

Calibrated bootstrap and saddlepoint approximations of finite population L -statistics

Andrius Čiginas^a and Dalius Pumputis^b

^a Institute of Data Science and Digital Technologies, Vilnius University,
Akademijos str. 4, LT-08663 Vilnius, Lithuania

^b Faculty of Fundamental Sciences, Vilnius Gediminas Technical University,
Saulėtekis ave. 11, LT-10223 Vilnius, Lithuania

(e-mail: andrius.ciginas@mii.vu.lt; dalius.pumputis@vgtu.lt)

Received July 26, 2018; revised January 22, 2019

Abstract. We propose two methods to approximate the distribution function of a Studentized linear combination of order statistics for a simple random sample drawn without replacement from a finite population. Using auxiliary data available for the population units, the first method modifies a nonparametric bootstrap approximation, and the second one corrects an empirical saddlepoint approximation based on the bootstrap. We conclude from simulations that, on the tails of distribution of interest, both approximations improve their initial versions and alternative Edgeworth approximations.

MSC: 62E20, 62D05

Keywords: sampling without replacement, auxiliary information, Studentized statistic, Hoeffding decomposition, synthetic estimation, jackknife

1 Introduction

Consider a study variable x with real values $\mathcal{X} = \{x_1, \dots, x_N\}$ in the population $\mathcal{U} = \{1, \dots, N\}$. Let $\mathbb{X} = \{X_1, \dots, X_n\}$ be the measurements of the simple random sample units $\{1, \dots, n\}$, $n < N$, drawn without replacement from \mathcal{U} . The L -statistic

$$L = L_n(\mathbb{X}) = \frac{1}{n} \sum_{j=1}^n c_{j,n} X_{j:n} \quad (1.1)$$

is a linear combination of the order statistics $X_{1:n} \leq \dots \leq X_{n:n}$ of \mathbb{X} with real coefficients

$$c_{j,n} = J\left(\frac{j}{n+1}\right), \quad J: (0, 1) \rightarrow \mathbb{R},$$

called weights. The sample mean, Gini's mean difference, and trimmed means are particular cases of (1.1).

We aim to estimate the distribution function

$$F_S(y) = P\{\hat{\sigma}_J^{-1}(L - EL) \leq y\} \quad (1.2)$$

of the Studentized L -statistic, where

$$\hat{\sigma}_J^2 = \hat{\sigma}_J^2(\mathbb{X}) = \left(1 - \frac{n}{N}\right) \frac{n-1}{n} \sum_{k=1}^n (L_{(k)} - \bar{L})^2, \quad \bar{L} = \frac{1}{n} \sum_{k=1}^n L_{(k)}, \quad (1.3)$$

is the jackknife estimator of the variance $\sigma^2 = \text{Var } L$. Here $L_{(k)} = L_{n-1}(\mathbb{X} \setminus \{X_k\})$, $1 \leq k \leq n$, are L -statistics with weights $c_{j,n-1} = J(j/n)$, $1 \leq j \leq n-1$. The knowledge of (1.2) allows us to test hypotheses and construct confidence intervals for the parameter EL . In practice, the standard normal approximation Φ is commonly applied to (1.2) for large sample sizes, but for small to moderate n , this approximation is quite inaccurate. Typically, it should have the absolute error $O(n^{-1/2})$ according to the Berry–Esseen theorems.

An improvement over the normal approximation is provided by the one-term Edgeworth expansion constructed by [4] for Studentized symmetric finite population statistics (including L -statistics). The jackknife estimators of unknown parameters of the Edgeworth expansion proposed by [3] lead to the empirical Edgeworth expansion (EEE), which approximates the distribution function of the Studentized statistic up to the error $o(n^{-1/2})$ in probability. For the particular case of L -statistics, variants of EEEs were considered in [12], and, assuming that values of an auxiliary variable are available for all units of the population, the calibration technique [16] was applied to estimate the parameters of Edgeworth expansion by [14, 28]. In particular, well-correlated auxiliary information, often accessible for finite populations, improves EEEs based on the sample \mathbb{X} only. In the case of independent and identically distributed (i.i.d.) observations, EEEs of Studentized L -statistics were derived by [19, 26, 29].

We focus ourselves on alternative methods to approximate the distribution function (1.2). We present two new approximations in Section 2. These methods complement the recent works of the authors on the use of the auxiliary information, which is specified as follows: denote by z the auxiliary variable with known real values $\mathcal{Z} = \{z_1, \dots, z_N\}$ in the population \mathcal{U} . Let $\mathbb{Z} = \{Z_1, \dots, Z_n\}$ be the corresponding values of the sample units, and $Z_{1:n} \leq \dots \leq Z_{n:n}$ let be order statistics of \mathbb{Z} . [14] introduced the approximation

$$F_{S_z}(y) = P\{\hat{\sigma}_J^{-1}(\mathbb{Z})(L_n(\mathbb{Z}) - EL_n(\mathbb{Z})) \leq y\} \quad (1.4)$$

to (1.2). We call it synthetic because it is based on the auxiliary data only. The numerical study of [14] showed that (1.4) is very efficient if the shapes of distributions of the variables x and z are similar. However, this naive approximation yields misleading results in practical situations.

