



Component of Statistics Canada
Catalogue no. 12-001-X Business Survey Methods Division

Article

A single frame multiplicity estimator for multiple frame surveys

by Fulvia Mecatti

December, 2007



 Statistics Canada Statistique Canada

Canada 

A single frame multiplicity estimator for multiple frame surveys

Fulvia Mecatti¹

Abstract

Multiple Frame Surveys were originally proposed to foster cost savings on the basis of an *optimality* approach. As surveys on *special*, *rare* and *difficult-to-sample* populations are becoming more prominent, a single list of population units to be used as a sampling frame is often unavailable in sampling practice. In recent literature multiple frame designs have been put forward in order to increase population coverage, to improve response rates and to capture differences and subgroups. Alternative approaches to multiple frame estimation have appeared, all of them relying upon the virtual partition of the set of the available overlapping frames into disjointed domains. Hence the correct classification of sampled units into the domains is required for practical applications. In this paper a multiple frame estimator is proposed using a *multiplicity* approach. Multiplicity estimators require less information about unit domain membership hence they are insensitive to misclassification. Moreover the proposed estimator is analytically simple so that it is easy to implement and its exact variance is given. Empirical results from an extensive simulation study comparing the multiplicity estimator with major competitors are also provided.

Key Words: Difficult-to-Sample populations; Dual frame survey; Misclassification; Raking ratio; Variance estimation.

1. Introduction

In classic finite population sampling a basic hypothesis is the availability of a unique and complete list of units forming the target population to be used as a sampling frame. In some cases a set of two or more lists is available for survey purposes. The general case of $Q \geq 2$ lists, singularly partial and possibly overlapping, is known as *Multiple Frame Survey*. Multiple frame surveys were originally introduced (Hartley 1974) as a device for reducing survey costs by achieving the same precision as a customary unique-frame survey. In modern sampling practice, as surveys of *special*, *rare* and *difficult-to-sample populations* are becoming more common (Kalton and Anderson 1986; Sudman and Kalton 1986; Sudman, Sirken and Cowan 1988) it is often the case that a unique list of units does not exist and the population size N is an unknown parameter to be estimated. Recent literature considers multiple frame surveys with the main aim of increasing population coverage, of improving response rates and of capturing differences and subgroups more accurately (Iachan and Dennis 1993; Carlson and Hall 1994; Haines and Pollock 1998; Eurostat 2000). In a recent paper Lohr and Rao (2006) stated: "As the U.S., Canada, and other nations grow in diversity, different sampling frames may better capture subgroups of the population. [...] We anticipate that modular sampling designs using multiple frames will be widely used in the future". A contemporary application could be found in web surveys: the population coverage can be improved and the bias due to the features of the site used for data collection can be reduced by using two or more independent web sites simultaneously. Since the

same unit can visit more than one site involved in the survey, the sites overlap configuring a multiple frame framework.

Estimation in multiple frame surveys, as first developed by Hartley (1962, 1974), is based on the virtual partition of the population (*i.e.*, the unknown union of the Q overlapping frame) into $2^Q - 1$ disjointed *domains* (*i.e.*, the mutually exclusive intersections of frames). Hence the total Y of a study variable y , taken as the parameter to be estimated, is expressed as a sum of *domain totals*. Sample data from the Q frames are used to produce estimates for the domain totals. Estimated domain totals are finally combined to provide estimation for the population total Y . A number of estimators have been developed according to alternative approaches to multiple frame estimation (see Section 2). Since all estimators appearing in literature rely on the partition into the domains as mentioned above, the correct identification of the domain membership of each sampled unit is required for their practical application. This is a strong assumption that may not always be true in practice, as argued for instance in Lohr and Rao (2006). Indeed this implies that every sampled unit should be questioned on both the survey value and on its membership to each frame involved in the survey, in order to be able to correctly classify them into the domains. In addition to the natural risk of misclassification there might also be a risk connected with *confidentiality* and with the *sensitivity* of units to the frame membership which could both increase the rate of non-response and affect the estimator precision. This situation could apply, for instance, when surveying *sensitive* characteristics (private behaviours, addictions ...) or when sampling *elusive* populations (illegal immigrants,

