

Improved Variance Estimation for Balanced Samples Drawn Via the Cube Method

F. Jay Breidt^a, G. Chauvet^{b,*}

^a*Department of Statistics, Colorado State University, Fort Collins, CO 80523-1877, USA*

^b*Ecole Nationale de la Statistique et de l'Analyse de l'Information, Campus de Ker Lann,
35170 Bruz, France*

Abstract

The Cube method proposed by Deville and Tillé (2004) enables the selection of balanced samples: that is, samples such that the Horvitz-Thompson estimators of auxiliary variables match the known totals of those variables. As an exact balanced sampling design often does not exist, the Cube method generally proceeds in two steps: a “flight phase” in which exact balance is maintained, and a “landing phase” in which the final sample is selected while respecting the balance conditions as closely as possible. Deville and Tillé (2005) derive a variance approximation for balanced sampling that takes account of the flight phase only, whereas the landing phase can prove to add non-negligible variance. This paper uses a martingale difference representation of the cube method to construct an efficient simulation-based method for calculating approximate second-order inclusion probabilities. The approximation enables nearly unbiased variance estimation, where the bias is primarily due to the limited number of simulations. In a Monte Carlo study, the proposed method has significantly less bias than the standard variance estimator, leading to improved confidence interval coverage.

Keywords: Balanced sampling, Cube method, Inclusion probabilities, Martingale

1. Introduction

A balanced sampling design has the attractive feature that the Horvitz-Thompson estimators of the totals for auxiliary variables, called balancing variables, exactly

*. Corresponding author

Email addresses: jbreidt@stat.colostate.edu (F. Jay Breidt), chauvet@ensai.fr (G. Chauvet)

Preprint submitted to Journal of Statistical Planning and Inference

June 17, 2010

match the known totals. Deville and Tillé (2004) introduced the *cube method*, which enables the selection of exact balanced samples if such samples may be found, or approximate balanced samples otherwise. The cube method proceeds in two phases: the *flight phase*, in which balancing constraints are maintained exactly, and the *landing phase*, in which the balancing constraints are successively relaxed until the sample is completed. In the case of balanced sampling with maximum entropy, Deville and Tillé (2005) proposed various variance approximations and associated variance estimators, and compared their performance through a set of simulations.

The variance approximations proposed by Deville and Tillé (2005) do not take into account the whole sampling process, as they ignore the contributions from the landing phase. In fact, the landing phase generates an additional term of variance that may not be negligible in many cases. For example, if the balancing variables have high predictive power for the variable of interest, then the contribution from the landing phase may represent a major part of the variance. In the extreme case when the variable of interest may be perfectly explained through a linear regression model including the balancing constraints as independent variables, the variance due to the flight phase is 0. As pointed out by a referee, the variance associated with the landing phase may also be important when the number of balancing constraints is high as compared to the (expected) sample size. The use of linear mixed models as a way to relax balancing constraints is considered in Breidt and Chauvet (2010).

In this paper, we propose a simulation-based approximation of the variance-covariance matrix of the sampling design. The use of simulation-based methods to approximate inclusion probabilities is discussed in Fattorini (2006) and Thompson and Wu (2008). The simulation method relies on a martingale difference representation of the cube method. The proposed simulation-based approximation is shown to be particularly efficient as compared to a naive simulation-based approximation, which ignores the martingale difference structure. The approximation enables the computation of a variance estimator that takes into account the whole sampling process. We show that the resulting variance estimator has significantly less bias than the standard variance estimator from Deville and Tillé (2005), leading to confidence intervals with better coverage rates. The paper is organized as follows. After defining notation in Section 2, we present the simulation-based variance approximation and the resulting variance estimator in Section 3. The martingale-difference based variance approximation is compared to the naive simulation-based approximation in Section 4.1, and the

variance estimator is compared to the standard estimator of Deville and Tillé (2005) in Section 4.2 through a set of simulations. A brief discussion follows in Section 5.

2. Balanced Sampling

Let U denote a finite, labeled population of size N . Let S denote a random sample selected from U by means of a sampling design $p(\cdot)$. Let $\pi_k = \Pr[k \in S]$ denote the inclusion probability of unit k , and $\pi_{kl} = \Pr[k, l \in S]$ denote the second-order inclusion probability. Write $\boldsymbol{\pi} = (\pi_1, \dots, \pi_k, \dots, \pi_N)'$ for the vector of inclusion probabilities and $\mathbf{I}(S) = (I_1, \dots, I_k, \dots, I_N)'$ for the vector of sample membership indicators, where $I_k = 1$ if $k \in S$ and 0 otherwise. The sampling design is assumed to be of fixed size. This implies that $\sum_{k \in U} \pi_k = n = \sum_{k \in U} I_k$, where n denotes the sample size.

The Horvitz-Thompson (1952) estimator,

$$\hat{t}_{z\pi} = \sum_{k \in U} \frac{\mathbf{z}_k}{\pi_k} I_k, \quad (1)$$

estimates without bias the finite population total $t_z = \sum_{k \in U} \mathbf{z}_k$ of the (vector) variable \mathbf{z}_k . The variance of the Horvitz-Thompson estimator may be obtained by means of the Yates-Grundy (1953) formula:

$$\text{Var}(\hat{t}_{z\pi}) = -\frac{1}{2} \sum_{k, l \in U: k \neq l} \Delta_{kl} \left(\frac{\mathbf{z}_k}{\pi_k} - \frac{\mathbf{z}_l}{\pi_l} \right) \left(\frac{\mathbf{z}_k}{\pi_k} - \frac{\mathbf{z}_l}{\pi_l} \right)', \quad (2)$$

where $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$. In the particular case of a scalar variable of interest $\mathbf{z}_k = y_k$, we obtain

$$\text{Var}(\hat{t}_{y\pi}) = -\frac{1}{2} \sum_{k, l \in U: k \neq l} \Delta_{kl} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2. \quad (3)$$