The errors of nonparametric bootstrap approximations to distribution functions of statistics are of similar order as for EEEs. For samples drawn without replacement, this fact is known at least for statistics that are smooth functions of multivariate sample means [8] and U -statistics [5]. The latter class of estimators includes the Gini mean difference statistic, and in this case the error of the bootstrap approximation is $o(n^{-1/2})$ in probability. The simulation study of [12] also suggests that the accuracy of the nonparametric bootstrap approximation to (1.2) is similar to that of EEEs. For i.i.d. observations, the quality of nonparametric bootstrap approximations of U -statistics was investigated by [23], and distributions of trimmed means were approximated by [20, 21]. In Section 2.1, we use an auxiliary information and employ synthetic approximation (1.4) to construct a new calibrated nonparametric approximation to (1.2) based on the finite population bootstrap variant proposed by [8]. This calibration appears to be related to the empirical likelihood estimation for finite populations by [9]. The use of auxiliary information in the construction of bootstrap estimators is not widely studied in the literature, but, for instance, the paper of [2] presents several algorithms that incorporate the auxiliary data into bootstrap procedures. In the i.i.d. setting, a bootstrapping with auxiliary information was proposed by [31], which is similar to the conditional bootstrap methods introduced by [25].

Saddlepoint approximations to distribution functions of statistics are known as very accurate for small sample sizes and, in particular, on the tails of distributions. In Section 2.2, we present empirical saddlepoint

approximations to (1.2) applied to the distribution function of a suitably Studentized linear part of the L -statistic. More specifically, the linear part of (1.1) is taken from Hoeffding’s decomposition results of [6, 11], and then the “true” saddlepoint approximation, constructed for the distribution function of the Studentized sample mean by [15], is applied directly. This methodology is similar to that outlined by [17] and applied by [7] for standardized L -statistics in the case of i.i.d. observations. Next, we derive two empirical saddlepoint approximations based on the bootstrap, without and with the auxiliary information. In the traditional statistics, [24] applied saddlepoint approximations to the distributions of Studentized trimmed means.

The accuracy of approximations to the distribution functions of the Studentized L -statistics is investigated by simulation experiments in Section 3. The conclusions, stated in Section 4, are based on these numerical comparisons.

2 Approximations to the distribution

2.1 Calibrated nonparametric bootstrap

We use the finite population bootstrap scheme from [8]. Write $N = mn + l$, where $0 \leq l < n$. Given the sample \mathbb{X} , the empirical set (bootstrap population) \mathcal{X} of size N is formed by taking m copies of \mathbb{X} and, in case $l > 0$, adding the remaining l values, which are the simple random sample $\mathbb{Y} = \{Y_1, \dots, Y_l\}$ drawn without replacement from the set \mathbb{X} . Next, the simple random sample $\tilde{\mathbb{X}} = \{\tilde{X}_1, \dots, \tilde{X}_n\}$ is drawn without replacement from \mathcal{X} , and $L_n(\tilde{\mathbb{X}})$ and $\hat{\sigma}_J^2(\tilde{\mathbb{X}})$ are the bootstrap estimators of statistics (1.1) and (1.3), respectively. The nonparametric bootstrap approximation to (1.2) is

$$F_{\text{SB}}(y) = P\{\hat{\sigma}_J^{-1}(\tilde{\mathbb{X}})(L_n(\tilde{\mathbb{X}}) - E(L_n(\tilde{\mathbb{X}}) \mid \mathbb{X}, \mathbb{Y})) \leq y \mid \mathbb{X}\}, \tag{2.1}$$

which averages over $\binom{n}{l}$ possible bootstrap populations. Here the quantity $\mu(\mathcal{X}) = E(L_n(\tilde{\mathbb{X}}) \mid \mathbb{X}, \mathbb{Y})$ is the expectation of $L_n(\tilde{\mathbb{X}})$ under the fixed empirical population \mathcal{X} . For the original set \mathcal{X} , assuming that $x_1 \leq \dots \leq x_N$ without any loss of generality, it is expressed as follows [10, Appendix A]:

$$\mu(\mathcal{X}) = EL_n(\mathbb{X}) = E(L_n(\mathbb{X}) \mid \mathcal{X}) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^n c_{j,n} \mathcal{H}_{N-1, n-1, i-1}(j-1)x_i$$

with the hypergeometric probabilities

$$\mathcal{H}_{N,n,i}(j) = \binom{i}{j} \binom{N-i}{n-j} / \binom{N}{n}$$

having the support $\max\{0, n + i - N\} \leq j \leq \min\{n, i\}$. To evaluate (2.1), we apply the following Monte Carlo approximation. First, we construct independently B bootstrap populations $\tilde{\mathcal{X}}^{(b)}$, $1 \leq b \leq B$. Second, for each b , we draw independently R simple random samples $\tilde{\mathbb{X}}^{(b,r)} = \{\tilde{X}_1^{(b,r)}, \dots, \tilde{X}_n^{(b,r)}\}$, $1 \leq r \leq R$, without replacement from $\tilde{\mathcal{X}}^{(b)}$. Then, as proposed by [8],