1. Fulvia Mecatti, Department of Statistics, University of Milan-Bicocca, Via Bicocca degli Arcimboldi, 8. Ed. U7, 20126 Milan, Italy. E-mail: fulvia.mecatti@unimib.it.

ex-prisoners, patients...). In the present paper a different approach to estimation in multiple frame surveys is adopted. The concept of unit *multiplicity*, corresponding to the *number* of frames to which units belong, is proposed in alternative to the existing approaches based on the domain membership, *i.e.*, to which frames units belong. An unbiased estimator, naturally insensitive to domain misclassification and applying to any number of frames, is presented. The proposed multiplicity estimator has a simple analytical structure so that it can be easily implemented, while its exact variance is given in a closed form and hence readily estimated for any sample size.

In Section 2 an overall discussion of the main contributions to multiple frame estimation is presented in a unified view and the necessary notation is introduced. In Section 3 a multiplicity estimator is proposed and variance estimation is analysed. An extensive simulation study comparing the proposed estimator with major competitors is presented in Section 4.

2. Optimum, pseudo-optimum and single frame estimation

Although literature has mostly dealt with the dual frame case ($Q = 2$), a general theoretical framework for multiple frame surveys ($Q \geq 2$) has been recently provided in Lohr and Rao (2006). By using their multiple frame notation, different estimation approaches are briefly reviewed and the available estimators are presented in a unified way which highlights their dependency on the domain membership of the sampled units.

Let $A_1 \cdots A_q \cdots A_Q$ be a collection of $Q \geq 2$ overlapping frames, the union of which offers a population coverage adequate for survey objectives. Let the index sets K be the subsets of the range of the frame index $q = 1 \cdots Q$. For every index set $K \subseteq \{1 \cdots q \cdots Q\}$ a domain is defined as the set $D_K = (\cap_{q \in K} A_q) \cap (\cap_{q \notin K} A_q^c)$, where c denotes complementation.

Let the *domain membership indicator* be the indicator $\delta_i(K)$ taking value 1 if unit i is included in domain D_K and 0 otherwise. The estimating total Y over the (unknown) union of the Q overlapping frames is then expressed as a sum over the set of $2^Q - 1$ disjointed domains

$$Y = \sum_{i \in \cup_q A_q} y_i = \sum_K \sum_{i \in \cup_q A_q} \delta_i(K) y_i \tag{1}$$

Let s_q be a sample selected from frame A_q under a given design, independently for $q = 1 \cdots Q$. A general expression of a multiple frame estimator based on the domain classification is then

$$\hat{Y} = \sum_K \sum_{q \in K} \sum_{i \in s_q} w_i^{(q)} \delta_i(K) y_i \tag{2}$$

Note that when a unbiased estimator for the total Y is given, an estimator for the population size N is also given by simply substituting sample values y_i by 1's.

Estimators available in literature result from setting weights $w_i^{(q)}$ in (2) according to three main approaches. Since multiple frame surveys were originally put forward with the aim of fostering cost savings by achieving equal or greater precision than a customary unique-frame survey, an *optimum* approach was first suggested by using optimum weights $w_{i, \text{opt}}^{(q)}$ in (2), *i.e.*, by minimizing the estimator variance (Hartley 1962, 1974; Lund 1968; Fuller and Burmeister 1972). Optimum estimators have optimal theoretical properties (Skinner 1991; Lohr and Rao 2000) but present practical problems due mainly to their complexity (explicit though complex formulae for optimum weights $w_{i, \text{opt}}^{(q)}$ with any number of frames are given in Lohr and Rao 2006, Section 3). Moreover, optimum weights depend on unknown population covariances so that they must be estimated from sample data. This is both computationally complex and affects optimality since the extra variability in estimating the covariances leads to larger mean square errors (Lohr and Rao 2006, Section 7).