We assume that a vector $\mathbf{x}_k = (x_{1k}, \dots, x_{qk})'$ of q auxiliary variables is known at the design stage for each unit k in the population. Such auxiliary information allows for possible gains in efficiency through the selection of balanced samples, drawn via the *cube method* (Deville and Tillé, 2004). The sampling design $p(\cdot)$ is said to be *balanced* on variables \mathbf{x} if the equations

$$\hat{t}_{\mathbf{x}\pi} = t_{\mathbf{x}} \quad (4)$$

hold exactly. That is, the Horvitz-Thompson estimator for the auxiliary vector \mathbf{x} exactly matches the known vector of totals. The variables \mathbf{x}_k are called the balancing variables. For example, the condition of fixed sample size is met if the inclusion probability π_k belongs to the vector \mathbf{x}_k of balancing variables. If the equations (4) are satisfied, the variance of the Horvitz-Thompson estimator is zero for any linear combination of the balancing variables \mathbf{x} . As an exact balanced sample generally cannot be found, the cube method enables the selection of approximately balanced samples, proceeding in a flight phase and a landing phase as noted above.

Several implementations of the cube method have been proposed in the literature. Algorithm 1 given in the appendix covers both the flight phase and the landing phase. The flight phase given in Algorithm 1 corresponds to the general flight phase as proposed by Deville and Tillé (2005) and described by Chauvet and Tillé (2006, Algorithm 1). The variance estimation strategy proposed in this paper may still be applied if the very fast implementation of the flight phase proposed by Chauvet and Tillé (2006, Algorithm 2) is used instead. The landing phase given in Algorithm 1 proceeds by successively relaxing the balancing constraints. A landing phase by means of an enumerative algorithm could alternatively be used (see Tillé, 2006, p. 163), but the variance estimation strategy proposed in §3.1 of this paper would no longer be valid, since the martingale difference representation of the cube method would not be maintained during the whole sampling process. Also, the landing phase given in Algorithm 1 has the advantage of permitting the selection of a balanced sample in a reasonable amount of time, even if the number of balancing variables is large (see Tillé, 2006, p. 164).

Algorithm 1 proceeds in steps $t = 0, 1, \dots, T$ from $\pi_0(S) = \pi$ to $\pi_T(S) = \mathbf{I}(S)$, the final sample. At each step, one or more coordinates of $\pi_t(S)$ are randomly rounded to 0 or 1, and remain there forever. Let $X = [\mathbf{x}'_k]_{k \in U}$, the $N \times q$ matrix of auxiliary variables. During the flight phase, the balancing equations remain exactly respected, that is:

$$X' \text{diag}\{\pi_k^{-1}\}_{k \in U} \pi_t(S) = X' \mathbf{1}_N$$

where $\mathbf{1}_N$ denotes the $N \times 1$ vector of ones. When exact balance is no longer possible, the constraints are relaxed successively in the landing phase. It is thus necessary to order the balancing variables with respect to their relative importance, since the last balancing variables are removed first. We assume that this

ordering has been done prior to implementation of Algorithm 1.

If the sample S is selected by means of Algorithm 1, the variance of the Horvitz-Thompson estimator may be decomposed as follows:

$$\begin{aligned} \text{Var}(\hat{t}_{y\pi}) &= \text{Var}\left(\mathbb{E}\left[\hat{t}_{y\pi} \mid \boldsymbol{\pi}_{T(F)}(S)\right]\right) + \mathbb{E}\left[\text{Var}\left(\hat{t}_{y\pi} \mid \boldsymbol{\pi}_{T(F)}(S)\right)\right] \\ &= V_F(\hat{t}_{y\pi}) + V_L(\hat{t}_{y\pi}), \end{aligned} \quad (5)$$

where $\boldsymbol{\pi}_{T(F)}(S)$ stands for the sequence $\boldsymbol{\pi}_t(S)$ at the end of the flight phase, $t = T(F)$. The term $V_F(\hat{t}_{y\pi})$ denotes the variance due to the flight phase, whereas the term $V_L(\hat{t}_{y\pi})$ denotes the variance due to the landing phase.

The cube method belongs to the family of martingale algorithms (Tillé, 2006, p. 32). That is, the sequence $\{\boldsymbol{\pi}_t(S)\}_{t=0,\dots,T}$, where T denotes the stopping time of the algorithm, defines a discrete time martingale.

3. Variance estimation for balanced sampling

The variance of the Horvitz-Thompson estimator is given by formula (3). Consequently, an unbiased variance estimator is given by

$$v(\hat{t}_{y\pi}) = -\frac{1}{2} \sum_{k,l \in S: k \neq l} \frac{\Delta_{kl}}{\pi_{kl}} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2 \quad (6)$$

if all $\pi_{kl} = \Delta_{kl} + \pi_k \pi_l$ are strictly positive. Second-order inclusion probabilities are, however, usually difficult to compute for a general balanced sampling design. A particular case of balancing variables for which these probabilities may be calculated is when $\mathbf{x}_k = x_k = \pi_k$, that is, when the constraint of fixed sample size is the only balancing constraint. If this balanced sampling design is performed with maximum entropy, Deville (2000) and Matei and Tillé (2005) proposed an algorithm that enables the computation of the exact second-order inclusion probabilities. In general, we resort to simulation for computation of Δ_{kl} .