$$\tilde{F}_{\text{SB}}(y) = \frac{1}{BR} \sum_{b=1}^B \sum_{r=1}^R \mathbb{I}\{\hat{\sigma}_J^{-1}(\tilde{\mathbb{X}}^{(b,r)})(L_n(\tilde{\mathbb{X}}^{(b,r)}) - \mu(\tilde{\mathcal{X}}^{(b)})) \leq y\} \tag{2.2}$$

is the formula to calculate (2.1) in practice. Here $\mathbb{I}\{\cdot\}$ is the indicator function.

Using representation (2.2), we define the calibrated nonparametric bootstrap approximation

$$F_{\text{SB}w}(y) = \frac{1}{BR} \sum_{b=1}^B \sum_{r=1}^R w_{br} \mathbb{I}\{\hat{\sigma}_J^{-1}(\tilde{\mathbb{X}}^{(b,r)})(L_n(\tilde{\mathbb{X}}^{(b,r)}) - \mu(\tilde{\mathcal{X}}^{(b)})) \leq y\} \tag{2.3}$$

to (1.2), where the weights $\mathbf{W} = (w_{br}) \in \mathbb{R}^{B \times R}$ minimize the function

$$d(\mathbf{W}) = \frac{1}{BR} \sum_{b=1}^B \sum_{r=1}^R (w_{br} - 1)^2 \quad (2.4)$$

and satisfy the calibration equations

$$\frac{1}{BR} \sum_{b=1}^B \sum_{r=1}^R w_{br} \mathbb{I}\{Z_{b,r} \leq y_i\} = F_{S_Z}(y_i), \quad 1 \leq i \leq T, \quad (2.5)$$

where

$$Z_{b,r} = \hat{\sigma}_J^{-1}(\tilde{Z}^{(b,r)})(L_n(\tilde{Z}^{(b,r)}) - \mu(\tilde{Z}^{(b)})), \quad 1 \leq b \leq B, 1 \leq r \leq R. \quad (2.6)$$

Here the sets $\tilde{Z}^{(b)}$ and $\tilde{Z}^{(b,r)}$, constructed from \mathbb{Z} , represent exactly the same sample units as the sets $\tilde{\mathcal{X}}^{(b)}$ and $\tilde{\mathbb{X}}^{(b,r)}$ selected from the given \mathbb{X} . The auxiliary function (1.4) is evaluated by drawing a large number of independent samples from \mathcal{Z} . The arbitrarily chosen points $y_1 < \dots < y_{T-1}$ are, for example, uniformly spaced quantiles of the distribution function of values (2.6), and the choice of the point

$$y_T = \max \left\{ \max_{1 \leq b \leq B, 1 \leq r \leq R} Z_{b,r}, \max_{s \in [0,1]} F_{S_Z}^{-1}(s) \right\}$$

means that the last equation in (2.5) is the requirement that the average of the calibrated weights is equal to 1. For simplicity, we can set $y_T = 10^3$. Explicit expressions of the weights are presented in the following proposition.

Proposition 1. *Let $y_1 < \dots < y_T$. Assume that $m_1 > 0$ and $m_i > m_{i-1}$ for $2 \leq i \leq T$, where m_i is the number of values in the set (2.6) that are smaller than or equal to the value y_i . Then the weights \mathbf{W} minimizing function (2.4) and satisfying calibration equations (2.5) are unique and expressed by*

$$w_{br} = 1 + \frac{1}{2} \sum_{j=1}^T \lambda_j \mathbb{I}\{Z_{b,r} \leq y_j\}, \quad 1 \leq b \leq B, 1 \leq r \leq R, \quad (2.7)$$

where $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_T)^\top = \mathbf{A}^{-1} \mathbf{b}$ with $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{T \times T}$ and $\mathbf{b} = (b_1, \dots, b_T)^\top$ given by

$$a_{ij} = \frac{1}{2BR} \sum_{b=1}^B \sum_{r=1}^R \mathbb{I}\{Z_{b,r} \leq y_i\} \mathbb{I}\{Z_{b,r} \leq y_j\} \quad \text{and} \quad b_i = F_{S_Z}(y_i) - \frac{1}{BR} \sum_{b=1}^B \sum_{r=1}^R \mathbb{I}\{Z_{b,r} \leq y_i\},$$

respectively.