In order to improve the applicability, a *single frame* (SF) approach has been proposed by using *fixed weights* which ensure design-unbiasedness. For simple random sampling in every frame, the SF estimator is given by substituting weights $w_i^{(q)}$ in (2) with $w_{i, \text{SF}}^{(q)} = w^{(K)} = (\sum_{q \in K} f_q)^{-1}$ where $f_q = n_q / N_q$ denotes the frame sampling fraction (Bankier 1986; Kalton and Anderson 1986; Skinner 1991; Skinner, Holmes and Holt 1994). Since fixed weights usually differ from optimum weights, the SF estimator is generally less efficient than an optimum estimator (Lohr and Rao 2000). Finally a *pseudo-optimum approach* was proposed (Skinner and Rao 1996; Lohr and Rao 2000) in order to achieve both a wider applicability than optimum estimators and to improve efficiency compared with the SF approach. A pseudo-maximum likelihood (PML) estimator for multiple frame surveys is given by substituting in (2): $w_{i, \text{PML}}^{(q)} = w^{(K)} = \hat{N}_K / \sum_{q \in K} \sum_{i \in s_q} \delta_i(K) = \hat{N}_K / n_K$ where the estimated domain sizes \hat{N}_K are the solution of a system of non linear equations. Although complex to implement for practical applications (an iterative linear approximation of \hat{N}_K under simple random sampling is given in Lohr and Rao 2006, Section 4.1) the PML estimator retains good theoretical properties from the optimum approach.

Note that formula (2) involves the domain membership indicator $\delta_i(K)$; hence optimum, pseudo-optimum and SF estimators apply only if the correct classification of sample data into the $2^Q - 1$ domains is accomplished.

In the next Section a multiple frame estimator is presented on the basis of a single frame *multiplicity* approach which does not require domain classification.

3. The single frame multiplicity estimator

The notion of multiplicity was first introduced in connection with Network Sampling (Casady and Sirken 1980; Sirken 2004). It is also a tool of the Generalized Weight Share Method (Lavallée 2002; 2007) as well as of the Center Sampling estimation theory (Mecatti 2004) since center sampling and multiple frame surveys are equivalent under certain conditions. In Lohr and Rao (2006), the multiplicity of domain D_K is defined as the cardinality of the index set K . Since domains are mutually exclusive, multiplicity is also a characteristic of every population unit, being the *number* of frames in which each unit is included among the Q involved in the survey.

Let m_i be the multiplicity of unit i . Note that unit multiplicity may be collected simply by asking sampled units *how many frames they belong to*.

Since clearly $\sum_q \sum_{i \in A_q} y_i = \sum_{i \in \cup_q A_q} m_i y_i$, it follows that

$$Y = \sum_{q=1}^Q \sum_{i \in A_q} y_i m_i^{-1}. \tag{3}$$

Notice that expression (3), which involves exclusively sums over the frames, represents a practical advantage with respect to equation (1). In fact the domains provide a virtual (unknown) partition of the population while the sample selection is actually performed in the Q overlapping frames. This leads to a SF multiplicity estimator as given by

$$\hat{Y}_M = \sum_{q=1}^Q \sum_{i \in s_q} w_i^{(q)} y_i m_i^{-1} \tag{4}$$

with fixed weights $w_i^{(q)}$ ensuring, for instance, design-unbiasedness. For simple random sampling of every frame we have $w_i^{(q)} = f_q^{-1}, \forall i \in s_q$.

Unlike the optimum, PML and SF estimators discussed in Section 2, estimator (4) does not involve the sample membership indicator and it is very simple to implement in practical applications. Furthermore, it is to be noted that for simple random sampling of every frame, the sampled values in multiplicity estimator (4) are weighted by $(f_q m_i)^{-1}$, *i.e.*, by a *specific* frame coefficient; vice versa, in the SF estimator sampled values are weighted by $w_i^{(K)} = (\sum_{q \in K} f_q)^{-1}$, *i.e.*, by an *average* coefficient over the frames involved in each domain. As a consequence \hat{Y}_M is expected to be more accurate than the SF estimator, as confirmed by simulation results. Moreover, owing to its Horvitz-Thompson structure, the exact variance of \hat{Y}_M can be derived in closed form. For

simple random sampling of every frame the estimator variance is given by

$$V(\hat{Y}_M) = \sum_{q=1}^Q \frac{N_q - n_q}{n_q(N_q - 1)} \left[N_q \sum_{i \in A_q} y_i^2 m_i^{-2} - \left(\sum_{i \in A_q} y_i m_i^{-1} \right)^2 \right]. \tag{5}$$