3.1. Simulation-based variance estimation

We consider two variance estimators, based on two different simulation-based approximations of the design variance-covariance matrix,

$$\Delta = [\Delta_{kl}]_{k,l \in U} = \text{Var}(\mathbf{I}(S)).$$

Since

$$\mathbb{E}[(\mathbf{I}(S) - \boldsymbol{\pi})(\mathbf{I}(S) - \boldsymbol{\pi})'] = \Delta,$$

an obvious, unbiased simulation-based estimator of Δ is given by

$$\Delta_{SIM} = \frac{1}{C} \sum_{c=1}^C (\mathbf{I}(S_c) - \boldsymbol{\pi})(\mathbf{I}(S_c) - \boldsymbol{\pi})', \quad (7)$$

where $S_1, \dots, S_c, \dots, S_C$ are C independent replicates of the sample. These replicates may be selected by Algorithm 1. As pointed out by a referee, the $\boldsymbol{\pi}$ vector may be replaced by its simulation-based estimator $\boldsymbol{\pi}_{SIM} = \frac{1}{C} \sum_{c=1}^C \mathbf{I}(S_c)$, to get an alternative estimator of Δ given by

$$\Delta_{SIM2} = \frac{1}{C} \sum_{c=1}^C (\mathbf{I}(S_c) - \boldsymbol{\pi}_{SIM})(\mathbf{I}(S_c) - \boldsymbol{\pi}_{SIM})'. \quad (8)$$

The corresponding variance estimator for a given sample S is then obtained by plugging (7) into (6),

$$v_{SIM}(\hat{t}_{y\pi}) = -\frac{1}{2} \sum_{k,l \in S: k \neq l} \frac{\Delta_{SIM,kl}}{\Delta_{SIM,kl} + \pi_k \pi_l} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2. \quad (9)$$

This estimator is no longer exactly unbiased, because the unbiased approximations of the Δ_{kl} enter non-linearly. The estimator also has greater variance than v due to the simulation. Both the bias and the additional variance can be made arbitrarily small for sufficiently large C .

We now propose an alternative simulation-based approximation for Δ that uses the martingale structure of the sampling algorithm. Specifically, note that

$$\mathbf{I}(S) = \boldsymbol{\pi}_T(S) = \boldsymbol{\pi} + \sum_{t=1}^T \boldsymbol{\delta}_t(S), \quad (10)$$

where the innovations $\{\boldsymbol{\delta}_t(S)\}_{t=1, \dots, T}$ are given in Algorithm 1. By construction, $\{\boldsymbol{\delta}_t(S)\}$ is a martingale difference (MD) sequence with respect to the sequence of sigma-fields $\mathcal{F}_{t-1} = \sigma(\boldsymbol{\delta}_0(S), \boldsymbol{\delta}_1(S), \dots, \boldsymbol{\delta}_{t-1}(S))$, and so these random vectors

are uncorrelated and have mean zero. Hence,

$$\begin{aligned}
\Delta = \text{Var}(\mathbf{I}(S)) &= \text{Var}(\boldsymbol{\pi}_T(S)) \\
&= \sum_{t=1}^T \text{Var}(\boldsymbol{\delta}_t(S)) \\
&= \sum_{t=1}^T \text{E}[\text{Var}(\boldsymbol{\delta}_t(S) | \mathcal{F}_{t-1})] \\
&= \text{E}\left[\sum_{t=1}^T \lambda_{1t}^*(S) \lambda_{2t}^*(S) \mathbf{u}_t(S) \mathbf{u}_t'(S)\right],
\end{aligned}$$

the last equality following from Step 3 of Algorithm 1. Consequently the Δ matrix is unbiasedly estimated by

$$\Delta_{MD} = \frac{1}{C} \sum_{c=1}^C \sum_{t=1}^T \lambda_{1t}^*(S_c) \lambda_{2t}^*(S_c) \mathbf{u}_t(S_c) \mathbf{u}_t'(S_c), \quad (11)$$

where $S_1, \dots, S_c, \dots, S_C$ are C independent replicates of the sample, selected by Algorithm 1.

The corresponding variance estimator for a given sample S is then obtained by plugging (11) into (6),

$$v_{MD}(\hat{t}_{y\pi}) = -\frac{1}{2} \sum_{k,l \in S: k \neq l} \frac{\Delta_{MD,kl}}{\Delta_{MD,kl} + \pi_k \pi_l} \left(\frac{y_k}{\pi_k} - \frac{y_l}{\pi_l} \right)^2. \quad (12)$$

A key question is whether the MD-based approximation is any better than the standard simulation-based approximation, SIM. Note that

$$\begin{aligned}
(\mathbf{I}(S) - \boldsymbol{\pi})(\mathbf{I}(S) - \boldsymbol{\pi})' &= \sum_{r=1}^T \sum_{t=1}^T \boldsymbol{\delta}_r(S) \boldsymbol{\delta}_t'(S) \\
&= \sum_{t=1}^T \text{E}[\boldsymbol{\delta}_t(S) \boldsymbol{\delta}_t'(S) | \mathcal{F}_{t-1}] \\
&\quad + \sum_{t=1}^T (\boldsymbol{\delta}_t(S) \boldsymbol{\delta}_t'(S) - \text{E}[\boldsymbol{\delta}_t(S) \boldsymbol{\delta}_t'(S) | \mathcal{F}_{t-1}]) \\
&\quad + \sum_{r \neq t} \boldsymbol{\delta}_r(S) \boldsymbol{\delta}_t'(S) \\
&= \sum_{t=1}^T \lambda_{1t}^*(S) \lambda_{2t}^*(S) \mathbf{u}_t(S) \mathbf{u}_t'(S) + \Lambda_1(S) + \Lambda_2(S),
\end{aligned}$$

so that

$$\begin{aligned}\Delta_{SIM} &= \Delta_{MD} + \frac{1}{C} \sum_{c=1}^C \Lambda_1(S_c) + \frac{1}{C} \sum_{c=1}^C \Lambda_2(S_c) \\ &= \Delta_{MD} + \Lambda_1 + \Lambda_2.\end{aligned}\tag{13}$$

Clearly $E[\Lambda_1] = E[\Lambda_2] = 0$. Equation (13) implies, in some sense, that one part of the variability of Δ_{SIM} vanishes in Δ_{MD} , since the two random terms Λ_1 and Λ_2 are omitted. The sign of the covariance between Δ_{MD} and $\Lambda_1 + \Lambda_2$ would be necessary to make this point rigorous; unfortunately we have not found a tractable approach for computation of this covariance. However, our simulation results on a small population suggest that our approximation performs particularly well: see Section 4.1.