Proof. Consider the Lagrange function

$$\mathcal{L} = \mathcal{L}(\mathbf{W}, \boldsymbol{\lambda}) = d(\mathbf{W}) - \sum_{j=1}^T \lambda_j \left(\frac{1}{BR} \sum_{b=1}^B \sum_{r=1}^R w_{br} \mathbb{I}\{Z_{b,r} \leq y_j\} - F_{S_Z}(y_j) \right).$$

Equating the partial derivatives $\partial \mathcal{L} / \partial w_{br}$, $1 \leq b \leq B$, $1 \leq r \leq R$, to zero, we derive expressions (2.7). Next, insert these expressions into calibration equations (2.5) and obtain the system of linear equations $\mathbf{A} \boldsymbol{\lambda} = \mathbf{b}$.

Write the matrix \mathbf{A} as

$$\mathbf{A} = \frac{1}{2BR} \begin{pmatrix} m_1 & m_1 & \cdots & m_1 \\ m_1 & m_2 & \cdots & m_2 \\ \vdots & \vdots & \ddots & \vdots \\ m_1 & m_2 & \cdots & m_T \end{pmatrix}.$$

By the properties of determinants,

$$\begin{aligned} \det(\mathbf{A}) &= \frac{1}{(2BR)^T} \det \begin{pmatrix} m_1 & m_1 & \cdots & m_1 \\ 0 & m_2 - m_1 & \cdots & m_2 - m_1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & m_T - m_{T-1} \end{pmatrix} \\ &= \frac{1}{(2BR)^T} m_1(m_2 - m_1) \cdots (m_T - m_{T-1}) > 0. \end{aligned}$$

Therefore the solution λ of the system of equations is unique, which leads to the complete knowledge of calibrated weights (2.7). \square

Remark 1. Our experiments show that the number T of points needs not be large to get an optimal calibrated bootstrap approximation. In the simulation study of Section 3, for the uniformly spaced quantiles $y_1 < \cdots < y_{T-1}$ of the distribution of (2.6), the choice $T = 10^2$ is better than $T = 10$, but $T = 10^3$ gives no significant further improvement of (2.3). Moreover, if T is large, then some of the quantiles can coincide. Then $\det(\mathbf{A}) = 0$, and there is no unique solution (2.7).

Remark 2. If we replace the minimization of distance (2.4) by the maximization of function

$$g(\mathbf{W}) = \sum_{b=1}^B \sum_{r=1}^R \log(w_{br}),$$

then the calibrated estimation becomes a finite population version of the empirical likelihood (EL) method of [27]. In the case of i.i.d. observations, see also the review by [22]. For the simple random sampling without replacement, an EL estimation that uses auxiliary information is discussed by [9]. Indeed, our version of EL is an extension of the latter methodology to the estimation of multivariate means (proportions). There are evidences that similar calibration procedures and EL estimators are asymptotically equivalent as the sample size tends to infinity [30]. In our situation, the “sample size” BR is large, and the numerical tests show that both resulting approximations to (1.2) are almost identical. The disadvantage of the EL method is that there is no explicit expression of the Lagrange multipliers that define the weights and are the solution of the system of T nonlinear equations.

2.2 Saddlepoint approximations

To employ saddlepoint techniques to approximate the distribution function (1.2), the idea of [18], later noted by [17], is linearization of complex statistic (1.1), and then applying the saddlepoint approximation to the distribution function of the linear part. For the simple random samples drawn without replacement, the general symmetric statistics are linearized using Hoeffding’s decomposition by [6],

$$L - EL = H + R, \quad \text{where } H = H_n(\mathbb{X}) = \frac{1}{n} \sum_{j=1}^n h(X_j) \tag{2.8}$$

is a linear statistic with influence function h , and $R = R_n(\mathbb{X})$ is a remainder term. Here the random variables $h(X_1), \dots, h(X_n)$ are identically distributed with $P\{h(X_1) = h(x_k)\} = N^{-1}$, $1 \leq k \leq N$, where, letting $x_1 \leq \dots \leq x_N$, the explicit expressions

$$h(x_k) = h(k; \mathcal{X}) = - \sum_{i=1}^{N-1} \left(\mathbb{I}\{i \geq k\} - \frac{i}{N} \right) d_i(x_{i+1} - x_i) \tag{2.9}$$

with

$$d_i = d_{i,N,n} = \sum_{j=1}^n c_{j,n} \mathcal{H}_{N-2,n-1,i-1}(j-1)$$

are available for the particular case of L -statistics [11]. In decomposition (2.8), the components H and R are centered and uncorrelated. Furthermore, $R = O(n^{-1/2})$ in probability for many commonly used statistics [6] and for L -statistics with sufficiently smooth weight functions J ; see [13].