An unbiased variance estimator for simple random sampling of every frame is then

$$\hat{v}(\hat{Y}_M) = \sum_{q=1}^Q \frac{N_q(N_q - n_q)}{n_q^2(N_q - 1)} \left[N_q \sum_{i \in s_q} y_i^2 m_i^{-2} - f_q^{-1} \left(\sum_{i \in s_q} y_i m_i^{-1} \right)^2 \right]. \tag{6}$$

The performance of the multiplicity estimator for finite sample sizes has been empirically studied under simple random sampling and compared to major competitors in a simulation study.

4. Simulation study

Several simulation results concerning dual frame estimators have appeared in literature (Bankier 1986; Skinner and Rao 1996; Lohr and Rao 2000). In the general case of $Q \geq 2$ frames, Lohr and Rao (2006) extensively investigated the empirical mean squared errors of a set of eight estimators under optimum, pseudo-optimum and single frame approaches, in a three-frame framework under a two-stage design. Their results suggest that optimum estimators are theoretically optimal but in practice the extra variability in estimating optimum weights leads to larger mean squared errors. Hence the PML estimator appears as the best performer in terms of empirical relative efficiency. Furthermore, their study regarded a case of about 10% of sampled units misclassified into the domains and more research on the effects of misclassification on the estimator performances is recommended.

In the present study pseudo-optimum and single frame estimators are compared with the multiplicity estimator (4), with three main objectives:

- i) to investigate empirical conditions in which the multiplicity estimator results more efficient than the SF estimator (Section 4.2);
- ii) to consider the raking ratio correction to known frame sizes N_q as already proposed in order to improve efficiency of the SF estimator (Section 4.3);
- iii) to explore the effects of increasing rates of misclassification upon the empirical properties of the PML and SF estimators (simple and raked)

versus the natural insensitivity of the multiplicity estimator (Section 4.4).

4.1 Implementation

The simulation study was performed in an artificial three-frame setup and implemented as follows. N population pseudo-values y_i are generated from a Gamma distribution. Some preliminary simulations indicated that both increasing values of the population size N and different values for the Gamma parameters (leading to an asymmetrical and almost symmetrical shape) do not produce significant differences in the pattern of the relative performance of the estimators considered. The study was then conducted by setting $N = 1,200$ and by generating from a Gamma distribution with parameters of 1.5 and 2. Every pseudo-value y_i is randomly assigned to the $Q = 3$ frames according to 3 independent Bernoulli trials with probability $\alpha_q = N_q / N$, $q = 1, 2, 3$. Different scenarios regarding both frame coverage and frame overlapping result from different choices for α_q , under the two constraints: a) $\sum_q \alpha_q \geq 1$ in order to ensure that the 3 frames cover the entire population and b) the 3 frames are non-empty. In some cases, the desired frame overlapping was produced by fixing the ratio N_K / N of the population units included in each domain.

Chosen a set of sampling fractions $f_q = n_q / N_q$, $q = 1, 2, 3$, a simple random sampling is selected independently from every frame, iteratively for 10,000 simulation runs. For a given estimator, say \hat{Y} , the collection of values $\{\hat{Y}_p, p = 1 \dots 10,000\}$ is assumed as its monte carlo distribution and the empirical mean $E_{mc}(\hat{Y}) = \sum_p \hat{Y}_p / 10,000$ and the empirical mean squared error $MSE_{mc}(\hat{Y}) = \sum_p [\hat{Y}_p - Y]^2 / 10,000$ are calculated. The monte carlo error is controlled by only accepting simulations giving empirical relative bias $RB_{mc}(\hat{Y}) = 100 \cdot |E_{mc}(\hat{Y}) - Y| / Y$ less than 1.5% for those estimators known to be unbiased. Furthermore, by using the exact variance of the multiplicity estimator as given by (6), simulations ensure $|MSE_{mc}(\hat{Y}_M) - V(\hat{Y}_M)| \leq 0.03$. Several different scenarios have been investigated by combining different levels of frame coverage, of frame overlapping and of sampling disproportion, leading to 29 simulated populations. In Figure 1 the simulated populations are represented as points in the plane formed by the two main simulation parameters, namely the (total) frame coverage on the horizontal axis (as given by $\sum_q \alpha_q$) and the sampling disproportion on the vertical axis, i.e., the dispersion among the sampling fractions f_q as measured by $\sum_q \sum_q |f_q - f_q| / 3^2$. The different shape of populations/points in Figure 1 indicates different levels of overlapping, namely the total rate of population units classified into the four overlapping domains.

