## Survey Methodology

# Linearization versus bootstrap for variance estimation of the change between Gini indexes

by Guillaume Chauvet and Camelia Goga

SURVEY
METHODOLOGY

A JOURNAL
PUBLISHED BY
STATISTICS CANADA

Statistics   Statistique
Canada     Canada

Canada

## How to obtain more information

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

You can also contact us by

**email at** STATCAN.infostats-infostats.STATCAN@canada.ca

**telephone,** from Monday to Friday, 8:30 a.m. to 4:30 p.m., at the following toll-free numbers:

| | |
|---|---|
| • Statistical Information Service | 1-800-263-1136 |
| • National telecommunications device for the hearing impaired | 1-800-363-7629 |
| • Fax line | 1-877-287-4369 |

**Depository Services Program**

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

## Standards of service to the public

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

## Note of appreciation

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

## Standard table symbols

The following symbols are used in Statistics Canada publications:

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

# Linearization versus bootstrap for variance estimation of the change between Gini indexes

**Guillaume Chauvet and Camelia Goga[1]**

## Abstract

This paper investigates the linearization and bootstrap variance estimation for the Gini coefficient and the change between Gini indexes at two periods of time. For the one-sample case, we use the influence function linearization approach suggested by Deville (1999), the without-replacement bootstrap suggested by Gross (1980) for simple random sampling without replacement and the with-replacement of primary sampling units described in Rao and Wu (1988) for multistage sampling. To obtain a two-sample variance estimator, we use the linearization technique by means of partial influence functions (Goga, Deville and Ruiz-Gazen, 2009). We also develop an extension of the studied bootstrap procedures for two-dimensional sampling. The two approaches are compared on simulated data.

**Key Words:** Composite estimator; Horvitz-Thompson estimator; Influence function; Intersection estimator; Replication weights; Two-sample survey; Two-dimensional sampling design; Union estimator; Variance estimation.

## 1 Introduction

The Gini coefficient (Gini, 1914) is one of the best known concentration measure often desired in economical studies. If $\mathcal{Y}_1$ denotes a quantitative positive variable such as the income and $F_1(\cdot)$ denotes its distribution function defined on $]-\infty, \infty[$, the Gini coefficient is

$$G_1 = \frac{1}{2} \frac{\iint |v - u| \, dF(u) \, dF(v)}{\int u \, dF(u)},$$

provided $\int u \, dF(u) \neq 0$. The Gini coefficient measures the dispersion of a quantitative positive variable within a population. Statistical institutes generally make use of the Gini coefficient to evaluate the income inequalities of a country at different periods of time, or of different countries at the same time. In the last decades, the Gini coefficient has also been considered in economic and sociodemographic fields (see for example Navarro, Muntaner, Borrell, Benach, Quiroga, Rodriguez-Sanz, Vergès and Pasarin, 2006; Bhattacharya, 2007; Lai, Huang, Risser and Kapadia, 2008; Barrett and Donald, 2009), biology (Graczyk, 2007), environment (Druckman and Jackson, 2008; Groves-Kirkby, Denman and Phillips, 2009) or astrophysics (Lisker, 2008).

There is an extensive literature on variance estimation for the Gini coefficient with observations obtained from survey data, see Langel and Tillé (2013) for a review. Glasser (1962) and Sandström, Wretman and Waldèn (1985) considered the case of simple random sampling. Sandström, Wretman and Waldèn (1988) listed possible variance estimators for a general sampling design, including a jackknife variance estimator. This latter approach was further investigated by Yitzhaki (1991), Karagiannis and Kovačević (2000) and

---

1. Guillaume Chauvet, Université Rennes, ENSAI, CNRS, IRMAR - UMR 6625, F-35000 Rennes, France. E-mail: chauvet@ensai.fr; Camelia Goga, Laboratoire des Mathématiques de Besançon, Université de Bourgogne Franche-Comté, UMR 6623, Besançon Cedex, France. E-mail: camelia.goga@univ-fcomte.fr.

Berger (2008). Linearization variance estimation was studied by Kovačević and Binder (1997), and Berger (2008) demonstrated the equivalence between linearization and a generalized jackknife technique first suggested by Campbell (1980). Qin, Rao and Wu (2010) proposed bootstrap and empirical likelihood based confidence intervals for the Gini coefficient. They studied these methods both theoretically and empirically in the particular case of stratified with replacement simple random sampling. However, bootstrap variance estimation has not been compared with alternative methods for the change between Gini indexes.

In this article, we consider linearization versus bootstrap to estimate the change between Gini indexes. The paper is structured as follows. In Section 2, we first consider the estimation of the Gini coefficient in the one-sample case. The notation is defined in Section 2.1, and the substitution estimator of the Gini coefficient is presented in Section 2.2. The linearization variance estimator is given in Section 2.3, with application to the simple random sampling (SI) design and to a multistage sampling design. The main principles of the weighted bootstrap are briefly reviewed at the beginning of Section 2.4, and the without-replacement bootstrap (BWO) suitable for SI sampling is introduced in Section 2.4.1, while the bootstrap of primary sampling units (BWR) suitable for multistage sampling is introduced in Section 2.4.2. In Section 3, we consider the estimation of the change between Gini indexes in the two-sample case. The notation is defined in Section 3.1, and we briefly review the principles of composite estimation which is applied in Section 3.1.1 for the two-dimensional SI design (SI2) and in Section 3.1.2 for a two-dimensional two-stage sampling design (MULT2). The composite estimator of the change between Gini indexes is presented in Section 3.2. The linearization variance estimator by means of the partial influence functions is given in Section 3.3, with application to the SI2 design and to the MULT2 design. An extension of the BWO for the SI2 design and of the BWR for the MULT2 design are then presented in Section 3.4. Linearization and the proposed bootstrap methods are compared in Section 4 through a simulation study. Section 5 concludes.

## 2  One sample case

### 2.1  Notation

Let $U$ denote some finite population of size $N$ whose units may be identified by the labels $k = 1, \ldots, N$. Suppose that the variable $\mathcal{Y}_1$ is measured on the population $U$, and let $y_{11}, \ldots, y_{1N}$ denote the values taken by $\mathcal{Y}_1$ on the units in the population. Let $M_1 = \sum_{k \in U} \delta_{y_{1k}}$ denote the discrete measure taking unit mass on any point $y_{1k}$ in the population and 0 elsewhere, with $\delta_{y_{1k}}$ the Dirac mass at $y_{1k}$. Most of the parameters of interest $\theta_1$ studied in surveys can be written as a functional $T$ of $M_1$, namely $\theta_1 = T(M_1)$. For instance, the total $t_{y1} = \sum_{k \in U} y_{1k}$ equals $\int \mathcal{Y}_1 dM_1$. In practice, a sample $s$ (with or without repetitions) is selected by means of a sampling design $p(\cdot)$, and we observe the values $y_{1k}$ for $k \in s$ only. A substitution principle is used for estimation (see Deville, 1999, and Goga, Deville and Ruiz-Gazen, 2009). Let $\pi_k$ denote the expected number of draws for unit $k$ in the sample; in case of without-replacement sampling, this is the probability that unit $k$ is selected in the sample. Let $\hat{M}_1 = \sum_{k \in s} w_k \delta_{y_{1k}}$ denote the discrete measure taking

mass $w_k$ on any point in the sample and 0 elsewhere, where $w_k = \pi_k^{-1}$ is the sampling weight. Substituting $\hat{M}_1$ into $\theta_1$ yields the estimator $\hat{\theta}_1 = T(\hat{M}_1)$.

For a without-replacement sampling design, the substitution estimator for a total is the so-called Horvitz-Thompson (HT) estimator $\hat{t}_{y1}^{\mathrm{HT}} = \sum_{k \in s} w_k y_{1k}$. The HT variance estimator is

$$v^{\mathrm{HT}}\left(\hat{t}_{y1}^{\mathrm{HT}}\right) = \sum_{k \in s}\sum_{l \in s} \frac{\Delta_{kl}}{\pi_{kl}} \frac{y_{1k}}{\pi_k} \frac{y_{1l}}{\pi_l}, \tag{2.1}$$

where $\pi_{kl} = \Pr(k, l \in s)$ denotes the probability that units $k$ and $l$ are selected jointly in the sample, and $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$. In the particular case of simple random sampling without replacement (SI) of size $n$, we have $\hat{t}_{y1}^{\mathrm{HT}} = N\bar{y}_{1,s}$ with $\bar{y}_{1,s} = n^{-1}\sum_{k \in s} y_{1k}$, and formula (2.1) yields

$$v^{\mathrm{HT}}\left(\hat{t}_{y1}^{\mathrm{HT}}\right) = N^2\left(\frac{1}{n} - \frac{1}{N}\right) S_{y_{1,s}}^2 \quad \text{where} \quad S_{y_{1,s}}^2 = \frac{1}{n-1}\sum_{k \in s}(y_{1k} - \bar{y}_{1,s})^2. \tag{2.2}$$

For a with-replacement sampling design, the substitution estimator for a total is the so-called Hansen-Hurwitz (HH) estimator $\hat{t}_{y1}^{\mathrm{HH}} = \sum_{k \in s} w_k y_{1k}$. We consider the important case of multistage sampling, where the $N$ units are grouped inside $N_I$ non-overlapping Primary Sampling Units (PSU) $U_1, \ldots, U_{N_I}$, and where a with-replacement first-stage sample $s_I$ of size $m$ is selected. Let $\pi_{Ii}$ denote the expected number of draws for the PSU $U_i$ in $s_I$. A second-stage sample $s_i$ is then selected inside any $i \in s_I$ by means of some sampling design $p_i(\cdot)$. Let $\pi_{k|i}$ denote the expected number of draws for unit $k$ in $s_i$. The estimated measure is then $\hat{M}_1 = \sum_{i \in s_I}\sum_{k \in s_i} \pi_{Ii}^{-1}\pi_{k|i}^{-1}\delta_{y_{1k}}$. We have $\hat{t}_{y1}^{\mathrm{HH}} = \sum_{i \in s_I} \pi_{Ii}^{-1}\hat{Y}_i$ where $\hat{Y}_i = \sum_{k \in s_i} \pi_{k|i}^{-1}y_{1k}$, and an unbiased variance estimator for $\hat{t}_{y1}^{\mathrm{HH}}$ is

$$v^{\mathrm{HH}}\left(\hat{t}_{y1}^{\mathrm{HH}}\right) = \frac{m}{m-1}\sum_{i \in s_I}\left(\frac{\hat{Y}_i}{\pi_{Ii}} - \frac{\hat{t}_{y1}^{\mathrm{HH}}}{m}\right)^2. \tag{2.3}$$

## 2.2 Estimating the Gini coefficient

If the variable $\mathcal{Y}_1$ is measured on the population $U$, the Gini coefficient is

$$G_1 = \frac{1}{2}\frac{\sum_{k \in U}\sum_{l \in U}|y_{1k} - y_{1l}|}{N\sum_{k \in U} y_{1k}},$$

see for example Nygård and Sandström (1985). It follows that $G_1$ is zero if $\mathcal{Y}_1$ is constant on the population, which occurs when the total of $\mathcal{Y}_1$ is equally distributed among all the population individuals. In the opposite case, when only one individual owns the whole amount of $\mathcal{Y}_1$, $G_1$ is maximized and equal to $1 - 1/N$: the total of $\mathcal{Y}_1$ is then concentrated in one point only, which means maximum inequality among members of the population.