The jackknife estimator of the variance $\sigma_H^2 = \text{Var } H$ of the linear statistic reduces from (1.3) to

$$\hat{\sigma}_{HJ}^2 = \hat{\sigma}_{HJ}^2(\mathbb{X}) = \left(1 - \frac{n}{N}\right) \frac{S^2}{n}, \quad \text{where } S^2 = \frac{1}{n-1} \sum_{j=1}^n (h(X_j) - H)^2$$

is the variance of the transformed sample. Thus, we approximate distribution function (1.2) by

$$\tilde{F}_S(y) = P\{\hat{\sigma}_{HJ}^{-1} H \leq y\} \tag{2.10}$$

and we next apply the saddlepoint approximation results of [15] to the latter distribution function. To this aim, introduce the random variables $A_j = h(X_j)/\sigma_1$, $1 \leq j \leq n$, where $\sigma_1^2 = N^{-1} \sum_{k=1}^N h^2(x_k)$. Then $\mathbb{A} = \{A_1, \dots, A_n\}$ is the simple random sample without replacement from the set $\mathcal{A} = \{a_1, \dots, a_N\}$ with $a_k = h(x_k)/\sigma_1$, $1 \leq k \leq N$. The random variable $\hat{\sigma}_{HJ}^{-1} H$ in (2.10) coincides with the Studentized mean of the sample \mathbb{A} considered by [15], and their assumptions $\sum_{k=1}^N a_k = 0$ and $N^{-1} \sum_{k=1}^N a_k^2 = 1$ hold as well.

Write shortly $p = n/N$ and $q = 1 - p$, and let $\mathbf{t} = (t_0, t_1, t_2)^\top \in \mathbb{R}^3$ and $\mathbf{u} = (u_0, u_1, u_2)^\top \in \mathbb{R}^3$. Consider the function

$$K(\mathbf{t}) = -p(t_0 + t_2) + \frac{1}{N} \sum_{k=1}^N \log(q + p \exp(t_0 + t_1 a_k + t_2 a_k^2)) \tag{2.11}$$

and define $\mathbf{t}(\mathbf{u})$ as the solution of the equation system

$$K'(\mathbf{t}) = \mathbf{u}, \tag{2.12}$$

which is solved numerically in practice. Then we introduce the functions

$$\Lambda(\mathbf{u}) = \mathbf{t}^\top(\mathbf{u})\mathbf{u} - K(\mathbf{t}(\mathbf{u})) \quad \text{and} \quad \Delta(\mathbf{u}) = \det(K''(\mathbf{t}(\mathbf{u}))).$$

For the fixed point $y \in \mathbb{R}$, we define $u_0 = 0$ and

$$u_1 = u_1(y, u_2) = y \left(\frac{u_2 + p}{c_n^2 + y^2/p} \right)^{1/2} \quad \text{with } c_n = \left(\frac{n-1}{pq} \right)^{1/2}. \tag{2.13}$$

Then we solve the equation

$$\frac{\partial \Lambda(0, u_1(y, u_2), u_2)}{\partial u_2} = 0, \tag{2.14}$$

that is, find numerically the number $u_2(y)$ minimizing the function $\Lambda(0, u_1(y, u_2), u_2)$, and then we denote $u_1(y) = u_1(y, u_2(y))$ according to (2.13). The calculation of the minimizers $u_2(\tilde{y})$ at the points \tilde{y} close to y is fast because the solutions of (2.14) and (2.12) vary slowly. Next, we evaluate the functions

$$G(y) = \frac{(pq)^{1/2}(\partial^2 \Lambda(0, u_1(y, u_2), u_2)/\partial u_2^2|_{(u_1(y), u_2(y))})^{-1/2}}{|\partial v(u_1, u_2)/\partial u_1|_{(u_1(y), u_2(y))}|\Delta^{1/2}(0, u_1(y), u_2(y))},$$

where $v(u_1, u_2) = c_n u_1/(u_2 + p - u_1^2/p)^{1/2}$ and

$$D(y) = \Lambda(0, u_1(y), u_2(y)).$$

Then the saddlepoint approximation to distribution function (2.10) of the Studentized linear part of (1.1) is [15]

$$\tilde{F}_{SS}(y) = \Phi\left(N^{1/2}\left(w(y) - N^{-1}\frac{\log(w(y)G(y)/D'(y))}{w(y)}\right)\right), \tag{2.15}$$

where $w(y) = \text{sgn}(y)(2D(y))^{1/2}$. Assuming that the variance of the remainder R in (2.8) is negligible, we use approximation (2.15) to distribution (1.2) of interest as well.

However, the “true” saddlepoint approximation (2.15) is useless in practice because function (2.11) depends on the set \mathcal{A} of unknown characteristics. We apply the bootstrap of [8] to estimate these $a_k, 1 \leq k \leq N$. A similar application of the same bootstrap to saddlepoint approximations is done in [1]. The explicit estimators of (2.9) for $1 \leq k \leq N$ are [12]

$$\hat{h}_B(k; \mathbb{X}) = -\sum_{j=1}^{n-1} \sum_{i=mj}^{mj+l} \left(\mathbb{I}\{i \geq k\} - \frac{i}{N}\right) d_i \mathcal{H}_{n,l,j}(i - mj)(X_{j+1:n} - X_{j:n}). \tag{2.16}$$

Then the bootstrap estimators $\hat{a}_B(k) = \hat{h}_B(k; \mathbb{X})/\hat{\sigma}_{1B}^2, 1 \leq k \leq N$, where $\hat{\sigma}_{1B}^2 = N^{-1} \sum_{k=1}^N \hat{h}_B^2(k; \mathbb{X})$, are plugged into (2.11), and the empirical saddlepoint approximation $F_{SSB}(y)$ to (1.2) is obtained following the formulas used to calculate (2.15).