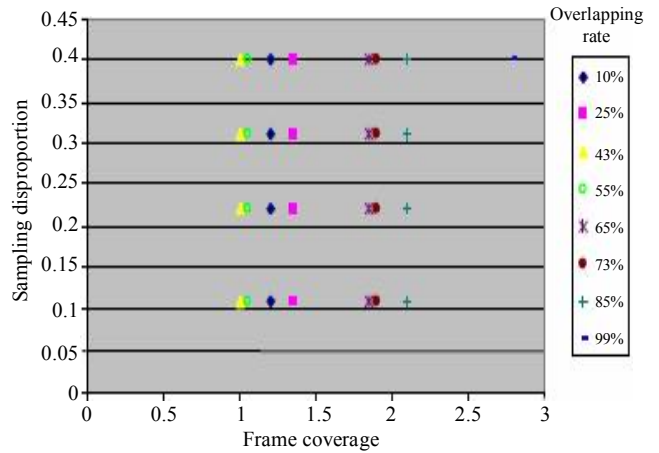


Figure 1 Simulated populations

4.2 Multiplicity versus simple single frame estimation

As noted in Section 3, the multiplicity estimator involves specific frame weights whereas the SF estimator is based on average coefficients. As a consequence the two estimators coincide for constant sample fraction $f_q = f$ in every frame, i.e., for proportionate sampling, and they offer different estimates for disproportionate sampling. Simulation results provide empirical evidence that the multiplicity estimator is more accurate than the simple SF estimator. Estimator \hat{Y}_M is shown to be more efficient in all the cases explored except in one extreme case in which the three frames are almost complete and hence the total overlapping is close to 100%. Neglecting this single case, efficiency gains of \hat{Y}_M over the SF estimator, as measured by a customary empirical efficiency ratio (see Table 1), range from 5% to 48%, and are never less than 26% in half of the simulations. Efficiency of the multiplicity estimator over the SF estimator increases as the sampling disproportion increases (see Table 2) whereas it has resulted as being essentially independent with respect to increasing levels of frame coverage and overlapping.

Table 1
Empirical efficiency ratio of \hat{Y}_M versus SF estimator: Elementary statistics over 28 simulated populations

| average | Max | min | median | 75 th quantile |
|---------|------|------|--------|---------------------------|
| 0.7425 | 0.95 | 0.52 | 0.74 | 0.89 |

Table 2
Empirical efficiency ratio of \hat{Y}_M versus the SF estimator for increasing levels of sampling disproportion

| | | | | |
|---|------|------|------|------|
| Sampling Disproportion | 0.11 | 0.22 | 0.31 | 0.40 |
| Empirical efficiency ratio averaged for different levels of frames coverage/overlapping | 0.92 | 0.81 | 0.68 | 0.57 |

4.3 Raking ratio adjustment

It has been suggested that the raking ratio adjustment using the known frame sizes N_q (Bankier 1986) be used in order to improve efficiency of the simple SF estimator. Theoretical and empirical results that have already appeared in literature confirm that the raking ratio SF estimator (SFrak) can be considerably more efficient than the simple SF estimator (Skinner 1991; Lohr and Rao 2000, 2006; Mecatti 2005).