If all individuals $k \neq l$ have different values for the variable $\mathcal{Y}_1$, the Gini coefficient $G_1$ is

$$G_1 = \frac{\sum_{k=1}^{N} y_{1(k)} \left(2k/N - 1\right)}{t_{y1}} - \frac{1}{N} = \frac{\sum_{k \in U} y_{1k} \left\{2F_{1N}\left(y_{1k}\right) - 1\right\}}{t_{y1}} - \frac{1}{N} \tag{2.4}$$

with $y_{1(1)} \leq \cdots \leq y_{1(N)}$ the ordered values and $F_{1N}(\cdot) = N^{-1} \sum_{k \in U} 1_{\{y_{1k} \leq \cdot\}}$ the finite population distribution function; see Sandström, Wretman and Waldèn (1988) and Deville (1997) for further details on the derivation of (2.4). Nygård and Sandström (1985) called the term $-1/N$ the Gini finite population correction and gave several reasons to make this correction, such as the non-negativity of the lower bound of $G_1$. As is frequently done in the literature (see for example Glasser, 1962), this correction is ignored in the sequel. We redefine the Gini coefficient as

$$G_1 = \frac{\sum_{k \in U} y_{1k} \left\{2F_{1N}\left(y_{1k}\right) - 1\right\}}{t_{y1}} = \frac{\int \left\{2F_{1N}(y) - 1\right\} y \, dM_1(y)}{\int y \, dM_1(y)} \tag{2.5}$$

where the finite population distribution function $F_{1N}(\cdot)$ is a functional family

$$F_{1N}(y) = \frac{1}{\int dM_1(y)} \int 1_{\{\xi \leq y\}} \, dM_1(\xi) \tag{2.6}$$

indexed by $y$. Substituting $\hat{M}_1$ into (2.5) and (2.6) yields the estimator

$$\hat{G}_1 = \frac{\int \left\{2\hat{F}_{1N}(y) - 1\right\} y \, d\hat{M}_1(y)}{\int y \, d\hat{M}_1(y)} = \frac{\sum_{k \in s} w_k \left\{2\hat{F}_{1N}(y_{1k}) - 1\right\} y_{1k}}{\sum_{k \in s} w_k y_{1k}}, \tag{2.7}$$

where

$$\hat{F}_{1N}(y) = \frac{1}{\int d\hat{M}_1(y)} \int 1_{\{\xi \leq y\}} \, d\hat{M}_1(\xi) = \frac{1}{\sum_{k \in s} w_k} \sum_{k \in s} w_k 1_{\{y_{1k} \leq y\}} \tag{2.8}$$

is the substitution estimator of the distribution function $F_{1N}$.

## 2.3 Linearization variance estimation

We give below some brief details about the influence function linearization (IFL) (Deville, 1999), which consists in giving a first-order expansion of the substitution estimator $\hat{\theta}_1 = T(\hat{M}_1)$ around the true value $\theta_1 = T(M_1)$, to approximate the error by a linear estimator of some artificial *linearized variable*. More precisely, the first derivatives of $T$ with respect to $M_1$ are the influence functions

$$\mathrm{IT}(M_1; y) = \lim_{h \to 0} \frac{T(M_1 + h\delta_y) - T(M_1)}{h},$$

and $u_{1k} = \mathrm{IT}(M_1; y_{1k})$ is the linearized variable for all $k \in U$. Suppose that $T(\cdot)$ is homogeneous, namely there exists some positive number $\beta$ dependent on $T$ such that $T(rM_1) = r^\beta T(M_1)$ for any real $r > 0$.

Assume also that $\lim_{N\to\infty} N^{-\beta} T(M_1) < \infty$. Under some additional regularity assumptions upon $T(\cdot)$ and the sampling design (e.g., Goga and Ruiz-Gazen, 2014), Deville (1999) establishes that

$$\hat{\theta}_1 - \theta_1 = \left( \sum_{k\in s} w_k u_{1k} - \sum_{k\in U} u_{1k} \right) + o_p \left( N^{\beta} n^{-1/2} \right),$$

so that the error $\hat{\theta}_1 - \theta_1$ can be approximated by the error of the HT estimator for the total of the linearized variable $u_{1k}$. For a without-replacement sampling design, using a sample-based estimator $\hat{u}_{1k}$ of the linearized variable $u_{1k}$ in the HT variance estimator yields the variance estimator

$$v_{\mathrm{LIN}}^{\mathrm{HT}} \left( \hat{\theta}_1 \right) = \sum_{k\in s}\sum_{l\in s} \frac{\Delta_{kl}}{\pi_{kl}} \frac{\hat{u}_{1k}}{\pi_k} \frac{\hat{u}_{1l}}{\pi_l}, \tag{2.9}$$

where $\pi_{kl} = \Pr(k, l \in s)$ denotes the probability that units $k$ and $l$ are selected jointly in the sample, and $\Delta_{kl} = \pi_{kl} - \pi_k \pi_l$. Several results of asymptotic normality have been proved for specific sampling designs, see Hájek (1960, 1961, 1964), Rosén (1972), Sen (1980), Krewski and Rao (1981), Gordon (1983), Ohlsson (1986, 1989), Chen and Rao (2007), Bränden and Jonasson (2012), Saegusa and Wellner (2013) and Chauvet (2015), among others. If the sampling design is such that the substitution estimator $\hat{\theta}_1$ satisfies a central-limit theorem, an approximately $(1 - 2\alpha)\%$ confidence interval is $\left[ \hat{\theta}_1 - z_\alpha \sqrt{v_{\mathrm{lin}}\left(\hat{\theta}_1\right)}, \hat{\theta}_1 + z_\alpha \sqrt{v_{\mathrm{lin}}\left(\hat{\theta}_1\right)} \right]$ where $z_\alpha$ is the upper $\alpha\%$ cutoff for the standard normal distribution.

In case of the Gini coefficient, we have $\beta = 0$ and the linearized variable is

$$u_{1k} = 2F_{1N}(y_{1k}) \frac{y_{1k} - \bar{y}_{1k,U<}}{t_{y1}} - y_{1k} \frac{G_1 + 1}{t_{y1}} + \frac{1 - G_1}{N}, \tag{2.10}$$

where $\bar{y}_{1k,U<} = \left( \sum_{l\in U} 1_{\{y_{1l}<y_{1k}\}} \right)^{-1} \sum_{j\in U} y_{1j} 1_{\{y_{1j}<y_{1k}\}}$ denotes the mean of the $y_{1j}$ lower than $y_{1k}$, see Deville (1999). Kovačević and Binder (1997) derived the same expression by means of the estimating equations linearization method; using the Demnati and Rao (2004) linearization approach also leads to the same result. The estimated linearized variable is

$$\hat{u}_{1k} = 2\hat{F}_{1N}(y_{1k}) \frac{y_{1k} - \bar{y}_{1k,s<}}{\hat{t}_{y1}} - y_{1k} \frac{\hat{G}_1 + 1}{\hat{t}_{y1}} + \frac{1 - \hat{G}_1}{\hat{N}} \tag{2.11}$$

where $\bar{y}_{1k,s<} = \left( \sum_{l\in s} w_l 1_{\{y_{1l}<y_{1k}\}} \right)^{-1} \sum_{j\in s} w_j y_{1j} 1_{\{y_{1j}<y_{1k}\}}$.

In the particular SI case, the linearization variance estimator for the Gini coefficient is

$$v_{\mathrm{LIN}}^{\mathrm{HT}} \left( \hat{G}_1 \right) = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) S_{\hat{u}_1,s}^2 \quad \text{where} \quad S_{\hat{u}_1,s}^2 = \frac{1}{n-1} \sum_{k\in s} \left( \hat{u}_{1k} - \bar{\bar{u}}_{1,s} \right)^2, \tag{2.12}$$

and where $\bar{\bar{u}}_{1,s} = n^{-1} \sum_{k\in s} \hat{u}_{1k}$. In the particular case of multistage sampling and with-replacement sampling of PSUs, the linearization variance estimator for the Gini coefficient is

$$v_{\text{LIN}}^{\text{HH}}\left(\hat{G}_1\right) \;=\; \frac{m}{m-1}\sum_{i\in s_I}\left(\frac{\hat{U}_{1i}}{\pi_{Ii}} - \frac{\hat{t}_{\hat{u}_1}^{\text{HH}}}{m}\right)^2 \quad \text{where} \quad \hat{U}_{1i} \;=\; \sum_{k\in s_i}\pi_{k|i}^{-1}\hat{u}_{1k}. \tag{2.13}$$

## 2.4 Bootstrap variance estimation

The use of bootstrap techniques in survey sampling has been extensively studied in the literature. The main bootstrap techniques may be thought as particular cases of the weighted bootstrap (Bertail and Combris, 1997; Antal and Tillé, 2011; Beaumont and Patak, 2012); see also Shao and Tu (1995, Chapter 6), Davison and Hinkley (1997, Section 3.7) and Davison and Sardy (2007) for detailed reviews. Under a weighted bootstrap procedure, the measure $\hat{M}_1 = \sum_s w_k \delta_{y_k}$ is estimated, conditionally on the sample $s$, by the bootstrap measure

$$\hat{M}_1^* \;=\; \sum_{k\in s} w_k D_k \delta_{y_k} \tag{2.14}$$

where $D = \{D_k\}_{k\in s}$ denotes a (random) vector of resampling weights. We note $E_*$ and $V_*$ for the expectation and variance with respect to the resampling scheme. In case of without-replacement sampling, the vector $D$ is generated in such a way that

$$E_*\left(\sum_s w_k D_k y_k\right) \;\simeq\; \hat{t}_{y1}^{\text{HT}} \quad \text{and} \quad V_*\left(\sum_s w_k D_k y_k\right) \;\simeq\; v^{\text{HT}}\left(\hat{t}_{y1}^{\text{HT}}\right) \tag{2.15}$$

so that the two first moments of the HT-estimator are approximately matched. In case of with-replacement sampling, the vector $D$ is generated in such a way that

$$E_*\left(\sum_s w_k D_k y_k\right) \;\simeq\; \hat{t}_{y1}^{\text{HH}} \quad \text{and} \quad V_*\left(\sum_s w_k D_k y_k\right) \;\simeq\; v^{\text{HH}}\left(\hat{t}_{y1}^{\text{HH}}\right) \tag{2.16}$$

so that the two first moments of the HH-estimator are approximately matched.

Under any weighted bootstrap technique, the plug-in estimator of $\theta_1 = T(M_1)$ is $\hat{\theta}_1^* = T\left(\hat{M}_1^*\right)$, and the variance of $\hat{\theta}_1 = T\left(\hat{M}_1\right)$ is estimated by

$$V_*\left(\hat{\theta}_1^*\right) \;=\; E_*\left\{\hat{\theta}_1^* - E_*\left(\hat{\theta}_1^*\right)\right\}^2. \tag{2.17}$$

Since the variance estimator (2.17) may be difficult to compute exactly, a simulation-based variance estimator may be used instead. More precisely, $C$ independent realizations $D_1,\ldots,D_C$ of the vector $D$ are generated, and we denote $\hat{\theta}_{1c}^* = T\left(\hat{M}_{1c}^*\right)$ with $\hat{M}_{1c}^*$ the Bootstrap measure associated to the vector $D_c$. Then $V\left(\hat{\theta}_1\right)$ is estimated by

$$v_B\left(\hat{\theta}_1\right) \;=\; \frac{1}{C-1}\sum_{c=1}^{C}\left\{\hat{\theta}_{1c}^* - \frac{1}{C}\sum_{c'=1}^{C}\hat{\theta}_{1c'}^*\right\}^2. \tag{2.18}$$

Two types of confidence intervals are usually computed. The percentile method makes use of the ordered bootstrap estimates $\hat{\theta}_{(1c)}^*$, $c = 1,\ldots,C$ to form a $(1 - 2\alpha)\%$ confidence interval $[\hat{\theta}_{(1L)}^*, \hat{\theta}_{(1U)}^*]$ with $L = \alpha C$

and $U = (1 - \alpha) C$. The bootstrap$-t$ involves the estimation of the pivotal statistic $t = (\hat{\theta}_1 - \theta_1) / \sqrt{v_{\text{BWO}}(\hat{\theta}_1)}$ by its bootstrap counterpart $t^* = (\hat{\theta}_1^* - \hat{\theta}_1) / \sqrt{v_{\text{BWO}}^*(\hat{\theta}_1^*)}$, where $v_{\text{BWO}}^*(\hat{\theta}_1^*)$ is obtained by applying the bootstrap procedure to the resample $s^*$. The bootstrap$-t$ is computationally very intensive since a double bootstrap is required, and is thus less attractive for a data user. Therefore, we do not pursue this approach further and we focus on the percentile method.

