65, No.4, 637-670.

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

Sen, A.R. (1953). On the estimate of the variance in sampling with varying probabilities. *Journal of the Indian Society of Agricultural Statistics*, 5, 119-127.

Skinner, C.J. (2015). Cross-classified sampling: some estimation theory. *Statistics and Probability Letters*, 104, 163-168.

Vos, J. W. E. (1964). Sampling in space and time. *Review of the International Statistical Institute*, 32, No. 3, 226-241.

Yates, F. and Grundy, P.M. (1953). Selection without replacement from within strata with probability proportional to size. *Journal of the Royal Statistical Society B*, **15**, 235-261.

8 Appendix

Proof of equations (3.10)-(3.12)

We can rewrite

$$V_1(\hat{t}_Y) = \sum_{k \in U_D} \frac{V(\hat{Y}_{\bullet k})}{\pi_k^D} + \sum_{k \neq l \in U_D} \frac{\pi_{kl}^D}{\pi_k^D \pi_l^D} \text{Cov}(\hat{Y}_{\bullet k}, \hat{Y}_{\bullet l}). \quad (8.1)$$

We have

$$V(\hat{Y}_{\bullet k}) = \sum_{i \in U_M} (1 - \pi_i^M) \frac{(Y_{ik})^2}{\pi_i^M} + \sum_{i \neq j \in U_M} \frac{\pi_{ij}^M - \pi_i^M \pi_j^M}{\pi_i^M \pi_j^M} Y_{ik} Y_{jk}. \quad (8.2)$$

From assumptions (H1), (H2) (H3) and Cauchy-Schwarz inequality, there exists some constant C such that for any $k \in U_D$,

$$V(\hat{Y}_{\bullet k}) \leq C \frac{N_M^2}{n_M}. \quad (8.3)$$

Also, from the Cauchy-Schwarz inequality, there exists some constant C such that for any $k \neq l \in U_D$:

$$\text{Cov}(\hat{Y}_{\bullet k}, \hat{Y}_{\bullet l}) \leq C \frac{N_M^2}{n_M}. \quad (8.4)$$

From equation (8.3) and assumption (H2), the first term in the right hand sum of (8.1) is $O(N_D^2 N_M^2 n_M^{-1} n_D^{-1})$. From equation (8.4) and assumptions (H2) and (H3), the absolute value of the second term in the RHS of (8.1) is $O(N_D^2 N_M^2 n_M^{-1})$. Therefore, there exists some constant C such that

$$V_1(\hat{t}_Y) \leq C \frac{N_D^2 N_M^2}{n_M}. \quad (8.5)$$

We can prove similarly that there exists some constant C such that

$$V_2(\hat{t}_Y) \leq C \frac{N_D^2 N_M^2}{n_D}. \quad (8.6)$$

From equation (2.13), the term $V_3(\hat{t}_Y)$ may be split into four terms according to the intersection of $\{i, j\}$ and $\{k, l\}$. From assumptions (H1)-(H3), it is easily shown that the absolute value of each of these four terms is $O(N_D^2 N_M^2 n_M^{-1} n_D^{-1})$. Therefore, there exists some constant C such that

$$V_3(\hat{t}_Y) \leq C \frac{N_D^2 N_M^2}{n_M n_D}. \quad (8.7)$$

Equations (3.10)-(3.12) follow immediately from equations (8.5)-(8.7) and assumption (H4).