In order to adjust the multiplicity estimator via raking ratio, knowledge of the domain membership of sampled units has to be assumed. By using this additional, though redundant, information \hat{Y}_M may be rewritten as

$$\hat{Y}_M = \sum_K \sum_{q \in K} (|K| f_q)^{-1} \sum_{i \in s_q} \delta_i(K) y_i \quad (7)$$

where $|K|$ indicates the number of frames involved in domain D_K and it equals unit multiplicity m_i for all $i \in D_K$. Setting the initial weights at $h_{Kq}^{(0)} = (|K| f_q)^{-1}$, the t^{th} iteration of the raking ratio multiplicity estimator (Mrak) is obtained by substituting the following raked weights in (7)

$$h_{Kq}^{(t)} = \begin{cases} \frac{N_q h_{Kq}^{(t-1)}}{\sum_{q \in K} h_{Kq}^{(t-1)} n_q} & \text{if } q \in K \\ h_{Kq}^{(t-1)} & \text{if } q \notin K \end{cases}$$

where $q = Q$ if t is a multiple of Q otherwise $q = t \bmod(Q)$, for $t = 1, 2, \dots$ until convergence.

Simulations regarded different levels of frames coverage combined with different sets of sampling fractions, leading to increasing sampling disproportion.

Empirical results show that Mrak is more efficient than SFrak in 38% of cases explored and it is equally or less efficient in the remaining cases. Efficiency gains range from 3% to 74% and occur for low levels of frame coverage. For increasing frame coverage (and hence increasing overlapping) Mrak estimator is superior to SFrak estimator for high sampling disproportion only. In the other cases, namely for increasing frame coverage/overlapping combined with low to medium sampling disproportion, Mrak can be considerably less efficient than SFrak (see Table 3 for the ten indicative cases) and also severely biased. Thus empirical results suggest that the raking ratio adjustment has better effects under a single frame approach than under a multiplicity approach, although there are conditions in which the latter is still superior. With this respect more research is needed. Particularly, since the raking ratio procedure is in fact a special case of calibration (Deville and Särndal 1992; Deville, Särndal and Sautory 1993), potential improvements might follow by applying the more general calibration to estimator \hat{Y}_M . Calibration of the multiplicity

estimator, as viewed as a particular case of the Generalized Weight Share Method, is outlined in Lavallée (2002, 2007).

Table 3 Efficiency of Mrak versus SFrak: Ten indicative simulation runs

| Frame coverage $\alpha_q = N_q/N$ | | | Sampling fractions $f_q = n_q/N_q$ | | | Empirical efficiency ratio Mrak versus SFrak |
|--------------------------------------|------|------|---------------------------------------|------|------|---|
| 0.60 | 0.60 | 0.60 | 0.01 | 0.95 | 0.15 | 0.26 |
| 0.35 | 0.35 | 0.35 | 0.80 | 0.20 | 0.50 | 0.54 |
| 0.85 | 0.85 | 0.85 | 0.01 | 0.95 | 0.15 | 0.71 |
| 0.35 | 0.40 | 0.50 | 0.70 | 0.80 | 0.60 | 0.96 |
| 0.85 | 0.85 | 0.85 | 0.80 | 0.20 | 0.50 | 1.01 |
| 0.60 | 0.60 | 0.60 | 0.70 | 0.80 | 0.60 | 1.09 |
| 0.80 | 0.50 | 0.35 | 0.01 | 0.95 | 0.15 | 1.22 |
| 0.35 | 0.40 | 0.50 | 0.01 | 0.95 | 0.15 | 1.63 |
| 0.70 | 0.05 | 0.95 | 0.70 | 0.80 | 0.60 | 2.09 |
| 0.70 | 0.05 | 0.95 | 0.80 | 0.20 | 0.50 | 5.79 |

4.4 Misclassification

The aim of the final part of the simulation study is to investigate the sensitivity of the pseudo-optimum (PML) and single frame estimators (simple and raked) to increasing levels of misclassification of sampled units into the domains, with respect to the structural insensitivity of the proposed multiplicity estimator. For a chosen rate of misclassification, the desired number of sampled units to be in exactly classified is taken from the domain with the largest size and randomly assigned to the remaining domains, independently for each frame.