If a well-correlated auxiliary information \mathcal{Z} is available, then the calibration of bootstrap estimators (2.16) by [28] can lead to more efficient estimators of the characteristics $a_k, 1 \leq k \leq N$. The calibrated bootstrap estimators of (2.9) for $1 \leq k \leq N$ are

$$\tilde{h}_{Bw}(k; \mathbb{X}, \mathcal{Z}) = -\sum_{j=1}^{n-1} w_j(k) \sum_{i=mj}^{mj+l} \left(\mathbb{I}\{i \geq k\} - \frac{i}{N}\right) d_i \mathcal{H}_{n,l,j}(i - mj)(X_{j+1:n} - X_{j:n}),$$

where the calibration weights are

$$w_j(k) = 1 + (h(k; \mathcal{Z}) - \hat{h}_B(k; \mathbb{Z})) \left(\sum_{t=1}^{n-1} r_t^2(k)\right)^{-1} r_j(k), \quad 1 \leq j \leq n-1,$$

with

$$r_j(k) = -\sum_{i=mj}^{mj+l} \left(\mathbb{I}\{i \geq k\} - \frac{i}{N}\right) d_i \mathcal{H}_{n,l,j}(i - mj)(Z_{j+1:n} - Z_{j:n}).$$

We define the calibrated estimators of a_k by $\hat{a}_{Bw}(k) = \hat{h}_{Bw}(k; \mathbb{X}, \mathcal{Z}) / \hat{\sigma}_{1Bw}$, $1 \leq k \leq N$, where

$$\hat{h}_{Bw}(k; \mathbb{X}, \mathcal{Z}) = \tilde{h}_{Bw}(k; \mathbb{X}, \mathcal{Z}) - \frac{1}{N} \sum_{k=1}^N \tilde{h}_{Bw}(k; \mathbb{X}, \mathcal{Z}) \quad \text{and} \quad \hat{\sigma}_{1Bw}^2 = \frac{1}{N} \sum_{k=1}^N \hat{h}_{Bw}^2(k; \mathbb{X}, \mathcal{Z}),$$

and so we get another empirical saddlepoint approximation to (1.2), which we denote by $F_{SSBw}(y)$.

3 Simulation study

We compare the calibrated nonparametric bootstrap F_{SBw} and saddlepoint F_{SSBw} approximations to the distribution F_S with approximations F_{SB} and F_{SSB} based on the data \mathbb{X} only. The comparison also includes the normal approximation Φ , the synthetic approximation F_{S_z} , the “true” saddlepoint approximation \tilde{F}_{SS} of a theoretical interest, and the empirical Edgeworth expansion G_{SBw} with calibrated bootstrap estimators of parameters by [28], which appears to be the most robust approximation in the cases simulated by [14]. Moreover, we use the populations \mathcal{U}_1 , \mathcal{U}_2 , and \mathcal{U}_3 of size $N = 120$ and the L -statistics taken from the latter simulation study, where various empirical Edgeworth approximations were compared.

The values of the auxiliary variable z of the first population \mathcal{U}_1 are generated from the Fisher distribution $\mathcal{F}(5, 4)$, and then the values of the study variable x are obtained by the relationship $x_i = 2 + z_i + 0.7\sqrt{z_i}\varepsilon_i$, where independent errors ε_i are from the normal distribution $\mathcal{N}(0, 1)$. The Pearson correlation coefficient between the fixed sets \mathcal{Z} and \mathcal{X} is equal approximately to 0.92. For the second population \mathcal{U}_2 , the values of the variables z and x are simulated, respectively, according to the marginal distributions $\mathcal{N}(600, 150)$ and $\mathcal{F}(5, 4)$ and by applying the bivariate Student’s t copula. The resulting coefficient of linear correlation is close to 0.83. The elements of the third population \mathcal{U}_3 are business enterprises. The variable z denotes an annual turnover derived from administrative Value Added Tax data, and the variable x is an annual survey turnover. The correlation coefficient is 0.81. Figure 1 presents the distributions of both variables and their relationship.