Linearization methods provide variance formulas applicable to general sampling designs, but involve possibly intricate computation of derivatives for complex parameters of interest such as the Gini coefficient. Unlike the linearization, the bootstrap avoids theoretical work by re-calculating the existing estimation system repeatedly. Replicate weights are supplied with the data set, and may be easily used to produce variance estimates for a wide range of statistics. However, a bootstrap technique is usually not suitable for general sampling designs. That is, a particular sampling design usually requires a tailor made resampling scheme. In this paper, we focus on two particular bootstrap techniques, which will be generalized in Section 3 to the two-sample context.

### 2.4.1 Without-replacement bootstrap for SI sampling

When the sample $s$ is selected by means of SI, we consider the without replacement bootstrap (BWO) introduced by Gross (1980). The approach is readily extended to stratified simple random sampling (STSI) with a finite number of strata. Suppose that $N/n$ is an integer. Then the vector $D$ is obtained by, first creating a pseudo-population $U^*$ of size $N$ by duplicating $N/n$ times each unit $k$ in the original sample $s$, and then by selecting a SI resample $s^*$ of size $n$ in $U^*$.

The bootstrap measure is given by (2.14), where the resampling weight $D_k$ is the number of times unit $k \in s$ is selected in $s^*$. The building of $U^*$ may be avoided by noting that under the BWO procedure, the vector $D$ follows a multivariate hypergeometric distribution. Therefore, the resampling weights may be directly generated. It can be shown that the BWO procedure leads to

$$E_* \left( \sum_s w_k D_k y_k \right) = \hat{t}_{y1}^{\text{HT}} \quad \text{and} \quad V_* \left( \sum_s w_k D_k y_k \right) = \frac{1 - n^{-1}}{1 - N^{-1}} v^{\text{HT}} \left( \hat{t}_{y1}^{\text{HT}} \right), \qquad (2.19)$$

where $v^{\text{HT}} \left( \hat{t}_{y1}^{\text{HT}} \right)$ is given in (2.2), so that equation (2.15) is approximately matched for a large sample size.

Several solutions have been proposed to handle the case when $N/n$ is not an integer, see Chao and Lo (1985), Bickel and Freedman (1984), Sitter (1992b), Booth, Butler and Hall (1994), Presnell and Booth (1994), among others. The generalization of BWO variance estimation for unequal probability sampling designs is considered in Särndal, Swensson and Wretman (1992) and Chauvet (2007).

### 2.4.2 With-replacement bootstrap for multistage sampling

When the sample $s$ is selected by means of multistage sampling and with-replacement unequal probability sampling of PSUs, we consider the bootstrap of PSUs (BWR) introduced by Rao and Wu (1988).

A with-replacement resample $s_I^*$ of size $m - 1$ is selected by means of simple random sampling with replacement (SIR) in the original first-stage sample $s_I$. The bootstrap measure is

$$\hat{M}_1^* \;=\; \frac{m}{m-1}\sum_{i\in s_I^*}\sum_{k\in s_i}\pi_{Ii}^{-1}\pi_{k|i}^{-1}\delta_{y_{1k}} \;=\; \sum_{k\in s}w_k D_k \delta_{y_k}, \tag{2.20}$$

where the resampling weight $D_k$ equals $m(m-1)^{-1}$ multiplied by the number of times the PSU containing $k$ is selected in $s_I^*$.

The resampling size $m - 1$ is used to reproduce the usual unbiased variance estimator in the linear case (see Rao and Wu, 1988). It can be shown that the BWR procedure leads to

$$E_*\left(\sum_s w_k D_k y_k\right) \;=\; \hat{t}_{y1}^{\mathrm{HH}} \quad\text{and}\quad V_*\left(\sum_s w_k D_k y_k\right) \;=\; v^{\mathrm{HH}}\left(\hat{t}_{y1}^{\mathrm{HH}}\right), \tag{2.21}$$

where $v^{\mathrm{HH}}\left(\hat{t}_{y1}^{\mathrm{HH}}\right)$ is given in (2.3), so that equation (2.16) is exactly matched. The BWR procedure is particulary simple, since involving a resampling for the first-stage of sampling only, the sub-samples of Secondary sampling Units (SSUs) being left unchanged inside the resampled PSUs.

# 3  Two-sample case

## 3.1  Notation and composite estimation

Suppose now that two variables $\mathcal{Y}_1$ and $\mathcal{Y}_2$ are measured on the population $U$, and let $y_{d1},\dots,y_{dN}$ denote the values taken by $\mathcal{Y}_d, d = 1, 2,$ on the units in the population. The variables $\mathcal{Y}_1$ and $\mathcal{Y}_2$ may typically refer to some characteristic of interest collected at two different times $\tau_1$ and $\tau_2$. We consider the estimation of parameters $\Delta\theta$ that can be written as a functional $\Delta\theta = T(M_1, M_2)$, where $M_d = \sum_{k\in U}\delta_{\{y_{dk}\}}$. For instance, the linear case $\Delta t = t_{y2} - t_{y1}$ corresponds to the difference between the totals $t_{y2} = \sum_{k\in U}y_{2k}$ and $t_{y1} = \sum_{k\in U}y_{1k}$.

Let $s_1$ and $s_2$ be two samples of sizes $n_1$ and $n_2$, respectively, selected from the same population $U$ according to some two-dimensional sampling design $p(\cdot,\cdot)$ (see Goga, 2003). The variable $\mathcal{Y}_1$ is measured on $s_1$, while the variable $\mathcal{Y}_2$ is measured on $s_2$. Plugging sample-based estimators $\hat{M}_d$ in $\Delta\theta$ yields the substitution estimator $\widehat{\Delta\theta} = T(\hat{M}_1, \hat{M}_2)$. Unlike the one-sample case, several estimators $\hat{M}_d$ are possible. In what follows, we focus on the general class of *composite estimators* introduced by Goga, Deville and Ruiz-Gazen (2009). We note $s_{1\bullet} = s_1 \setminus s_2$, $s_3 = s_1 \cap s_2$ and $s_{2\bullet} = s_2 \setminus s_1$. For $\Diamond \in \{1\bullet, 3, 2\bullet\}$, we note $\pi_{\Diamond,k}$ the expected number of draws for unit $k$ in $s_\Diamond$ and $\hat{M}_{d,\Diamond} = \sum_{k\in s_\Diamond}w_{\Diamond,k}\delta_{y_{dk}}$, where $w_{\Diamond,k} = \pi_{\Diamond,k}^{-1}$. The composite estimators of $M_1$ and $M_2$ are

$$\hat{M}_1^{\mathrm{co}}(a) \;=\; a\,\hat{M}_{1,1\bullet} + (1-a)\,\hat{M}_{1,3} \quad\text{and}\quad \hat{M}_2^{\mathrm{co}}(b) \;=\; b\,\hat{M}_{2,2\bullet} + (1-b)\,\hat{M}_{2,3}, \tag{3.1}$$

where $a$ and $b$ are some known constants. The choice $a = b = 0$ leads to the *intersection estimator* with $\hat{M}_1^{\mathrm{int}} = \hat{M}_{1,3}$ and $\hat{M}_2^{\mathrm{int}} = \hat{M}_{2,3}$, where the overlapping sample $s_3$ only is used.

When estimating the parameter $\Delta t = t_{y2} - t_{y1}$, the composite estimator is

$$\widehat{\Delta t}^{\text{co}}(a,b) = \hat{t}_{y_2}^{\text{co}} - \hat{t}_{y_1}^{\text{co}}, \tag{3.2}$$

where $\hat{t}_{y_1}^{\text{co}} = \int y d\hat{M}_1^{\text{co}}(y)$ and $\hat{t}_{y_2}^{\text{co}} = \int y d\hat{M}_2^{\text{co}}(y)$. It may be rewritten as

$$\widehat{\Delta t}^{\text{co}}(a,b) = b(\hat{t}_{y_2,s_{2\bullet}} - \hat{t}_{y_2,s_3}) - a(\hat{t}_{y_1,s_{1\bullet}} - \hat{t}_{y_1,s_3}) + (\hat{t}_{y_2,s_3} - \hat{t}_{y_1,s_3}), \tag{3.3}$$

where $\hat{t}_{y_d,s_\diamond} = \sum_{k\in s_\diamond} w_{\diamond,k} y_{dk}$. The variance of the composite estimator is

$$V\left\{\widehat{\Delta t}^{\text{co}}(a,b)\right\} = (b,-a,1) V\left\{(\hat{t}_{y_2,s_{2\bullet}} - \hat{t}_{y_2,s_3}, \hat{t}_{y_1,s_{1\bullet}} - \hat{t}_{y_1,s_3}, \hat{t}_{y_2,s_3} - \hat{t}_{y_1,s_3})^\top\right\}(b,-a,1)^\top. \tag{3.4}$$

Finding the vector $(a_{\text{opt}}, b_{\text{opt}})^\top$ which minimizes the variance in (3.4) leads to the *optimal composite estimator* (Goga, Deville and Ruiz-Gazen, 2009, Section 3.6). Note that this is not an estimator per se, since it depends on unknown quantities which need to be estimated in practice. However, this is a useful benchmark which we will use for the appraisal of simpler composite estimators.

A variance estimator is obtained by substituting in (3.4) an estimator of the variance-covariance matrix. The derivation of variance estimators is detailed in Sections 3.1.1 and 3.1.2 for two examples of two-dimensional sampling designs.

## 3.1.1 Two-dimensional SI design

The two-dimensional SI design (SI2) of fixed size $(n_{1\bullet}, n_3, n_{2\bullet})$ assigns equal probabilities to all $s = (s_1, s_2)$ for which the associated subsamples $s_{1\bullet}$, $s_3$ and $s_{2\bullet}$ have the required sizes $n_{1\bullet}$, $n_3$ and $n_{2\bullet}$, see Goga (2003) and Qualité and Tillé (2008). The SI2 design has the attractive property that the marginal samples $s_{1\bullet}$, $s_3$ and $s_{2\bullet}$ are SI samples from the population $U$. Similarly, $s_1$ is a SI sample of size $n_1 = n_{1\bullet} + n_3$, and $s_2$ is a SI sample of size $n_2 = n_{2\bullet} + n_3$. For the SI2 sampling design, the composite estimator in (3.3) yields

$$\widehat{\Delta t}^{\text{co}}(a,b) = Nb(\bar{y}_{2,s_{2\bullet}} - \bar{y}_{2,s_3}) - Na(\bar{y}_{1,s_{1\bullet}} - \bar{y}_{1,s_3}) + N(\bar{y}_{2,s_3} - \bar{y}_{1,s_3}), \tag{3.5}$$

and the variance of the composite estimator is

$$V\left\{\widehat{\Delta t}^{\text{co}}(a,b)\right\} = N^2\left\{c_1(a) S_{y_1,U}^2 - 2c_{12}(a,b) S_{y_1 y_2,U} + c_2(b) S_{y_2,U}^2\right\}, \tag{3.6}$$

with

$$c_1(a) = \frac{(1-a)^2}{n_3} + \frac{a^2}{n_1 - n_3} - \frac{1}{N},$$

$$c_2(b) = \frac{(1-b)^2}{n_3} + \frac{b^2}{n_2 - n_3} - \frac{1}{N},$$

$$c_{12}(a,b) = \frac{(1-a)(1-b)}{n_3} - \frac{1}{N},$$

see Appendix for a proof.