Tables 4 and 5 show elementary statistics summarizing simulation results in the case of exact classification and in the case of slight misclassification equal to 1% of sampled units. Note that for exact classification all the estimators appear unbiased (or nearly unbiased). As regards efficiency, according to other simulation results (Lohr and Rao 2006) SFrak and PML estimators show similar performances. As expected, for exact classification they are more efficient than \hat{Y}_M in all the cases explored (except for two isolated cases) as a consequence of the different amount of information used in the estimation process. However the SF (simple and raked) and PML estimators tends to become biased and less efficient than \hat{Y}_M in presence of just a small amount of misclassification.

Table 4 Relative bias in case of 1% of misclassification: Elementary statistics over the 29 simulated populations

| (absolute) RB _{mc} 1% of sampled units misclassified | Average | Min | Max | Median | 75 th quantile |
|--|---------|------|------|--------|------------------------------|
| \hat{Y}_M | 0 | 0 | 0 | 0 | 0 |
| SF | 2.5880 | 0.83 | 7.02 | 2.65 | 2.13 |
| SFrak | 1.7632 | 0.73 | 4 | 1.97 | 1.65 |
| PML | 2.7352 | 0.23 | 4.67 | 3.46 | 2.87 |

Table 5
Empirical efficiency ratio of \hat{Y}_M versus SF and PML estimators: Elementary statistics over the 29 simulated populations for exact classification and for slight misclassification

| Empirical efficiency ratio | Average | Min | Max | Median | 75 th quantile |
|-----------------------------|---------|------|------|--------|---------------------------|
| <i>Exact classification</i> | | | | | |
| SFrak | 1.43 | 0.69 | 3.21 | 1.51 | 1.28 |
| PML | 1.41 | 0.72 | 3.30 | 1.47 | 1.25 |
| <i>1% misclassification</i> | | | | | |
| SF | 0.39 | 0.13 | 0.71 | 0.54 | 0.34 |
| SFrak | 0.78 | 0.13 | 1.98 | 0.95 | 0.74 |
| PML | 0.77 | 0.14 | 1.94 | 0.98 | 0.70 |

Finally we focused on the case of maximum efficiency of SF, SFrak and PML over \hat{Y}_M for exact classification, namely the case of high frame overlapping/coverage and low sampling disproportion. In this set up, increasing rates of misclassification of sampled units into domains (from 0 to 50%) were investigated. Table 6 and 7 show respectively the relative bias and the efficiency ratio of \hat{Y}_M versus SF, SFrak and PML estimators, for increasing levels of misclassification. It is to be noticed that although the negative effects of misclassification are rapid and severe for all the competitors, the PML estimator emerges as the least affected.

As a conclusion the proposed multiplicity estimator, besides being simple, is recommended when the risk of (even slight) misclassification of sampled units into the domains is a concrete possibility.

Table 6 (absolute) Relative bias for increasing rate of misclassification

| % misclassification | \hat{Y}_M | SF | SFrak | PML |
|---------------------|-------------|-------|-------|------|
| 0 | 0 | 0 | 0 | ≈ 0 |
| 1% | 0 | 2.57 | 1.38 | 4.3 |
| 5% | 0 | 13.57 | 7.15 | 2.75 |
| 10% | 0 | 17.80 | 14.14 | 4.56 |
| 20% | 0 | 25 | 25 | 6 |
| 50% | 0 | 144 | 68 | 39 |

Table 7
Empirical efficiency ratio of \hat{Y}_M versus SF, SFrak and PML estimators for increasing rate of misclassification

| % misclassification | \hat{Y}_M versus SF | \hat{Y}_M versus SFrak | \hat{Y}_M versus PML |
|---------------------|-----------------------|--------------------------|------------------------|
| 0 | 0.640 | 3.210 | 3.300 |
| 1% | 0.260 | 1.040 | 1.100 |
| 5% | 0.020 | 0.060 | 0.370 |
| 10% | 0.010 | 0.020 | 0.150 |
| 20% | 0.004 | 0.004 | 0.080 |
| 50% | ≈ 0 | 0.001 | 0.006 |

Acknowledgements

This work was partially supported by a grant from the Italian Ministry of University and Research. The author wishes to thank Jon N.K. Rao for the useful discussion and resourceful advice. Thanks are also due to the editor, two associate editors and two anonymous referees for their constructive comments and suggestions.