Six different scenarios for the simulation study are obtained by combining these populations with two L -statistics. The first one is the Gini mean difference defined using the weights

$$c_{j,n} = \frac{n+1}{n-1} J\left(\frac{j}{n+1}\right), \quad 1 \leq j \leq n,$$

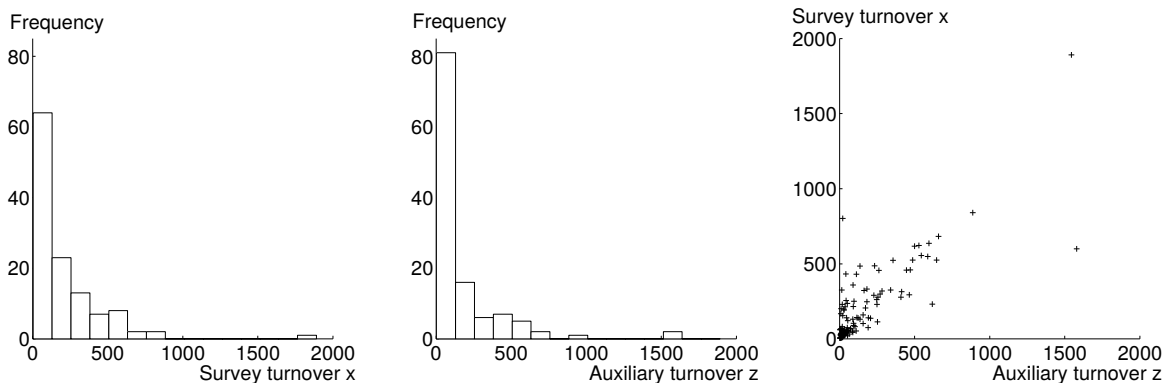


Figure 1. The annual survey turnover variable x and the annual administrative turnover variable z in the population \mathcal{U}_3 of business enterprises.

with smooth function $J(s) = 4s - 2$. Given the fixed numbers $0 \leq a < b \leq 1$, the trimmed mean

$$M_{n;a;b}(\mathbb{X}) = \frac{1}{[bn] - [an]} \sum_{j=[an]+1}^{[bn]} X_{j:n}, \tag{3.1}$$

where $[\cdot]$ is the greatest integer function, is represented asymptotically by the nonsmooth function $J(s) = (b - a)^{-1} \mathbb{I}\{a < s < b\}$. The second statistic is the trimmed mean (3.1) with values $a = 0$ and $b = 0.95$ of trimming proportions. The samples are of size $n = 40$ in all cases.

We compare all the approximations to the distribution F_S of interest by taking their s -quantiles with $s = 0.01, 0.05, 0.10, 0.90, 0.95, 0.99$. For each empirical quantile of the sample-based approximations, we evaluate its expectation (denoted by the operator E_m) and the root mean square error (R_m) using 10^3 samples, drawn independently from a particular population. For the quantiles of the population-based approximations Φ , F_{S_z} , and F_{SS} , we also calculate the characteristics R_m , but these constitute of the bias component only. Tables 1–6 present the results.

Almost all approximations improve the normal approximation Φ . The synthetic approximation F_{S_z} is the best one under the population \mathcal{U}_1 , but its accuracy is similar to that of Φ in the population \mathcal{U}_2 . Moreover, in the population \mathcal{U}_3 , the results of F_{S_z} are perfect for the Gini mean difference statistic, but they are bad for the trimmed mean. The calibrated bootstrap approximation F_{SBw} improves the bootstrap F_{SB} , and this improvement is significant in the population \mathcal{U}_1 . The calibrated saddlepoint approximation F_{SSBw} is better compared to the saddlepoint F_{SSB} based on the bootstrap, and the root mean square errors of F_{SSBw} are much closer to that of the “true” saddlepoint approximation F_{SS} . Comparing the calibrated nonparametric approximations F_{SBw} and F_{SSBw} with the calibrated (parametric) Edgeworth expansion G_{SBw} , the main difference between them is that the latter approximation tends to have smaller variances but larger biases. In many cases, the biases of G_{SBw} are particularly large for quantiles 0.01 and 0.99, where the root mean square errors of F_{SBw} and F_{SSBw} are thus much smaller. In turn, for far quantiles of F_S , the approximation F_{SBw} appears to be more accurate than or similar to F_{SSBw} . We conclude that the calibrated bootstrap F_{SBw} is the best approximation among the sample-based approximations F_{SB} , F_{SSBw} , F_{SSB} , and G_{SBw} , and it is robust compared to the synthetic approximation F_{S_z} .

Table 1. Approximations to $F_S^{-1}(s)$ of the Gini mean difference under the population \mathcal{U}_1