We consider two examples. The choice $a = b = 0$ leads to the intersection estimator

$$\widehat{\Delta t}^{\text{int}} = \widehat{\Delta t}^{\text{co}}(0,0) = \frac{N}{n_3} \sum_{k \in s_3} (y_{2k} - y_{1k}), \tag{3.7}$$

and the variance simplifies as

$$V\left\{\widehat{\Delta t}^{\text{int}}\right\} = N^2 \left(\frac{1}{n_3} - \frac{1}{N}\right) S^2_{y_2 - y_1, U}. \tag{3.8}$$

The choice $a = n_1^{-1} n_{1\bullet}$ and $b = n_2^{-1} n_{2\bullet}$ leads to the *union estimator*

$$\widehat{\Delta t}^{\text{uni}} = \widehat{\Delta t}^{\text{co}}(n_1^{-1} n_{1\bullet}, n_2^{-1} n_{2\bullet}) = \frac{N}{n_2} \sum_{k \in s_2} y_{2k} - \frac{N}{n_1} \sum_{k \in s_1} y_{1k} \tag{3.9}$$

where the complete samples are used, and the variance may be written as

$$V\left\{\widehat{\Delta t}^{\text{uni}}\right\} = N^2 \left\{\left(\frac{1}{n_1} - \frac{1}{N}\right) S^2_{y_1, U} - 2\left(\frac{n_3}{n_1 n_2} - \frac{1}{N}\right) S_{y_1 y_2, U} + \left(\frac{1}{n_2} - \frac{1}{N}\right) S^2_{y_2, U}\right\}. \tag{3.10}$$

The variances of the union estimator and of the intersection estimator were derived by Qualité and Tillé (2008), see also Tam (1984).

The choice of $a$ and $b$ is of practical importance to obtain an efficient composite estimator. After some algebra, the vector $(a_{\text{opt}}, b_{\text{opt}})^\top$ which minimizes the variance of $\widehat{\Delta t}^{\text{co}}(a, b)$ is given by

$$(a_{\text{opt}}, b_{\text{opt}})^\top = A^{-1} X \tag{3.11}$$

with

$$A = \begin{pmatrix} \dfrac{n_1}{n_1 - n_3} & -\dfrac{S_{y_1 y_2, U}}{S^2_{y_1, U}} \\[2ex] -\dfrac{S_{y_1 y_2, U}}{S^2_{y_2, U}} & \dfrac{n_2}{n_2 - n_3} \end{pmatrix} \quad \text{and} \quad X = \left(1 - \frac{S_{y_1 y_2, U}}{S^2_{y_1, U}}, 1 - \frac{S_{y_1 y_2, U}}{S^2_{y_2, U}}\right)^\top. \tag{3.12}$$

For two variables $\mathcal{Y}_1$ and $\mathcal{Y}_2$ related to a same characteristic collected at two different times, $S_{y_1 y_2, U}$ is expected to be close to $S^2_{y_1, U}$ and $S^2_{y_2, U}$. The vector $X$ in (3.12) is in turn close to the null vector, and if the size of the overlapping sample $s_3$ is comparable to that of $s_{1\bullet}$ and $s_{2\bullet}$ we obtain $a_{\text{opt}} \simeq 0$ and $b_{\text{opt}} \simeq 0$. Therefore, using the intersection estimator where $a = b = 0$ seems reasonable in practice. On the contrary, the union estimator can be very inefficient; see Section 4.2 for an illustration. These conclusions are consistent with that of Qualité and Tillé (2008), Section 2.2.2.

Several variance estimators may be used for the composite estimator. Estimating the dispersions on the overlapping sample only yields the unbiased variance estimator

$$v_{\text{int}}^{\text{HT}}\left\{\widehat{\Delta t}^{\text{co}}(a, b)\right\} = N^2 \left\{c_1(a) S^2_{y_1, s_3} - 2c_{12}(a, b) S_{y_1 y_2, s_3} + c_2(b) S^2_{y_2, s_3}\right\}, \tag{3.13}$$

while an estimation on the whole samples yields

$$v_{\text{uni}}^{\text{HT}}\left\{\widehat{\Delta t}^{\text{co}}(a,b)\right\} = N^2\left\{c_1(a)S_{y_1,s_1}^2 - 2c_{12}(a,b)S_{y_1 y_2, s_3} + c_2(b)S_{y_2,s_2}^2\right\}. \tag{3.14}$$

Berger (2004) considered variance estimation for the union estimator under a maximum entropy rotating sampling scheme, by estimating separately the three components in (3.6).

### 3.1.2 Two-dimensional multistage design

We now consider a two-dimensional two-stage sampling design (MULT2). We assume that a with-replacement first-stage sample $s_I$ of size $m$ is first selected among the PSUs $U_1,\ldots,U_{N_I}$. Inside each PSU $i \in s_I$, a SI2 sample of size $(n_{1\bullet}^i, n_3^i, n_{2\bullet}^i)$ is then selected. This type of sampling design emerges in particular in case of a self-weighted two-stage design in two waves, with a partial replacement at the second wave of the SSUs selected at the first wave. The composite estimator in (3.3) yields

$$\widehat{\Delta t}^{\text{co}}(a,b) = \sum_{i \in s_I} \pi_{Ii}^{-1} \widehat{\Delta t}^{i,\text{co}}(a,b) \tag{3.15}$$

where

$$\widehat{\Delta t}^{i,\text{co}}(a,b) = N_i b\left(\bar{y}_{2,s_{2\bullet}^i} - \bar{y}_{2,s_3^i}\right) - N_i a\left(\bar{y}_{1,s_{1\bullet}^i} - \bar{y}_{1,s_3^i}\right) + N_i\left(\bar{y}_{2,s_3^i} - \bar{y}_{1,s_3^i}\right), \tag{3.16}$$

where $\bar{y}_{d,s_\diamond^i} = (n_\diamond^i)^{-1}\sum_{k \in s_\diamond^i} y_{\diamond k}$, where $s_\diamond^i = s_\diamond \cap U_i$, and where $N_i$ denotes the number of SSUs inside the PSU $u_i$.

For example, using the overlapping samples only inside the PSUs yields the intersection estimator

$$\widehat{\Delta t}^{\text{int}} = \sum_{i \in s_I} \pi_{Ii}^{-1} \widehat{\Delta t}^{i,\text{int}} \qquad \text{with} \qquad \widehat{\Delta t}^{i,\text{int}} = N_i\left(\bar{y}_{2,s_3^i} - \bar{y}_{1,s_3^i}\right). \tag{3.17}$$

Using the complete samples inside the PSUs yields the union estimator

$$\widehat{\Delta t}^{\text{uni}} = \sum_{i \in s_I} \pi_{Ii}^{-1} \widehat{\Delta t}^{i,\text{uni}} \qquad \text{with} \qquad \widehat{\Delta t}^{i,\text{uni}} = N_i\left(\bar{y}_{2,s_2^i} - \bar{y}_{1,s_1^i}\right). \tag{3.18}$$

We note that for any vector of values $(a,b)^\top$, the variance due to the first-stage of sampling for $\widehat{\Delta t}^{\text{co}}(a,b)$ is the same. The possible composite estimators thus differ with respect to the second-stage variance only. In view of the discussion in Section 3.1.1, we therefore expect the intersection estimator to be close to the optimal composite estimator; see Section 4.2 for an illustration. An unbiased variance estimator for $\widehat{\Delta t}^{\text{co}}(a,b)$ is given by

$$v^{\text{HH}}\left\{\widehat{\Delta t}^{\text{co}}(a,b)\right\} = \frac{m}{m-1}\sum_{i \in s_I}\left(\frac{\widehat{\Delta t}^{i,\text{co}}(a,b)}{\pi_{Ii}} - \frac{\widehat{\Delta t}^{\text{co}}(a,b)}{m}\right)^2. \tag{3.19}$$

## 3.2 Estimation of the change between Gini indexes

The change between Gini indexes $\Delta G = G_2 - G_1$ may be written as

$$\Delta G \;=\; \frac{\int \{2F_{2N}(y)-1\}\, y dM_2(y)}{\int y dM_2(y)} - \frac{\int \{2F_{1N}(y)-1\}\, y dM_1(y)}{\int y dM_1(y)} \tag{3.20}$$

where $F_{dN}(y) = N^{-1}\sum_{k\in U} 1_{\{y_{dk}\le y\}},\ d=1,2.$ Using composite estimation leads to

$$\widehat{\Delta G}^{\,\mathrm{co}}(a,b) \;=\; \frac{\int \{2\hat{F}_{2N}^{\,\mathrm{co}}(y)-1\}\, y d\hat{M}_2^{\,\mathrm{co}}(y)}{\int y d\hat{M}_2^{\,\mathrm{co}}(y)} - \frac{\int \{2\hat{F}_{1N}^{\,\mathrm{co}}(y)-1\}\, y d\hat{M}_1^{\,\mathrm{co}}(y)}{\int y d\hat{M}_2^{\,\mathrm{co}}(y)} \tag{3.21}$$

where $\hat{F}_{dN}^{\,\mathrm{co}}(y) = \left\{\int d\hat{M}_d^{\,\mathrm{co}}(y)\right\}^{-1} \int 1_{\{\xi\le y\}}\, d\hat{M}_d^{\,\mathrm{co}}(\xi).$

Usually, in a temporal sampling framework, the samples $s_1$ and $s_2$ are not independent. Consequently, our set-up differs from the usual estimation of functionals depending on distribution functions estimated with independent samples; see for example Pires and Branco (2002) and Reid (1981), who give the first-order expansion of a two-sample functional using the partial influence functions. Davison and Hinkley (1997, page 71) give bootstrap methods under a similar framework. Using a general two-dimensional sampling design $p(\cdot,\cdot)$, Goga, Deville and Ruiz-Gazen (2009) give a two-sample linearization technique of bivariate functionals that will be used in what follows.

## 3.3  Linearization variance estimation

To obtain the asymptotic variance of $\widehat{\Delta\theta}^{\,\mathrm{co}}(a,b)$, we adopt the asymptotic framework introduced by Goga, Deville and Ruiz-Gazen (2009), which is an extension to the two-sample case of the asymptotic framework of Isaki and Fuller (1982). Define, when they exist, the *partial influence functions* of a functional $T(M_1,M_2)$ at point $y$ as

$$I_1 T(M_1,M_2;y) \;=\; \lim_{h\to 0} \frac{T(M_1+h\delta_y,M_2)-T(M_1,M_2)}{h},$$

$$I_2 T(M_1,M_2;y) \;=\; \lim_{h\to 0} \frac{T(M_1,M_2+h\delta_y)-T(M_1,M_2)}{h}.$$

We define the *linearized variables* $u_{dk} = I_d T(M_1,M_2;y_{dk})$ for $d=1,2$ as the partial influence functions of $T$ at $(M_1,M_2)$ and $y=y_{dk}$. For the change between Gini indexes $\Delta G$, the linearized variables $u_{dk}$ may be computed using (2.10), namely

$$u_{dk} \;=\; 2F_{dN}(y_{dk})\frac{y_{dk}-\bar{y}_{dk,U<}}{t_{y_d}} - y_{dk}\frac{G_d+1}{t_{y_d}} + \frac{1-G_d}{N}, \tag{3.22}$$

where $\bar{y}_{dk,U<} = \left(\sum_{l\in U} 1_{\{y_{dl}<y_{dk}\}}\right)^{-1} \sum_{j\in U} y_{dj} 1_{\{y_{dj}<y_{dk}\}}.$ The estimated linearized variable is

$$\hat{u}_{dk} \;=\; 2\hat{F}_{dN}^{\,\mathrm{co}}(y_{dk})\frac{y_{dk}-\bar{y}_{dk,s<}^{\,\mathrm{co}}}{\hat{t}_{y1}^{\,\mathrm{co}}} - y_{dk}\frac{\hat{G}_d^{\,\mathrm{co}}+1}{\hat{t}_{y1}^{\,\mathrm{co}}} + \frac{1-\hat{G}_d^{\,\mathrm{co}}}{\hat{N}}. \tag{3.23}$$