References

- Bankier, M.D. (1986). Estimators based in several stratified samples with applications to multiple frame surveys. *Journal of American Statistical Association*, 81, 1074-1079.
- Casady, R.J., and Sirken, M.G. (1980). A multiplicity estimator for multiple frame sampling. *Proceedings of the Survey Research Methods Section*, American Statistical Association, 601-605.
- Carlson, B.L., and Hall, J.W. (1994). Weighting sample data when multiple sample frames are used. *Proceedings of the Survey Research Methods Section*, American Statistical Association, 882-887.
- Deville, J.C. and Särndal, C.-E. (1992). Calibration Estimators in Survey Sampling. *Journal of American Statistical Association*, 87, 376-382.
- Deville, J.-C., Särndal, C.-E. and Sautory, O. (1993). Generalized raking procedures in survey sampling. *Journal of American Statistical Association*, 88, 1013-1020.
- Eurostat (2000). *Push and Pull Factors of International Migration*. Country Report-Italy, 3/2000/E/n.5, Bruxelles: European Communities Printing Office.
- Fuller, W.A., and Burmeister, L.F. (1972). Estimators of samples selected from two overlapping frames. *Proceedings of the Social Statistics Sections*, American Statistical Association, 245-249.
- Haines, D.E., and Pollock, K.H. (1998). Combining multiple frame to estimate population size and totals. *Survey Methodology*, 24, 79-88.
- Hartley, H.O. (1962). Multiple frame surveys. *Proceedings of the Social Statistics Sections*, American Statistical Association, 203-206.
- Hartley, H.O. (1974). Multiple frame methodology and selected applications. *Sankhyā*, C, 36, 99-118.
- Iachan, R., and Dennis, M.L. (1993). A multiple frame approach to sampling the homeless and transient population. *Journal of Official Statistics*, 9, 747-764.
- Kalton, G., and Anderson, D.W. (1986). Sampling rare populations. *Journal of Royal Statistical Society*, A, 149, 65-82.
- Lavallée, P. (2002). *Le sondage indirect ou la méthode généralisée du partage des poids*. Editions de l'Université de Bruxelles, Editions Ellipses.
- Lavallée, P. (2007). *Indirect Sampling*. New York: Springer.
- Lohr, S.L., and Rao, J.N.K. (2000). Inference from dual frame surveys. *Journal of American Statistical Association*, 95, 271-280.

- Lohr, S., and Rao, J.N.K. (2006). Multiple frame surveys: Point estimation and inference. *Journal of American Statistical Association*, 101, 1019-1030.
- Lund, R.E. (1968). Estimators in multiple frame surveys. *Proceedings of the Social Statistics Sections*, American Statistical Association, 282-288.
- Mecatti, F. (2004). Center sampling: A strategy for surveying difficult-to-sample populations. *Proceedings of the Statistics Canada Symposium*.
- Mecatti, F. (2005). Single frame estimation in multiple frame survey. *Proceedings of the Statistics Canada Symposium*.
- Sirken, M.G. (2004). Network sample surveys of rare and elusive populations: A historical review. *Proceedings of Statistics Canada Symposium, Keynote Address*.
- Skinner, C.J. (1991). On the efficiency of raking ratio estimator for multiple frame surveys. *Journal of American Statistical Association*, 86, 779-784.
- Skinner, C.J., Holmes, D.J. and Holt, D. (1994). Multiple frame sampling for multivariate stratification. *International Statistical review*, 62, 333-347.
- Skinner, C.J., and Rao, J.N.K. (1996). Estimation in dual frame surveys with complex designs. *Journal of American Statistical Association*, 91, 349-356.
- Sudman, S., and Kalton, G. (1986). New developments in the sampling of special populations. *Annual Review of Sociology*, 12, 401-429.
- Sudman, S., Sirken, M.G. and Cowan, C.D. (1988). Sampling rare and elusive populations. *Science*, 2, 991-996.