3.2. Maximum entropy variance estimation

Another variance approximation is provided by Deville and Tillé (2005). They assume that the balanced sampling is performed with maximum entropy, and that the sampling design is exactly balanced. Then, under an assumption of asymptotic normality of the multivariate Horvitz-Thompson estimator $\hat{t}_{z\pi}$ under Poisson sampling, they derive the following variance approximation:

$$\text{Var}_{DT}(\hat{t}_{y\pi}) = \frac{N}{N-q} \sum_{k \in U} \pi_k (1 - \pi_k) \left(\frac{y_k}{\pi_k} - \frac{y_k^*}{\pi_k} \right)^2,\tag{14}$$

where

$$y_k^* = \mathbf{x}'_k \left(\sum_{l \in U} \pi_l (1 - \pi_l) \frac{\mathbf{x}_l \mathbf{x}'_l}{\pi_l^2} \right)^{-1} \sum_{l \in U} \pi_l (1 - \pi_l) \frac{\mathbf{x}_l y_l}{\pi_l^2}$$

is a weighted prediction of y_k obtained with the q balancing variables \mathbf{x}_k . Other slightly different approximations are proposed in Deville and Tillé (2005), but their simulation results suggest that approximation (14) performs well among variance approximations that may be computed in case of any set of inclusion probabilities. A variance estimator is obtained through a substitution of each total in (14) by its Horvitz-Thompson estimator, using a plug-in principle. The resulting estimator is

$$v_{DT}(\hat{t}_{y\pi}) = \frac{n}{n-q} \sum_{k \in S} (1 - \pi_k) \left(\frac{y_k}{\pi_k} - \frac{\tilde{y}_k^*}{\pi_k} \right)^2,\tag{15}$$

where

$$\tilde{y}_k^* = \mathbf{x}'_k \left(\sum_{l \in S} (1 - \pi_l) \frac{\mathbf{x}_l \mathbf{x}'_l}{\pi_l^2} \right)^{-1} \sum_{l \in S} (1 - \pi_l) \frac{\mathbf{x}_l y_l}{\pi_l^2}.$$

The hypothesis of exact balancing assumed by Deville and Tillé (2005) implies that approximation (14) accounts for the variance due to the flight phase only. Consequently, the variance estimator given in (15) may lead to serious bias in variance estimation if the variance due to the landing phase is appreciable. With the method we propose, the whole sampling process is taken into account in the variance estimation.

4. Results of the numerical studies

4.1. Comparison of the two simulation-based approximations

We first consider the two simulation-based approximations, Δ_{SIM} and Δ_{MD} , for a small finite population in which exact computation of Δ is possible. The population U is of size $N = 10$. Two balanced designs are considered, one of size $n = 3$ and one of size $n = 5$. For each design, the inclusion probability vector π is obtained by drawing z_1, \dots, z_N as independent and identically distributed Uniform(0,0.2) random variables. Then define $\pi_k = nz_k / \sum_{k \in U} z_k$. We consider the sampling design with $x_{k1} = 1$ and $x_{k2} = \pi_k$ as balancing variables. This is not an exact balanced sampling design.

Since the population is of very small size, the complete support of the design can be enumerated, and the design probabilities can be computed exactly. In the case $n = 3$, there are 58 possible samples, with probabilities of selection ranging from 3.12×10^{-4} to 7.92×10^{-2} . In the case $n = 5$, there are 61 possible samples with probabilities of selection ranging from 8.74×10^{-4} to 1.42×10^{-1} . Further, Δ can be computed exactly.

The two simulation-based approximations Δ_{SIM2} and Δ_{MD} given by formulas (8) and (11), respectively, are compared. The results obtained with the simulation-based approximation Δ_{SIM} given by formula (7) are very similar to those obtained with Δ_{SIM2} given by formula (8), and are thus omitted. Following Deville and Tillé (2005), we note that if we approximate Δ by Δ_{approx} , the worst-case variance approximations are given when the ratio

$$\frac{\mathbf{y}' \Delta_{approx} \mathbf{y}}{\mathbf{y}' \Delta \mathbf{y}}$$

is far from one, with extreme values given by

$$\min_{\mathbf{y}} \mathbf{y}' \Delta_{approx} \mathbf{y} \text{ subject to } \mathbf{y}' \Delta \mathbf{y} = 1,$$

and

$$\max_{\mathbf{y}} \mathbf{y}' \Delta_{approx} \mathbf{y} \text{ subject to } \mathbf{y}' \Delta \mathbf{y} = 1.$$

In turn, these optimization problems correspond to choosing the smallest and largest nonnegative general eigenvalues α_{min} and α_{max} , respectively, where the general eigenvalues are the roots of the polynomial $\det[\Delta_{approx} - \alpha \Delta]$. See Deville and Tillé (2005, p. 576) and Harville (1997, p. 562, 581). A good approximation of the Δ matrix should give general eigenvalues α_{min} and α_{max} close to 1.