### 3.3.1 Two-dimensional SI design

In case of the SI2 design presented in Section 3.1.1, plugging the variables $u_{dk}$ derived in (3.22) into the variance formula in (3.6) yields the variance approximation

$$V\left\{\widehat{\Delta G}^{\text{co}}(a,b)\right\} \simeq N^2 \left\{c_1(a) S^2_{u_1,U} - 2c_{12}(a,b) S_{u_1 u_2,U} + c_2(b) S^2_{u_2,U}\right\},$$

see Theorem 1 in Goga, Deville and Ruiz-Gazen (2009). To obtain a variance estimator, the linearized variables may be estimated in several ways. If the overlapping sample $s_3$ only is used, the estimated linearized variables $\hat{u}_d$ are obtained from (3.23) by taking $\hat{M}_1^{\text{co}} = \hat{M}_{1,3}$ and $\hat{M}_2^{\text{co}} = \hat{M}_{2,3}$. A variance estimator is then obtained by plugging these linearized variables into (3.13). This leads to

$$v_{\text{int}}^{\text{HT}}\left\{\widehat{\Delta G}^{\text{co}}(a,b)\right\} = N^2 \left\{c_1(a) S^2_{\hat{u}_1,s_3} - 2c_{12}(a,b) S_{\hat{u}_1\hat{u}_2,s_3} + c_2(b) S^2_{\hat{u}_2,s_3}\right\}. \tag{3.24}$$

If the whole samples $s_1$ and $s_2$ are used, the estimated linearized variable $\hat{u}_d$ are obtained from (3.23) by taking $\hat{M}_1^{\text{co}} = \hat{M}_{1,1}$ and $\hat{M}_2^{\text{co}} = \hat{M}_{2,2}$. A variance estimator is then obtained by plugging these linearized variables into (3.14). This leads to

$$v_{\text{uni}}^{\text{HT}}\left\{\widehat{\Delta G}^{\text{co}}(a,b)\right\} = N^2 \left\{c_1(a) S^2_{\hat{u}_1,s_1} - 2c_{12}(a,b) S_{\hat{u}_1\hat{u}_2,s_3} + c_2(b) S^2_{\hat{u}_2,s_2}\right\}. \tag{3.25}$$

### 3.3.2 Two-dimensional multistage design

In case of the MULT2 design presented in Section 3.1.2, the linearized variables may also be estimated in several ways. For the sake of simplicity, we consider using the overlapping sample $s_3$ only so that the estimated linearized variables $\hat{u}_d$ are obtained from (3.23) by taking $\hat{M}_1^{\text{co}} = \hat{M}_{1,3}$ and $\hat{M}_2^{\text{co}} = \hat{M}_{2,3}$. A variance estimator is then obtained by plugging these linearized variables into (3.19). This leads to

$$v^{\text{HH}}\left\{\widehat{\Delta G}^{\text{co}}(a,b)\right\} = \frac{m}{m-1}\sum_{i\in s_I}\left(\frac{\widehat{\Delta u}^{i,\text{co}}(a,b)}{\pi_{Ii}} - \frac{\widehat{\Delta u}^{\text{co}}(a,b)}{m}\right)^2, \tag{3.26}$$

where $\widehat{\Delta u}^{\text{co}}(a,b)$ and $\widehat{\Delta u}^{i,\text{co}}(a,b)$ are obtained from (3.15) and (3.16), respectively, by replacing $y_{dk}$ with $\hat{u}_{dk}$.

## 3.4 Bootstrap variance estimation

Bootstrap methods have not yet been studied for the change between Gini indexes. The principles of the weighted bootstrap technique can be extended to the two-sample context, i.e. each measure $\hat{M}_{d,\diamond}$ with $d = 1,2$ and $\diamond \in \{1\bullet, 3, 2\bullet\}$ is estimated, conditionally on the samples originally selected, by some weighted bootstrap measure $\hat{M}^*_{d,\diamond}$ which enables to match, at least approximately, the two first moments of an unbiased estimator in the linear case. In Section 3.4.1, we consider a generalization of the BWO to the SI2 design. In Section 3.4.2, we propose a generalisation of the BWR to the MULT2 design.

### 3.4.1  A generalization of the BWO to the SI2 design

We first consider the SI2 design. Building a pseudo-population $U^*$ is more intricate in the two-sample case, since the variables of interest measured at waves $\tau_1$ and $\tau_2$ need to be available for each unit in $U^*$. We therefore describe a bootstrap algorithm where the overlapping sample $s_3$ only is used to build the pseudo-population $U^*$, in the spirit of the intersection variance estimator in (3.24).

Suppose that $N/n_3$ is an integer. The vectors $D_\Diamond$ are obtained by, first creating a pseudo-population $U^*$ of size $N$ by duplicating $N/n_3$ times each unit $k$ in the original sample $s_3$. A SI2 resample $s^* = (s^*_{1\bullet}, s^*_3, s^*_{\bullet 2})$ of size $(n_{1\bullet}, n_3, n_{2\bullet})$ is then selected in $U^*$. The bootstrap measures are then

$$\hat{M}^*_{d,\Diamond} \;=\; \sum_{k\in s_3} w_{\Diamond,k} D_{\Diamond,k}\delta_{y_{dk}},\tag{3.27}$$

with $D_{\Diamond,k}$ the number of times that unit $k$ is selected in the resample $s^*_\Diamond$. In the linear case, the bootstrap estimator of the parameter $\Delta t$ is then

$$\widehat{\Delta t}^{\,co*}(a,b) \;=\; b\big(\hat{t}_{y_2,s^*_{2\bullet}} - \hat{t}_{y_2,s^*_3}\big) - a\big(\hat{t}_{y_1,s^*_{1\bullet}} - \hat{t}_{y_1,s^*_3}\big) + \big(\hat{t}_{y_2,s^*_3} - \hat{t}_{y_1,s^*_3}\big),\tag{3.28}$$

where $\hat{t}_{y_d,s^*_\Diamond} = \sum_{k\in s_3} w_{\Diamond,k} D_{\Diamond,k} y_{dk}$. After some algebra, we obtain

$$E_*\left\{\widehat{\Delta t}^{\,co*}(a,b)\right\} \;=\; \widehat{\Delta t}^{\,\text{int}} \qquad\text{and}\qquad V_*\left\{\widehat{\Delta t}^{\,co*}(a,b)\right\} \;=\; \frac{1-n_3^{-1}}{1-N^{-1}}\, v_{\text{int}}^{\text{HT}}\left\{\widehat{\Delta t}^{\,co}(a,b)\right\},\tag{3.29}$$

where $\widehat{\Delta t}^{\,\text{int}}$ is given in (3.7), and $v_{\text{int}}^{\text{HT}}\big(\hat{t}_{y1}^{\text{HT}}\big)$ is given in (3.13). The proposed generalization of the BWO therefore enables to exactly match the intersection estimator of the first moment, and to approximately match the intersection estimator of the second moment for a large $n_3$.

The building of $U^*$ may be avoided by noting that under the BWO procedure, each vector $D_\Diamond$ follows a multivariate hypergeometric distribution. Therefore, the resampling weights may be directly generated. The algorithm may be adapted to the general case when $N/n_3$ is not an integer by means of any of the techniques mentioned in Section 2.4.

### 3.4.2  A generalization of the BWR for the two-dimensional multistage design

We now consider the two-dimensional two-stage sampling design with a common first-stage sample $s_I$ presented in Section 3.1.2. The proposed bootstrap procedure is similar to that described in Rao and Wu (1988). A with-replacement resample $s^*_I$ of size $m-1$ is selected by means of simple random sampling with replacement (SIR) in the original first-stage sample $s_I$. The bootstrap measures are then

$$\hat{M}^*_{d,\Diamond} \;=\; \frac{m}{m-1}\sum_{i\in s^*_I}\sum_{k\in s^i_\Diamond}\pi_{Ii}^{-1}\pi_{\Diamond k|i}^{-1}\delta_{y_{dk}} \qquad\text{where}\qquad \pi_{\Diamond k|i} = \frac{n^i_\Diamond}{N_i}.\tag{3.30}$$

It may be rewritten as

$$\hat{M}^*_{d,\Diamond} \;=\; \sum_{k\in s_\Diamond} w_{\Diamond,k} D_{\Diamond,k}\delta_{y_{dk}},\tag{3.31}$$

with $s_\diamond$ the union of the samples $s_\diamond^i$ for $i \in s_I$, and where the resampling weight $D_{\diamond,k}$ equals $m(m-1)^{-1}$ multiplied by the number of times the PSU containing $k$ is selected in $s_I^*$.

In the linear case, the bootstrap estimator of the parameter $\Delta t$ is then

$$\widehat{\Delta t}^{\,co*}(a,b) = \frac{m}{m-1} \sum_{i \in s_I^*} \pi_{Ii}^{-1} \widehat{\Delta t}^{\,i,co}(a,b) \tag{3.32}$$

where $\widehat{\Delta t}^{\,i,co}(a,b)$ is defined in (3.16). After some algebra, we obtain

$$E_*\left\{\widehat{\Delta t}^{\,co*}(a,b)\right\} = \widehat{\Delta t}^{\,co}(a,b) \quad \text{and} \quad V_*\left\{\widehat{\Delta t}^{\,co*}(a,b)\right\} = v^{HH}\left\{\widehat{\Delta t}^{\,co}(a,b)\right\}, \tag{3.33}$$

where $\widehat{\Delta t}^{\,co}(a,b)$ is given in (3.15), and $v^{HH}\left\{\widehat{\Delta t}^{\,co}(a,b)\right\}$ is given in (3.19). The proposed generalization of the BWR therefore enables to exactly match the composite estimator of the first moment, and the associated estimator of the second moment.

# 4 Simulation study

In this section, five artificial populations are first generated as described in Section 4.1. In Section 4.2, the union estimator is compared with the intersection estimator in terms of asymptotic variance. A Monte Carlo experiment is then presented in Section 4.3, and the performances of the linearization and the bootstrap are compared in case of a SI2 sampling design. A similar comparison is made in Section 4.4, in case of the bi-dimensional two-stage sampling design.

## 4.1 Simulation set-up

We generated 5 finite populations of size $N = 40,000$, each containing two study variables $y_1$ and $y_2$. The $y_{1k}$ values and the $y_{2k}$ values were generated according to the lognormal model

$$y_{dk} = \exp(\alpha_d \, \varepsilon_k). \tag{4.1}$$

The $\varepsilon_k$'s were generated according to a standard normal distribution. The values of the Gini coefficients for the five populations are presented in Table 4.1.

**Table 4.1**
**Gini coefficients for 5 populations**

| Population | Pop. 1 | Pop. 2 | Pop. 3 | Pop. 4 | Pop. 5 |
|:---:|:---:|:---:|:---:|:---:|:---:|
| $G_1$ | 0.249 | 0.298 | 0.348 | 0.397 | 0.447 |
| $G_2$ | 0.259 | 0.318 | 0.378 | 0.437 | 0.496 |
| $\Delta G$ | 0.010 | 0.020 | 0.030 | 0.040 | 0.049 |

In each of the 5 populations, the units were grouped into $M = 500$ clusters of equal size $N_0 = 80$. The clusters were built so that the intra-cluster correlation coefficient with respect to the variable $y_1$ was approximately equal to 0.20 in each population.

## 4.2 Comparison of the union estimator and of the intersection estimator

In this section, we compare the union estimator with the intersection estimator for the change between Gini indexes in terms of asymptotic variance. We consider two sampling designs: the SI2 design presented in Section 3.1.1 with $(n_{1\bullet}, n_3, n_{2\bullet}) = (1{,}000; 1{,}000; 1{,}000)$, $(1{,}000; 2{,}000; 1{,}000)$ or $(1{,}000; 4{,}000; 1{,}000)$; the MULT2 design presented in Section 3.1.2 with $m = 300$ and $(n_{1\bullet}^i, n_3^i, n_{2\bullet}^i) = (10; 10; 10)$, $(10; 20; 10)$ or $(10; 40; 10)$.