These general eigenvalues for Δ_{SIM2} and Δ_{MD} , computed with varying numbers of simulations ranging from $C = 50$ to $C = 10,000$, are given in Figure 1 with a logarithmic scale for the eigenvalues. We note that with the proposed MD approximation, α_{min} and α_{max} converge quickly to 1 as the number of simulations grows. That is, the matrix Δ_{MD} converges quickly to Δ . Moreover, the matrix Δ_{MD} systematically outperforms the matrix Δ_{SIM2} , regardless of the number of simulations used.

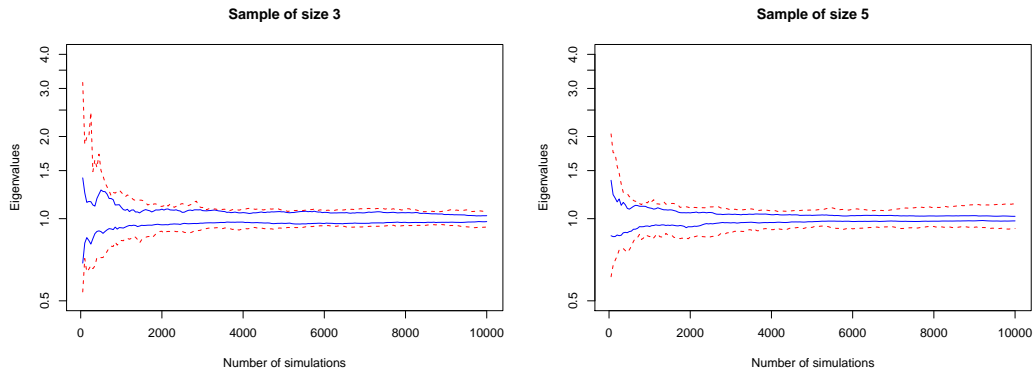


Figure 1: Largest and smallest nonnegative general eigenvalues versus number of simulations, C . The blue, solid curves correspond to the proposed approximation (MD), and the red, dashed curves to the usual simulation-based approximation (SIM). Vertical axis is logarithmic.

4.2. Comparison with the Deville-Tillé variance estimator

Some numerical studies adapted from Deville and Tillé (2005) are considered to assess the proposed MD variance estimator. We consider balanced sampling designs in three different populations. In the two first populations, the inclusion probabilities are generated by using a uniform distribution.

Balanced design 1: The population U_1 is of size $N = 40$. The design considered has unequal probabilities and $n = 15$. For $k \in U_1$, define

$$z_{2k} = k, z_{3k} = k^{-1}, z_{4k} = k^{-2}$$

and let $\bar{z}_i = N^{-1} \sum_{k \in U_1} z_{ik}$, $s_{zi}^2 = (N-1)^{-1} \sum_{k \in U_1} (z_{ik} - \bar{z}_i)^2$ denote the usual empirical mean and variance for $i = 2, 3, 4$. Then $\mathbf{x}'_k = (x_{1k}, x_{2k}, x_{3k}, x_{4k})$, where $x_{1k} = \pi_k$ and

$$x_{ik} = \frac{z_{ik} - \bar{z}_i}{s_{zi}}$$

for $i = 2, 3, 4$. This is not an exact balanced sampling design.

Balanced design 2: The population U_2 is of size $N = 30$. The design considered has unequal probabilities and $n = 10$. Define $x_{k1} = 1$, $x_{k2} = 1$ if $k \in \{1, \dots, 15\}$ and 0 otherwise, $x_{k3} = 1$ if $k \in \{11, \dots, 25\}$ and 0 otherwise, and $x_{k4} = 1$ if $k \in \{1, \dots, 5\} \cup \{21, \dots, 30\}$ and 0 otherwise. This is not an exact balanced sampling design.

Balanced design 3: The population U_3 is of size $N = 45$. The design considered has equal probabilities and $n = 15$. Define $x_{k1} = 1$, $x_{k2} = 1$ if $k \in \{1, \dots, 15\}$ and 0 otherwise, $x_{k3} = 1$ if $k \in \{16, \dots, 30\}$ and 0 otherwise, and $x_{k4} = 1$ if $k \in \{1, \dots, 9\} \cup \{16, \dots, 21\} \cup \{31, \dots, 39\}$ and 0 otherwise. This is an exact balanced sampling design.

In each population, five variables of interest y_1, \dots, y_5 are generated according to the linear regression model

$$y_{jk} = \beta_1 + \beta_2 x_{2k} + \beta_3 x_{3k} + \beta_4 x_{4k} + \sigma_j \epsilon_k, \quad (16)$$

for $j = 1, \dots, 5$. The ϵ_k 's are generated according to a normal distribution with mean 0 and variance 1. For populations U_1 and U_3 , $\beta_1 = \beta_2 = \beta_3 = \beta_4 = 1$. For population U_2 , $\beta_1 = 0$ and $\beta_2 = \beta_3 = \beta_4 = 1$. In each case, the coefficient σ_j was chosen to give a model R^2 (coefficient of determination) approximately equal to 0.1 for y_1 , 0.2 for y_2 , and so on.

In population U_3 , five other variables of interest $y_1^{int}, \dots, y_5^{int}$ are generated according to the linear regression model with interactions

$$y_{jk}^{int} = \gamma_1 + \gamma_2 x_{2k} + \gamma_3 x_{3k} + \gamma_4 x_{4k} + \gamma_5 x_{2k} \times x_{4k} + \gamma_6 x_{3k} \times x_{4k} + \eta_j \epsilon_k, \quad (17)$$

for $j = 1, \dots, 5$. The ϵ_k 's are generated according to a normal distribution with mean 0 and variance 1, and we used $\gamma_1 = \gamma_2 = \gamma_3 = \gamma_4 = \gamma_5 = \gamma_6 = 1$. The coefficient η_j was chosen to give a model R^2 (coefficient of determination) approximately equal to 0.1 for y_1^{int} , 0.2 for y_2^{int} , and so on. These variables of interest are meant to evaluate (to some extent) the performance of the variance estimators when the auxiliary information used is not fully adequate.