For each population, we compute the asymptotic variance $V_{\text{lin}}(\widehat{\Delta G}^{\text{uni}})$ of the union estimator, and the asymptotic variance $V_{\text{lin}}(\widehat{\Delta G}^{\text{int}})$ of the intersection estimator. So as to compare them, we compute the relative efficiency defined as

$$\text{RE}\left\{\widehat{\Delta G}^{\cdot}\right\} = \frac{V_{\text{lin}}\left\{\widehat{\Delta G}^{(\cdot)}\right\}}{V_{\text{lin}}\left\{\widehat{\Delta G}^{\text{opt}}\right\}}, \tag{4.2}$$

with $\widehat{\Delta G}^{\text{opt}}$ the optimal estimator.

The results are presented in Table 4.2. The union estimator is highly inefficient. Its asymptotic variance is 15 to 244 times higher than that of the intersection estimator for SI2, and 2 to 44 times higher than that of the intersection estimator for MULT2. The difference between both estimators tends to decrease when the sample size of the common sample increases and/or when $\Delta G$ increases. On the other hand, the intersection estimator is slightly less efficient than the optimal estimator for SI2, with RE ranging from 1.33 to 2.46, and approximately as efficient as the optimal estimator for MULT2, with RE ranging from 1.02 to 1.12. This supports the heuristic reasoning in Section 3.1.1. In view of the poor performance of the union estimator, and of the good performance of the intersection estimator, we confine our attention to the latter in the remainder of the simulation study.

**Table 4.2**
**Relative efficiency of the union estimator and of the intersection variance estimator for 5 populations**

| Design | Sample size | Pop. 1 | | Pop. 2 | | Pop. 3 | | Pop. 4 | | Pop. 5 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $\widehat{\Delta G}^{\text{uni}}$ | $\widehat{\Delta G}^{\text{int}}$ | $\widehat{\Delta G}^{\text{uni}}$ | $\widehat{\Delta G}^{\text{int}}$ | $\widehat{\Delta G}^{\text{uni}}$ | $\widehat{\Delta G}^{\text{int}}$ | $\widehat{\Delta G}^{\text{uni}}$ | $\widehat{\Delta G}^{\text{int}}$ | $\widehat{\Delta G}^{\text{uni}}$ | $\widehat{\Delta G}^{\text{int}}$ |
| SI2 | $n_3 = 1{,}000$ | 600.22 | 2.46 | 200.23 | 2.27 | 96.72 | 2.10 | 58.73 | 1.96 | 39.35 | 1.85 |
| | $n_3 = 2{,}000$ | 410.23 | 1.84 | 141.71 | 1.76 | 70.71 | 1.68 | 44.18 | 1.61 | 30.33 | 1.54 |
| | $n_3 = 4{,}000$ | 250.02 | 1.47 | 88.40 | 1.43 | 45.17 | 1.40 | 28.86 | 1.36 | 20.23 | 1.33 |
| MULT2 | $n_3^i = 10$ | 49.10 | 1.12 | 19.89 | 1.13 | 11.83 | 1.14 | 8.84 | 1.15 | 7.28 | 1.16 |
| | $n_3^i = 20$ | 23.08 | 1.05 | 9.75 | 1.05 | 6.08 | 1.05 | 4.73 | 1.06 | 4.04 | 1.07 |
| | $n_3^i = 40$ | 9.15 | 1.02 | 4.25 | 1.02 | 2.90 | 1.02 | 2.41 | 1.02 | 2.16 | 1.02 |

## 4.3 Comparison of linearization and bootstrap for the SI2 design

In this section, we compare the linearization and bootstrap for variance estimation and for producing confidence intervals, in case of the intersection estimator for the change between Gini indexes under the SI2 sampling design. From each population, we selected $B = 10,000$ two-dimensional samples by means of the SI2 design indexed by $(n_{1\bullet}, n_3, n_{2\bullet}) = (1,000; 1,000; 1,000)$, $(n_{1\bullet}, n_3, n_{2\bullet}) = (1,000; 2,000; 1,000)$ or $(n_{1\bullet}, n_3, n_{2\bullet}) = (1,000; 4,000; 1,000)$. In each sample, we computed the intersection estimator $\widehat{\Delta G}^{\text{int}}$ of the change between Gini indexes. For this estimator, we computed (i) the linearization variance estimator $v_{\text{int}}(\widehat{\Delta G}^{\text{int}})$ given in (3.24), and (ii) the Bootstrap variance estimator $v_{\text{BWO}}(\widehat{\Delta G})$, following the Bootstrap procedure described in Section 3.4.1.

To measure the bias of a variance estimator $v(\widehat{\Delta G})$, we used the Monte Carlo Percent Relative Bias

$$\text{RB}\left\{v\left(\widehat{\Delta G}\right)\right\} = 100 \times \frac{B^{-1}\sum_{b=1}^{B} v\left(\widehat{\Delta G}_b\right) - \text{MSE}\left(\widehat{\Delta G}\right)}{\text{MSE}\left(\widehat{\Delta G}\right)}, \tag{4.3}$$

where $v(\widehat{\Delta G}_b)$ denotes the estimator $v(\widehat{\Delta G})$ in the $b^{\text{th}}$ sample, and $\text{MSE}(\widehat{\Delta G})$ is a simulation-based approximation of the true mean square error of $\widehat{\Delta G}$, obtained from an independent run of 100,000 simulations. As a measure of stability of $v(\widehat{\Delta G})$, we used the Relative Stability

$$\text{RS}\left\{v\left(\widehat{\Delta G}\right)\right\} = \frac{\left[B^{-1}\sum_{b=1}^{B}\left\{v\left(\widehat{\Delta G}\right) - \text{MSE}\left(\widehat{\Delta G}\right)\right\}^2\right]^{1/2}}{\text{MSE}\left(\widehat{\Delta G}\right)}. \tag{4.4}$$

Finally, we compared the coverage rates of (i) the normality-based confidence interval with use of the linearization variance estimator and (ii) the confidence interval associated to the percentile Bootstrap. The bootstrap variance estimators and the bootstrap confidence intervals are based on $C = 1,000$ bootstrap replications. Error rates of the confidence intervals (with nominal one-tailed error rate of 2.5% in each tail) are compared. The comparison with nominal error rate of 5% gave no qualitative difference and is thus omitted.

The results are presented in Table 4.3. Both variance estimators are negatively biased. This bias is moderate (less than 5% ) in most cases, except for the smaller sample size $n = 1,000$, and for the population $U_5$ with the highest value of $\Delta G$. The bootstrap variance estimator is systematically slightly more biased than the linearization variance estimator, but the difference decreases as the sample size increases. For both variance estimators, the instability increases with $\Delta G$. The Bootstrap variance estimator is slightly more stable for the smaller sample size $n = 1,000$, but the situation is reversed when the sample size increases. Turning to the coverage of the confidence intervals, both methods lead to under-coverage which is consistent with the negative bias of both variance estimators. The normality-based confidence intervals show a slightly better coverage than the bootstrap percentile confidence intervals. For both confidence intervals, the under-coverage is more acute when $\Delta G$ increases, and reduces when the sample size increases.

**Table 4.3**
**Relative Bias, Relative Stability and Nominal One-Tailed Error Rates for linearization and Bootstrap variance estimation of the intersection estimator of the change between Gini indexes for 5 populations and with the SI2 sampling design**

| Pop. | Linearization | | | | | Bootstrap | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | **RB** | **RS** | **L** | **U** | **L+U** | **RB** | **RS** | **L** | **U** | **L+U** |
| | Sample size $(n_{1\bullet}, n_3, n_{2\bullet}) = (1{,}000; 1{,}000; 1{,}000)$ | | | | | | | | | |
| Pop. 1 | -1.41 | 24.6 | 1.8 | 4.5 | 6.3 | -1.83 | 24.6 | 1.8 | 4.9 | 6.7 |
| Pop. 2 | -1.98 | 32.4 | 1.6 | 5.2 | 6.8 | -2.64 | 32.1 | 1.7 | 5.9 | 7.6 |
| Pop. 3 | -2.80 | 41.9 | 1.3 | 6.3 | 7.7 | -3.83 | 40.9 | 1.3 | 7.0 | 8.3 |
| Pop. 4 | -4.00 | 52.5 | 1.0 | 7.7 | 8.7 | -5.57 | 50.6 | 1.1 | 8.2 | 9.3 |
| Pop. 5 | -5.80 | 64.0 | 1.0 | 9.2 | 10.1 | -8.11 | 60.6 | 0.8 | 9.9 | 10.7 |
| | Sample size $(n_{1\bullet}, n_3, n_{2\bullet}) = (1{,}000; 2{,}000; 1{,}000)$ | | | | | | | | | |
| Pop. 1 | -1.38 | 17.3 | 1.6 | 3.7 | 5.3 | -1.67 | 17.8 | 1.8 | 4.1 | 5.9 |
| Pop. 2 | -1.64 | 23.0 | 1.4 | 4.3 | 5.8 | -2.05 | 23.2 | 1.4 | 4.7 | 6.1 |
| Pop. 3 | -1.99 | 30.1 | 1.2 | 5.0 | 6.2 | -2.58 | 30.0 | 1.1 | 5.3 | 6.4 |
| Pop. 4 | -2.50 | 38.4 | 1.0 | 6.0 | 6.9 | -3.38 | 37.9 | 1.0 | 6.3 | 7.3 |
| Pop. 5 | -3.30 | 47.9 | 0.7 | 7.2 | 7.9 | -4.62 | 46.7 | 0.7 | 7.5 | 8.2 |
| | Sample size $(n_{1\bullet}, n_3, n_{2\bullet}) = (1{,}000; 4{,}000; 1{,}000)$ | | | | | | | | | |
| Pop. 1 | -0.60 | 11.9 | 2.0 | 3.4 | 5.3 | -0.68 | 12.8 | 2.1 | 3.4 | 5.5 |
| Pop. 2 | -0.67 | 15.9 | 1.8 | 3.7 | 5.6 | -0.80 | 16.5 | 2.0 | 3.9 | 5.9 |
| Pop. 3 | -0.83 | 20.8 | 1.8 | 4.4 | 6.2 | -1.03 | 21.3 | 1.9 | 4.4 | 6.3 |
| Pop. 4 | -1.13 | 26.7 | 1.5 | 5.0 | 6.6 | -1.46 | 26.9 | 1.6 | 5.0 | 6.6 |
| Pop. 5 | -1.64 | 33.4 | 1.4 | 5.8 | 7.1 | -2.18 | 33.5 | 1.4 | 5.8 | 7.1 |

## 4.4 Comparison of linearization and bootstrap for the MULT2 design

In this section, we compare the linearization and bootstrap for variance estimation and for producing confidence intervals, in case of the intersection estimator for the change between Gini indexes under the MULT2 sampling design presented in Section 3.1.2. From each population, we selected $B = 10{,}000$ two-dimensional two-stage samples by means of the MULT2 design indexed by $m = 300$ and $(n_{1\bullet}^i, n_3^i, n_{2\bullet}^i) = (10; 10; 10)$, $(10; 20; 10)$ or $(10; 40; 10)$. In each sample, we computed the intersection estimator $\widehat{\Delta G}^{\text{int}}$ of the change between Gini indexes. For this estimator, we computed (i) the linearization variance estimator $v^{\text{HH}}\{\widehat{\Delta G}^{\text{co}}(a, b)\}$ given in (3.26), and (ii) the Bootstrap variance estimator $v_{\text{BWR}}(\widehat{\Delta G}^{\text{int}})$, following the Bootstrap procedure described in Section 3.4.2.