Our objective is to estimate the variance of the Horvitz-Thompson estimator of the totals of the five y -variables of interest for populations U_1 and U_2 , and of the ten y -variables of interest for population U_3 . The methods considered included the proposed Martingale Difference variance estimator (MD) and the approximation given by Deville and Tillé (DT). The Δ_{MD} matrix in (11) was obtained from a separate simulation run of $C = 100,000$ samples. The same simulation run was used to compute the variance of the Horvitz-Thompson estimators given in (5), as well as the variance due to the flight phase which is given by the first term in the right-hand side of (5). The importance of the variance due to the flight phase in the overall variance is measured by the relative flight effect (FE) given by

$$FE(\hat{t}_{y\pi}) = 100 \times \frac{V_F(\hat{t}_{y\pi})}{\text{Var}(\hat{t}_{y\pi})}.$$

Since an exact balanced sampling design is considered in population U_3 , there is no landing phase and the variance is only due to the flight phase. Consequently, the FE is omitted in this latter case.

As a measure of bias of a point estimator $\hat{\theta}$ of a parameter θ , we used the Monte Carlo percent relative bias (RB) given by

$$RB_{MC}(\hat{\theta}) = 100 \times \frac{B^{-1} \sum_{b=1}^B \hat{\theta}_{(b)} - \theta}{\theta}$$

where $B = 2,000$ independent samples are selected according to the balanced sampling design, and $\hat{\theta}_{(b)}$ gives the value of the estimator for the b^{th} sample. As a measure of variance of an estimator $\hat{\theta}$ we used the Monte Carlo percent relative

stability (RS) given by

$$RS_{MC}(\hat{\theta}) = 100 \times \frac{\sqrt{B^{-1} \sum_{b=1}^B (\hat{\theta}_{(b)} - \theta)^2}}{\theta}.$$

When $\theta = \text{Var}(\hat{t}_{y\pi})$, we have $\hat{\theta}$ equal to either $v_{MD}(\hat{t}_{y\pi})$ or $v_{DT}(\hat{t}_{y\pi})$. We also assess the coverage of confidence intervals computed using the t distribution with $n-q$ degrees of freedom, where q denotes the number of balancing variables. The results are given in Tables 1, 2 and 3 for populations U_1 , U_2 and U_3 , respectively.

In each of the first two populations, the relative bias of the DT estimator increases as the R^2 of the model increases, whereas the MD estimator remains approximately unbiased. As a consequence, the DT estimator is also outperformed by the simulation-based variance estimator in terms of confidence interval error rates, especially when R^2 is high. Our interpretation is that R^2 measures the relative importance between the variance due to the flight phase and the variance due to the landing phase in the overall variance, which is in agreement with the FE . When R^2 is small, the balancing variables have little explanatory power for the variable y . The variance due to the flight phase is then high (close to the variance obtained in case of a simple fixed-size sampling strategy) and dominates the overall variance. When R^2 becomes larger, the balancing variables have more and more explanatory power for the variable of interest, and the importance of the variance due to the flight phase in the overall variance decreases. The extreme case is when $R^2 = 1$. The variable of interest y is then perfectly explained by the balancing variables. In this case, the variance due to the flight phase is 0. Because only the flight phase is taken into account in the DT variance estimator, it is therefore more and more negatively biased as R^2 increases. This effect is seen in both U_1 and U_2 . The results obtained with the third population are in agreement with this interpretation. In this latter case, the balancing is performed exactly, and the variance due to the landing phase is zero. The two variance estimators then give very similar results, even if the auxiliary information is not fully adequate. The MD variance estimator is less biased, while the DT variance estimator is more stable and gives slightly better results in terms of coverage (better in 12 out of 20 cases).

The MD variance estimator has considerable variability, in terms of percent relative stability, partly due to the small sample size considered in these examples. The DT variance estimator has low variability, but its extreme bias makes

Var.	Method	Nominal one-tailed error rate						% Rel. Bias RB_{MC}	% Rel. Stab. RS_{MC}	% Rel. Fil. Eff. FE
		2.5 %			5 %					
		L	U	$L+U$	L	U	$L+U$			
$n = 5$										
y_1	MD	1.15	0.85	2.00	1.15	1.65	2.80	11.8	656.6	56.5
	DT	3.70	6.85	10.55	7.35	11.75	19.10	-74.1	104.6	
	MD	1.20	0.90	2.10	1.20	2.00	3.20	9.4	581.8	51.0
y_2	DT	4.10	7.40	11.50	7.65	12.40	20.05	-77.3	102.3	
	MD	1.20	1.05	2.25	1.20	2.75	3.95	7.1	522.4	45.6
y_3	DT	4.20	7.80	12.00	8.30	13.05	21.35	-80.4	100.4	
	MD	1.30	1.10	2.40	1.30	3.90	5.20	5.0	479.4	40.3
y_4	DT	4.55	8.10	12.65	9.00	14.05	23.05	-83.3	98.9	
	MD	1.25	1.30	2.55	1.35	4.40	5.75	3.1	454.4	35.2
y_5	DT	4.70	8.60	13.30	9.45	15.50	24.95	-86.1	97.8	
$n = 15$										
y_1	MD	0.80	4.75	5.55	2.15	7.20	9.35	-0.8	37.7	76.2
	DT	2.00	5.65	7.65	4.35	8.90	13.25	-21.0	35.8	
	MD	0.65	5.60	6.25	1.75	8.45	10.20	-0.9	45.8	68.5
y_2	DT	2.55	6.35	8.90	5.30	9.75	15.05	-29.0	39.1	
	MD	0.40	6.10	6.50	1.60	9.35	10.95	-1.1	56.7	60.8
y_3	DT	3.50	7.00	10.50	6.90	10.90	17.80	-37.1	43.9	
	MD	0.35	7.10	7.45	1.30	10.40	11.70	-1.3	69.7	52.9
y_4	DT	4.85	8.10	12.95	8.75	11.60	20.35	-45.3	49.7	
	MD	0.25	8.15	8.40	0.80	12.30	13.10	-1.5	84.2	44.8
y_5	DT	6.80	9.55	16.35	11.40	13.25	24.65	-53.6	56.4	

Table 1: Monte-Carlo error rates in the upper and lower tails, percent relative biases and relative stabilities of two variance estimators for balanced sampling in population U_1

Var.	Method	Nominal one-tailed error rate										% Rel. Bias	% Rel. Stab.	% Rel. Fli. Eff.
		2.5 %					5 %							
		L	U	L + U	L	U	L	U	L + U	L	U			
<i>n</i> = 5														
y_1	MD	7.15	5.40	12.55	7.30	7.75	15.05	3.0	171.0	57.0				
	DT	3.35	4.90	8.25	6.05	11.15	17.20	-63.3	94.1					
y_2	MD	7.50	6.30	13.80	7.55	8.10	15.65	3.4	173.0	51.1				
	DT	3.55	5.00	8.55	6.80	11.55	18.35	-67.1	91.7					
y_3	MD	7.70	7.25	14.95	7.85	8.40	16.25	3.7	175.9	45.5				
	DT	3.95	5.20	9.15	7.55	12.25	19.80	-70.7	90.0					
y_4	MD	7.80	7.35	15.15	8.05	8.55	16.60	3.9	179.5	40.0				
	DT	4.50	5.90	10.40	8.20	12.80	21.00	-74.2	88.9					
y_5	MD	8.75	7.55	16.30	9.00	8.50	17.50	4.2	183.7	34.4				
	DT	4.45	6.55	11.00	9.50	13.65	23.15	-77.8	88.5					
<i>n</i> = 10														
y_1	MD	1.20	7.40	8.60	2.80	10.95	13.75	-1.6	57.0	83.3				
	DT	0.80	7.50	8.30	2.60	11.50	14.10	-12.7	54.4					
y_2	MD	1.65	7.60	9.25	3.25	11.35	14.60	-1.5	57.5	80.6				
	DT	1.05	7.70	8.75	2.90	11.75	14.65	-15.6	53.4					
y_3	MD	1.80	8.15	9.95	3.40	11.15	14.55	-1.4	58.3	77.4				
	DT	1.15	7.65	8.80	3.85	12.15	16.00	-18.9	52.6					
y_4	MD	2.50	8.60	11.10	4.10	11.05	15.15	-1.3	59.6	73.7				
	DT	1.35	8.15	9.50	4.10	12.40	16.50	-22.8	52.0					
y_5	MD	3.05	9.00	12.05	4.80	11.25	16.05	-1.1	61.7	69.3				
	DT	1.85	8.90	10.75	4.95	12.60	17.55	-27.4	51.8					

Table 2: Monte-Carlo error rates in the upper and lower tails, percent relative biases and relative stabilities of two variance estimators for balanced sampling in population U_2

Variable	Method	Nominal one-tailed error rate						RB_{MC}	% Rel. Bias	RS_{MC}
		2.5 %			5 %					
		L	U	$L+U$	L	U	$L+U$			
y_1	MD	2.05	3.20	5.25	4.90	5.60	10.50	0.2	34.4	
	DT	2.20	3.05	5.25	4.85	5.60	10.45	-0.6	33.7	
y_2	MD	1.30	4.65	5.95	3.30	7.45	10.75	-1.8	34.6	
	DT	1.25	4.30	5.55	3.05	7.80	10.85	-2.3	33.4	
y_3	MD	2.00	3.30	5.30	3.35	5.80	9.15	0.2	31.7	
	DT	2.00	3.10	5.10	3.55	5.80	9.35	-2.5	30.9	
y_4	MD	4.30	1.85	6.15	7.45	4.55	12.00	0.4	39.6	
	DT	5.45	1.35	6.80	8.50	3.65	12.15	-1.7	38.5	
y_5	MD	3.65	3.60	7.25	5.75	6.30	12.05	0.3	44.2	
	DT	2.60	2.35	4.95	4.75	4.55	9.30	-0.5	34.0	
y_1^{int}	MD	2.10	3.05	5.15	4.70	5.85	10.55	0.2	34.7	
	DT	2.15	3.05	5.20	4.85	5.55	10.40	-0.7	33.9	
y_2^{int}	MD	1.45	3.80	5.25	3.10	6.90	10.00	-1.6	32.9	
	DT	1.15	4.00	5.15	2.75	7.25	10.00	-1.3	32.8	
y_3^{int}	MD	1.70	4.00	5.70	3.05	6.50	9.55	-0.1	34.9	
	DT	1.70	4.10	5.80	3.25	7.00	10.25	-3.2	34.5	
y_4^{int}	MD	4.70	2.10	6.80	7.80	5.10	12.90	0.8	41.4	
	DT	4.60	2.05	6.65	7.75	4.70	12.45	-1.7	38.2	
y_5^{int}	MD	4.65	2.60	7.25	6.75	5.15	11.90	0.0	44.3	
	DT	2.80	2.05	4.85	5.15	4.35	9.50	-0.9	33.0	

Table 3: Monte-Carlo error rates in the upper and lower tails, percent relative biases and relative stabilities of two variance estimators for balanced sampling in population U_3

this point irrelevant in all cases. In the decomposition of the variance given in (5), the term $V_L(\hat{t}_{y\pi})$ associated with the flight phase is ignored in the DT estimator, which results in a negatively-biased estimator. However, $\text{Var}(\hat{t}_{y\pi} | \pi_{T(F)}(S))$ is highly unstable, since the units remaining at the end of the flight phase as well as their number are random. As a consequence, the MD estimator which helps in tracking the term $V_L(\hat{t}_{y\pi})$ is also highly variable. This variability is needed to achieve approximately correct confidence interval coverage.