To measure the bias of a variance estimator $v(\widehat{\Delta G})$, we used the Monte Carlo Percent Relative Bias defined in equation (4.3), and the Relative Stability defined in equation (4.4). The true mean square error of $\widehat{\Delta G}$ was obtained from an independent run of 100,000 simulations. Also, we compared the coverage rates of (i) the normality-based confidence interval with use of the linearization variance estimator and (ii) the confidence interval associated to the percentile Bootstrap. The bootstrap variance estimators and the bootstrap confidence intervals are based on $C = 1{,}000$ bootstrap replications. Error rates of the confidence intervals (with nominal one-tailed error rate of 2.5% in each tail) are compared. The comparison with nominal error rate of 5% gave no qualitative difference and is thus omitted.

The results are presented in Table 4.4. Both variance estimators are approximately unbiased for small values of $\Delta G$, but show a moderate negative bias which increases with $\Delta G$. The bootstrap variance estimator is more biased than the linearization variance estimator. For both variance estimators, the

instability increases with $\Delta G$. The Bootstrap variance estimator is slightly more stable than the linearization variance estimator. Both methods lead to an under-coverage which is consistent with the negative bias of both variance estimators. The normality-based confidence intervals perform slightly better. For both confidence intervals, the under-coverage is more acute when $\Delta G$ increases, and reduces when the sample size increases.

**Table 4.4**
**Relative Bias, Relative Stability and Nominal One-Tailed Error Rates for linearization and Bootstrap variance estimation of the intersection estimator of the Gini Coefficient Change for 5 populations and with the MULT2 sampling design**

| Pop. | Linearization | | | | | Bootstrap | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | **RB** | **RS** | **L** | **U** | **L+U** | **RB** | **RS** | **L** | **U** | **L+U** |
| | Sample sizes $m = 300$ and $(n_{1\bullet}^i, n_{3}^i, n_{2\bullet}^i) = (10; 10; 10)$ | | | | | | | | | |
| Pop. 1 | 1.23 | 33.8 | 0.6 | 4.9 | 5.5 | 1.09 | 33.2 | 0.6 | 6.0 | 6.6 |
| Pop. 2 | 0.64 | 41.1 | 0.8 | 5.5 | 6.3 | -0.20 | 39.7 | 0.6 | 6.5 | 7.1 |
| Pop. 3 | -0.42 | 48.7 | 0.7 | 7.1 | 7.8 | -2.05 | 46.6 | 0.7 | 8.4 | 9.1 |
| Pop. 4 | -2.07 | 56.4 | 0.8 | 8.4 | 9.2 | -4.47 | 53.3 | 0.6 | 9.6 | 10.2 |
| Pop. 5 | -4.44 | 63.7 | 0.9 | 9.2 | 10.1 | -7.56 | 59.5 | 0.4 | 10.3 | 10.7 |
| | Sample sizes $m = 300$ and $(n_{1\bullet}^i, n_{3}^i, n_{2\bullet}^i) = (10; 20; 10)$ | | | | | | | | | |
| Pop. 1 | 1.70 | 32.6 | 1.5 | 4.9 | 6.4 | -1.70 | 32.3 | 1.5 | 6.0 | 7.5 |
| Pop. 2 | 1.10 | 39.0 | 1.4 | 5.4 | 6.8 | -1.91 | 38.3 | 1.5 | 6.9 | 8.4 |
| Pop. 3 | 0.17 | 45.6 | 1.2 | 7.4 | 8.6 | -2.49 | 44.4 | 1.1 | 7.7 | 8.8 |
| Pop. 4 | -1.17 | 52.0 | 1.0 | 9.0 | 10.0 | -3.58 | 50.3 | 0.8 | 9.7 | 10.5 |
| Pop. 5 | -3.03 | 57.9 | 0.9 | 10.4 | 11.3 | -5.35 | 55.4 | 0.7 | 11.0 | 11.7 |
| | Sample sizes $m = 300$ and $(n_{1\bullet}^i, n_{3}^i, n_{2\bullet}^i) = (10; 40; 10)$ | | | | | | | | | |
| Pop. 1 | -0.99 | 32.1 | 1.2 | 6.1 | 7.3 | -3.21 | 32.2 | 1.7 | 6.7 | 8.4 |
| Pop. 2 | -1.68 | 38.3 | 1.4 | 6.7 | 8.1 | -3.70 | 38.3 | 1.4 | 7.6 | 9.0 |
| Pop. 3 | -2.58 | 44.6 | 1.3 | 7.5 | 8.8 | -4.40 | 44.5 | 1.2 | 8.9 | 10.1 |
| Pop. 4 | -3.78 | 50.6 | 1.1 | 8.9 | 10.0 | -5.50 | 50.1 | 0.9 | 10.6 | 11.5 |
| Pop. 5 | -5.39 | 55.9 | 0.8 | 10.9 | 11.7 | -7.16 | 54.8 | 0.6 | 12.8 | 13.4 |

# 5 Conclusion

In this paper, we considered the estimation of the change between Gini indexes. We presented the class of composite estimators introduced by Goga, Deville and Ruiz-Gazen (2009), and studied more particularly the intersection estimator which makes use of the common sample only, and the union estimator which makes use of the whole available samples. We justified both heuristically and through the simulation study in Section 4.2 that the intersection estimator can be close to the optimal estimator, while the union estimator exhibits poor performances in all the scenarios considered. The intersection estimator is also easy to compute, while the optimal estimator involves unknown quantities which need to be estimated in practice. We therefore advocate for the use of the intersection estimator for estimating the change between Gini indexes.

We also compared linearization and bootstrap for variance estimation and for producing confidence intervals. In the scenarios that we considered in the simulation study, the linearization performed better with usually smaller relative biases for the variance estimator, and better coverage rates with normality-based

confidence intervals than with percentile confidence intervals. Bootstrap $-t$ confidence intervals (not considered in the simulation study) would be a competitor of interest, but due to the intensive computational work involved, they are less attractive for a data user. Linearization has also the advantage to offer a unified approach suitable for any sampling design, while a specific sampling design usually requires a specific bootstrap procedure, as illustrated with the BWO for SI sampling and the BWR for multistage sampling.

From the simulation study, we note that the coverage rates may not be well respected neither with linearization nor bootstrap, particularly in the multistage context and even with large sample sizes. There is a need for confidence intervals with better coverage rates under a reasonable computational burden. This is a matter for further research.

## Acknowledgements

## Appendix

### Proof of equation (3.6)

From (3.3), we have $\widehat{\Delta t}^{co} = N(A^\top X)$, where $X = (\overline{y}_{2,s_{2\bullet}} - \overline{y}_{2,s_3}, \overline{y}_{1,s_{1\bullet}} - \overline{y}_{1,s_3}, \overline{y}_{2,s_3} - \overline{y}_{1,s_3})^\top$ and $A = (b, -a, 1)^\top$. This leads to

$$V\left\{\widehat{\Delta t}^{co}\right\} = N^2\{A^\top V(X)A\}. \tag{A.1}$$

We compute the elements in $V(X)$ separately. We have

$$\begin{aligned}V(\overline{y}_{2,s_3} - \overline{y}_{1,s_3}) &= \left(\frac{1}{n_3} - \frac{1}{N}\right)S^2_{y_2 - y_1, U} \\ &= \left(\frac{1}{n_3} - \frac{1}{N}\right)\left(S^2_{y_2,U} + S^2_{y_1,U} - 2S_{y_1y_2,U}\right).\end{aligned}$$

Also, since $E(\overline{y}_{2,s_{2\bullet}} - \overline{y}_{2,s_3} \mid s_2) = 0$, we have

$$\begin{aligned}V(\overline{y}_{2,s_{2\bullet}} - \overline{y}_{2,s_3}) &= \mathrm{EV}(\overline{y}_{2,s_{2\bullet}} - \overline{y}_{2,s_3} \mid s_2), \\ &= \mathrm{EV}\left(\frac{n_2}{n_{2\bullet}}\overline{y}_{2,s_2} - \frac{n_3}{n_{2\bullet}}\overline{y}_{2,s_3} - \overline{y}_{2,s_3} \mid s_2\right) \\ &= \left(1 + \frac{n_3}{n_{2\bullet}}\right)^2 \mathrm{EV}(\overline{y}_{2,s_3} \mid s_2) \\ &= \left(1 + \frac{n_3}{n_{2\bullet}}\right)^2\left(\frac{1}{n_3} - \frac{1}{n_2}\right)S^2_{y_2,U} \\ &= \frac{n_2}{n_3(n_2 - n_3)}S^2_{y_2,U}\end{aligned}$$

and

$$\text{Cov}\left(\bar{y}_{2,s_2\bullet} - \bar{y}_{2,s_3}, \bar{y}_{2,s_3} - \bar{y}_{1,s_3}\right) = \text{ECov}\left(\bar{y}_{2,s_2\bullet} - \bar{y}_{2,s_3}, \bar{y}_{2,s_3} - \bar{y}_{1,s_3} \mid s_2\right)$$

$$= \text{ECov}\left(\frac{n_2}{n_{2\bullet}} \bar{y}_{2,s_2} - \frac{n_3}{n_{2\bullet}} \bar{y}_{2,s_3} - \bar{y}_{2,s_3}, \bar{y}_{2,s_3} - \bar{y}_{1,s_3} \mid s_2\right)$$

$$= -\left(1 + \frac{n_3}{n_{2\bullet}}\right) \text{ECov}\left(\bar{y}_{2,s_3}, \bar{y}_{2,s_3} - \bar{y}_{1,s_3} \mid s_2\right)$$

$$= \left(1 + \frac{n_3}{n_{2\bullet}}\right)\left(\frac{1}{n_3} - \frac{1}{n_2}\right)\left(S^2_{y_2,U} - S_{y_1 y_2,U}\right)$$

$$= -\frac{1}{n_3}\left(S^2_{y_2,U} - S_{y_1 y_2,U}\right).$$

Similar arguments lead to

$$V\left(\bar{y}_{1,s_1\bullet} - \bar{y}_{1,s_3}\right) = \frac{n_1}{n_3\left(n_1 - n_3\right)} S^2_{y_1,U},$$

$$\text{Cov}\left(\bar{y}_{1,s_1\bullet} - \bar{y}_{1,s_3}, \bar{y}_{2,s_3} - \bar{y}_{1,s_3}\right) = \frac{1}{n_3}\left(S^2_{y_1,U} - S_{y_1 y_2,U}\right).$$

Finally, we consider $\text{Cov}\left(\bar{y}_{2,s_2\bullet} - \bar{y}_{2,s_3}, \bar{y}_{1,s_1\bullet} - \bar{y}_{1,s_3}\right)$. We first compute $\text{Cov}\left(\bar{y}_{2,s_2\bullet}, \bar{y}_{1,s_1\bullet}\right)$, which may be written as

$$\text{Cov}\left(\bar{y}_{2,s_2\bullet}, \bar{y}_{1,s_1\bullet}\right) = \text{Cov}\left(E\left(\bar{y}_{2,s_2\bullet} \mid s_1\bullet\right), E\left(\bar{y}_{1,s_1\bullet} \mid s_1\bullet\right)\right)$$

$$= \text{Cov}\left(\bar{y}_{2,U\setminus s_1\bullet}, \bar{y}_{1,s_1\bullet}\right)$$

$$= \text{Cov}\left(\frac{N}{N - n_{1\bullet}} \bar{y}_{2,U} - \frac{n_{1\bullet}}{N - n_{1\bullet}} \bar{y}_{2,s_1\bullet}, \bar{y}_{1,s_1\bullet}\right)$$

$$= -\frac{n_{1\bullet}}{N - n_{1\bullet}} \text{Cov}\left(\bar{y}_{2,s_1\bullet}, \bar{y}_{1,s_1\bullet}\right)$$

$$= -\frac{n_{1\bullet}}{N - n_{1\bullet}}\left(\frac{1}{n_{1\bullet}} - \frac{1}{N}\right) S_{y_1 y_2,U}$$

$$= -\frac{1}{N} S_{y_1 y_2,U}.$$

Similar arguments lead to

$$\text{Cov}\left(\bar{y}_{2,s_2\bullet}, \bar{y}_{1,s_3}\right) = \text{Cov}\left(\bar{y}_{2,s_3}, \bar{y}_{1,s_1\bullet}\right) = -\frac{1}{N} S_{y_1 y_2,U}.$$