5. Conclusion

In this paper, a new variance estimator for balanced sampling, using a martingale difference representation of the cube method, is proposed. This estimator is shown to perform well as compared to a naive simulation-based variance estimator that does not use the martingale structure. Unlike the Deville and Tillé (2005) variance estimator, which ignores sampling variation due to the landing phase, the proposed estimator takes into account the complete sampling process and leads to essentially unbiased variance estimators, and corresponding confidence intervals with proper coverage. Numerical results support our findings.

Acknowledgements

We thank a referee for helpful comments. Breidt's research was supported in part by the US National Science Foundation (SES-0922142).

Appendix A. Algorithm 1 for implementation of the cube method

Define the balancing matrix $A = (\mathbf{a}_1, \dots, \mathbf{a}_k, \dots, \mathbf{a}_N)$, where $\mathbf{a}_k = \mathbf{x}_k/\pi_k$. First initialize with $\pi_0(S) = \pi$ and $A(0) = A$. Next, at time $t = 0, \dots, T$, repeat the three following steps.

Step 1:

Let $E(t) = F(t) \cap \text{Ker}A(t)$, where

$$F(t) = \{\mathbf{v} \in \mathbb{R}^N : v_k = 0 \text{ if } \pi_k(t) \text{ is an integer.}\}$$

- If $E(t) \neq \{0\}$, then generate any vector $\mathbf{u}_t(S) \neq \mathbf{0}$ in $E(t)$, random or not. Put $A(t+1) = A(t)$.

- If $E(t) = \{0\}$, let m_t denote the largest integer such that $F(t) \cap \text{Ker}A_{m_t}(t) \neq \{0\}$, where $A_{m_t}(t)$ denotes the matrix given by the first m_t rows of $A(t)$. Generate any vector $\mathbf{u}_t(S) \neq 0$ in $F(t) \cap \text{Ker}A_{m_t}(t)$, random or not. Put $A(t+1) = A_{m_t}(t)$.

Step 2:

Compute the scalars $\lambda_{1t}^*(S)$ and $\lambda_{2t}^*(S)$, which are the largest values of λ_{1t} and λ_{2t} such that

$$0 \leq \pi_t(S) + \lambda_{1t}\mathbf{u}_t(S) \leq 1, \quad 0 \leq \pi_t(S) - \lambda_{2t}\mathbf{u}_t(S) \leq 1,$$

where the inequalities are interpreted element-wise. Note that $\lambda_{1t}^*(S) > 0$ and $\lambda_{2t}^*(S) > 0$.

Step 3:

Select $\pi_{t+1}(S) = \pi_t(S) + \delta_t(S)$, where

$$\delta_t(S) = \begin{cases} \lambda_{1t}^*(S)\mathbf{u}_t(S) & \text{with probability } q(t) \\ -\lambda_{2t}^*(S)\mathbf{u}_t(S) & \text{with probability } 1 - q(t) \end{cases}$$

and $q(t) = \lambda_{2t}^*(S)/(\lambda_{1t}^*(S) + \lambda_{2t}^*(S))$.

The procedure ends at step T , when $\pi_T(S)$ has only integer (0–1) components.

References

- Breidt, F., Chauvet, G., 2010. Penalized balanced sampling. Tech. rep., submitted.
- Chauvet, G., Tillé, Y., 2006. A fast algorithm for balanced sampling. *Comput. Statist.* 21 (1), 53–62.
- Deville, J.-C., 2000. Note sur l’algorithme de Chen, Dempster et Liu. Tech. rep., CREST-ENSAI.
- Deville, J.-C., Tillé, Y., 2004. Efficient balanced sampling: the cube method. *Biometrika* 91 (4), 893–912.
- Deville, J.-C., Tillé, Y., 2005. Variance approximation under balanced sampling. *J. Statist. Plann. Inference* 128 (2), 569–591.
- Fattorini, L., 2006. Applying the Horvitz-Thompson criterion in complex designs: a computer-intensive perspective for estimating inclusion probabilities. *Biometrika* 93 (2), 269–278.
- Harville, D. A., 1997. *Matrix algebra from a statistician’s perspective*. Springer-Verlag, New York.
- Horvitz, D. G., Thompson, D. J., 1952. A generalization of sampling without replacement from a finite universe. *J. Amer. Statist. Assoc.* 47, 663–685.

- Matei, A., Tillé, Y., 2005. Evaluation of variance approximations and estimators in maximum entropy sampling with unequal probability and fixed sample size. *Journal of Official Statistics* 21 (4), 543–570.
- Thompson, M., Wu, C., 2008. Simulation-based randomized systematic pps sampling under substitution of units. *Survey Methodology* 34 (1), 3–11.
- Tillé, Y., 2006. *Sampling algorithms*. Springer Series in Statistics. Springer, New York.
- Yates, F., Grundy, P., 1953. Selection without replacement from within strata with probability proportional to size. *JRSS B* 15, 253–261.