We obtain

$$\mathrm{Cov}\left(\bar{y}_{2,s_{2\bullet}} - \bar{y}_{2,s_3}, \bar{y}_{1,s_{1\bullet}} - \bar{y}_{1,s_3}\right) = \frac{1}{N} S_{y_1 y_2, U} + \mathrm{Cov}\left(\bar{y}_{2,s_3}, \bar{y}_{1,s_3}\right)$$

$$= \frac{1}{N} S_{y_1 y_2, U} + \left(\frac{1}{n_3} - \frac{1}{N}\right) S_{y_1 y_2, U}$$

$$= \frac{1}{n_3} S_{y_1 y_2, U}.$$

In summary, we obtain

$$V(X) = \begin{pmatrix} \dfrac{n_2}{n_3(n_2 - n_3)} S^2_{y_2, U} & \dfrac{1}{n_3} S_{y_1 y_2, U} & -\dfrac{1}{n_3}\left(S^2_{y_2, U} - S_{y_1 y_2, U}\right) \\[3ex] & \dfrac{n_1}{n_3(n_1 - n_3)} S^2_{y_1, U} & \dfrac{1}{n_3}\left(S^2_{y_1, U} - S_{y_1 y_2, U}\right) \\[3ex] & & \left(\dfrac{1}{n_3} - \dfrac{1}{N}\right)\left(S^2_{y_2, U} + S^2_{y_1, U} - 2 S_{y_1 y_2, U}\right) \end{pmatrix}$$

which, along with (A.1), leads to (3.6).

# References

Antal, E., and Tillé, Y. (2011). A direct bootstrap method for complex sampling designs from a finite population. *Journal of the American Statistical Association*, 106, 534-543.

Barrett, G.F., and Donald, S.G. (2009). Statistical inference with generalized Gini indices of inequality, poverty, and welfare. *Journal of Business and Economic Statistics*, 27, 1-17.

Beaumont, J.-F., and Patak, Z. (2012). On the generalized bootstrap for sample surveys with special attention to Poisson sampling. *International Statistical Review*, 80, 127-148.

Berger, Y.G. (2004). Variance estimation for measures of change in probability sampling. *Canadian Journal of Statistics*, 32, 451-467.

Berger, Y.G. (2008). A note on the asymptotic equivalence of jackknife and linearization variance estimation for the Gini coefficient. *Journal of Official Statistics*, 24, 541-555.

Bertail, P., and Combris, P. (1997). Bootstrap généralisé d'un sondage. *Annales d'Économie et de Statistique*, 46, 49-83.

Bhattacharya, D. (2007). Inference on inequality from household survey data. *Journal of Econometrics*, 137, 674-707.

Bickel, P.J., and Freedman, D.A. (1984). Asymptotic normality and the bootstrap in stratified sampling. *The Annals of Statistics*, 12, 470-482.

Booth, J.G., Butler, R.W. and Hall, P. (1994). Bootstrap methods for finite populations. *Journal of the American Statistical Association*, 89, 1282-1289.

Brändén, P., and Jonasson, J. (2012). Negative dependence in sampling. *Scandinavian Journal of Statistics*, 39, 830-838.

Campbell, C. (1980). A different view of finite population estimation. *Proceedings of the Survey Research Methods Section,* American Statistical Association, 319-324.

Chao, M.-T., and Lo, S.-H. (1985). A Bootstrap method for finite population. *Sankhyā, Series A*, 47, 3, 399-405.

Chauvet, G. (2007). Méthodes de Bootstrap en population finie. Ph.D. dissertation, Université Rennes 2.

Chauvet, G. (2015). Coupling methods for multistage sampling. *The Annals of Statistics*, 43(6), 2484-2506.

Chen, J., and Rao, J.N.K. (2007). Asymptotic normality under two-phase sampling designs. *Statistica Sinica*, 17, 1047-1064.

Davison, A.C., and Hinkley, D.V. (1997). *Bootstrap Methods and their Application*. Cambridge University Press.

Davison, A.C., and Sardy, S. (2007). Resampling variance estimation in surveys with missing data. *Journal of Official Statistics*, 23, 3, 371-386.

Demnati, A., and Rao, J.N.K. (2004). Linearization variance estimators for survey data. *Survey Methodology*, 30, 1, 17-26. Paper available at https://www150.statcan.gc.ca/n1/pub/12-001-x/2004001/article/6991-eng.pdf.

Deville, J.-C. (1997). Estimation de la variance du coefficient de Gini mesurée par sondage. *Actes des Journées de Méthodologie Statistique*, *Insee Méthodes*.

Deville, J.-C. (1999). Variance estimation for complex statistics and estimators: Linearization and residual techniques. *Survey Methodology*, 25, 2, 193-203. Paper available at https://www150.statcan.gc.ca/n1/pub/12-001-x/1999002/article/4882-eng.pdf.

Druckman, A., and Jackson, T. (2008). Measuring resource inequalities: The concepts and methodology for an area-based Gini coefficient. *Ecological Economics*, 65, 242-252.

Gini, C. (1914). Sulla misura della concentrazione e della variabilità dei caratteri. *Atti del Reale Istituto Veneto di Scienze Lettere ed Arti*.

Glasser, G.J. (1962). Variance formulas for the mean difference and coefficient of concentration. *Journal of the American Statistical Association*, 57, 648-654.

Goga, C. (2003). Estimation de la variance dans les sondages à plusieurs échantillons et prise en compte de l'information auxiliaire par des modèles nonparamétriques. Ph.D. dissertation, Université Rennes 2.

Goga, C., and Ruiz-Gazen, A. (2014). Efficient estimation of nonlinear finite population parameters using nonparametrics. *Journal of the Royal Statistical Society B*, 76, 113-140.

Goga, C., Deville, J.-C. and Ruiz-Gazen, A. (2009). Composite estimation and linearization method for two-sample survey data. *Biometrika*, 96, 691-709.

Gordon, L. (1983). Successive sampling in large finite populations. *Annals of Statistics*, 11, 702-706.

Graczyk, P.P. (2007). Gini coefficient: A new way to express selectivity of kinase inhibitors against a family of Kinases. *Journal of Medicinal Chemistry*, 50, 5773-5779.

Gross, S.T. (1980). Median estimation in sample surveys. *ASA Proceedings of Survey Research*, 181-184.

Groves-Kirkby, C.J., Denman, A.R. and Phillips, P.S. (2009). Lorenz Curve and Gini coefficient: Novel tools for analysing seasonal variation of environmental radon gas. *Journal of Environmental Management*, 90, 2480-2487.

Hájek, J. (1960). Limiting distributions in simple random sampling from a finite population. *Tud. Akad. Mat. Kutatò Int. Közl.*, 5, 361-374.

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

Hájek, J. (1964). Asymptotic theory of rejective sampling with varying probabilities from a finite population. *Annals of Mathematical Statistics*, 35, 1491-1523.

Isaki, C.T., and Fuller, W.A. (1982). Survey design under the regression superpopulation model. *Journal of the American Statistical Association*, 77, 89-96.

Karagiannis, E., and Kovačević, M.S. (2000). A method to calculate the jackknife variance estimator for the Gini coefficient. *Oxford Bulletin of Economics and Statistics*, 62, 119-122.

Kovačević, M.S., and Binder, D.A. (1997). Variance estimation for measures of income inequality and polarization - The estimating equation approach. *Journal of Official Statistics,* 13, 41-58.

Krewski, D., and Rao, J.N.K. (1981). Inference from stratified samples: Properties of the linearization, jackknife and balanced repeated replication methods. *Annals of Statistics,* 9, 1010-1019.

Lai, D., Huang, J., Risser, J.M. and Kapadia, A.S. (2008). Statistical properties of generalized Gini coefficient with application to health inequality measurement. *Social Indicator Research,* 87, 249-258.

Langel, M., and Tillé, Y. (2013). Variance estimation of the Gini index: Revisiting a result several times published. *Journal of the Royal Statistical Society, Series A,* 176, 521-540.

Lisker, T. (2008). Is the Gini coefficient a stable measure on galaxy structure? *The Astrophysical Journal Supplement Series,* 179, 319-325.

Navarro, V., Muntaner, C., Borrell, C., Benach, J., Quiroga, A., Rodríguez-Sanz, M., Vergès, N. and Pasarín, M.I. (2006). Politics and health outcomes. *The Lancet*, 18, 1033-1037.

Nygård, F., and Sandström, A. (1985). The estimation of the Gini and the entropy inequality parameters in finite populations. *Journal of Official Statistics*, 1, 4, 399-412.

Ohlsson, E. (1986). Asymptotic normality of the Rao-Hartley-Cochran estimator: An application of the martingale CLT. *Scandinavian Journal of Statistics*, 13, 17-28.

Ohlsson, E. (1989). Asymptotic normality for two-stage sampling from a finite population. *Probability Theory and Related Fields*, 81, 341-352.

Pires, A.M., and Branco, J.A. (2002). Partial influence functions. *Journal of Multivariate Analysis,* 83, 451-468.

Presnell, B., and Booth, J.G. (1994). *Resampling Methods for Sample Surveys*. Technical report.

Qin, Y., Rao, J.N.K. and Wu, C. (2010). Empirical likelihood confidence intervals for the Gini measure of income inequality. *Economic Modelling,* 27, 1429-1435.

Qualité, L., and Tillé, Y. (2008). Variance estimation of changes in repeated surveys and its application to the Swiss survey of value added. *Survey Methodology*, 34, 2, 173-181. Paper available at https://www150.statcan.gc.ca/n1/pub/12-001-x/2008002/article/10758-eng.pdf.

Rao, J.N.K., and Wu, C.F.J. (1988). Resampling inference with complex survey data. *Journal of the American Statistical Association,* 83, 231-241.

Reid, N. (1981). Influence functions for censored data. *The Annals of Statistics*, 9, 78-92.

Rosén, B. (1972). Asymptotic theory for successive sampling with varying probabilities without replacement. I, II. *Annals of Mathematical Statistics*, 43, 373-397, 748-776.

Saegusa, T., and Wellner, J.A. (2013). Weighted likelihood estimation under two-phase sampling. *The Annals of Statistics*, 41, 269-295.

Sandström, A., Wretman, J.H. and Waldèn, B. (1985). Variance estimators of the Gini coefficient - Simple random sampling. *Metron,* 43, 41-70.

Sandström, A., Wretman, J.H. and Waldèn, B. (1988). Variance estimators of the Gini coefficient - Probability sampling. *Journal of Business and Economic Statistics,* 6, 113-119.

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

Sen, P.K. (1980). Limit theorems for an extended coupon collector's problem and for successive subsampling with varying probabilities. *Calcutta Statistical Association Bulletin,* 29, 113-132.

Shao, J., and Tu, D. (1995). *The Jackknife and the Bootstrap*. Springer.

Sitter, R.R. (1992a). A resampling procedure for complex survey data. *Journal of the American Statistical Association,* 87, 755-765.

Sitter, R.R. (1992b). Comparing three bootstrap methods for survey data. *Canadian Journal of Statistics,* 20, 135-154.

Tam, S.M. (1984). On covariances from overlapping samples. *The American Statistician,* 38, 288-289.

Yitzhaki, S. (1991). Calculating jackknife variance estimators for parameters of the Gini method. *Journal of Business and Economic Statistics,* 9, 235-239.