s	0.01	0.05	0.10	0.90	0.95	0.99
F_S^{-1}	-6.454	-4.401	-3.457	0.982	1.181	1.494
Φ^{-1}	-2.326	-1.645	-1.282	1.282	1.645	2.326
$F_{S_z}^{-1}$	-7.258	-5.048	-3.822	1.004	1.218	1.580
\tilde{F}_{SS}^{-1}	-8.092	-5.384	-4.200	0.953	1.145	1.447
$E_m F_{SBw}^{-1}$	-6.601	-4.422	-3.344	1.020	1.252	1.664
$E_m F_{SB}^{-1}$	-6.551	-4.191	-2.709	1.095	1.352	1.839
$E_m \tilde{F}_{SSBw}^{-1}$	-8.512	-5.494	-4.150	0.992	1.207	1.552
$E_m F_{SSB}^{-1}$	-8.005	-4.647	-3.342	1.038	1.278	1.680
$E_m G_{SBw}^{-1}$	-3.053	-2.288	-1.821	0.967	1.151	1.355
$R_m \Phi^{-1}$	4.127	2.756	2.175	0.299	0.464	0.832
$R_m F_{S_z}^{-1}$	0.805	0.647	0.365	0.021	0.036	0.086
$R_m \tilde{F}_{SS}^{-1}$	1.638	0.983	0.744	0.030	0.036	0.047
$R_m F_{SBw}^{-1}$	1.998	1.412	1.151	0.057	0.106	0.218
$R_m F_{SB}^{-1}$	3.026	2.222	1.604	0.146	0.226	0.417
$R_m \tilde{F}_{SSBw}^{-1}$	3.177	2.091	1.692	0.031	0.052	0.100
$R_m F_{SSB}^{-1}$	3.991	2.286	1.998	0.079	0.128	0.238
$R_m G_{SBw}^{-1}$	3.402	2.115	1.638	0.050	0.073	0.175

Table 2. Approximations to $F_S^{-1}(s)$ of the trimmed mean under the population \mathcal{U}_1

s	0.01	0.05	0.10	0.90	0.95	0.99
F_S^{-1}	-3.270	-2.100	-1.558	1.155	1.451	2.004
Φ^{-1}	-2.326	-1.645	-1.282	1.282	1.645	2.326
$F_{S_z}^{-1}$	-3.206	-2.103	-1.599	1.104	1.388	1.900
\tilde{F}_{SS}^{-1}	-2.964	-1.928	-1.448	1.192	1.511	2.095
$E_m F_{SBw}^{-1}$	-3.211	-2.098	-1.577	1.117	1.411	1.942
$E_m F_{SB}^{-1}$	-3.395	-2.231	-1.682	1.101	1.391	1.924
$E_m \tilde{F}_{SSBw}^{-1}$	-2.996	-1.943	-1.457	1.189	1.507	2.087
$E_m F_{SSB}^{-1}$	-3.151	-2.001	-1.489	1.180	1.492	2.062
$E_m G_{SBw}^{-1}$	-2.782	-1.988	-1.540	1.078	1.338	1.706
$R_m \Phi^{-1}$	0.944	0.455	0.277	0.127	0.194	0.322
$R_m F_{S_z}^{-1}$	0.064	0.003	0.041	0.051	0.063	0.104
$R_m \tilde{F}_{SS}^{-1}$	0.306	0.172	0.111	0.037	0.060	0.091
$R_m F_{SBw}^{-1}$	0.255	0.131	0.090	0.050	0.060	0.104
$R_m F_{SB}^{-1}$	0.495	0.337	0.269	0.095	0.121	0.184
$R_m \tilde{F}_{SSBw}^{-1}$	0.301	0.166	0.106	0.037	0.060	0.093
$R_m F_{SSB}^{-1}$	0.435	0.189	0.114	0.042	0.067	0.113
$R_m G_{SBw}^{-1}$	0.493	0.138	0.078	0.092	0.135	0.324

Table 3. Approximations to $F_S^{-1}(s)$ of the Gini mean difference under the population \mathcal{U}_2

s	0.01	0.05	0.10	0.90	0.95	0.99
F_S^{-1}	-4.304	-2.366	-1.621	1.164	1.461	2.016
Φ^{-1}	-2.326	-1.645	-1.282	1.282	1.645	2.326
F_{Sz}^{-1}	-2.654	-1.742	-1.317	1.259	1.610	2.272
\tilde{F}_{SS}^{-1}	-5.347	-2.686	-1.836	1.079	1.342	1.792
$E_m F_{SBw}^{-1}$	-4.238	-2.385	-1.632	1.180	1.486	2.058
$E_m F_{SB}^{-1}$	-4.244	-2.393	-1.634	1.185	1.492	2.066
$E_m F_{SSBw}^{-1}$	-5.146	-2.663	-1.827	1.085	1.349	1.805
$E_m F_{SSB}^{-1}$	-5.267	-2.741	-1.868	1.085	1.349	1.804
$E_m G_{SBw}^{-1}$	-2.735	-1.917	-1.456	1.157	1.427	1.829
$R_m \Phi^{-1}$	1.977	0.721	0.340	0.118	0.184	0.310
$R_m F_{Sz}^{-1}$	1.650	0.624	0.304	0.095	0.149	0.256
$R_m \tilde{F}_{SS}^{-1}$	1.043	0.320	0.214	0.084	0.119	0.224
$R_m F_{SBw}^{-1}$	1.005	0.546	0.237	0.053	0.082	0.158
$R_m F_{SB}^{-1}$	1.074	0.592	0.259	0.062	0.092	0.171
$R_m F_{SSBw}^{-1}$	1.212	0.416	0.250	0.081	0.116	0.218
$R_m F_{SSB}^{-1}$	1.634	0.671	0.394	0.084	0.120	0.225
$R_m G_{SBw}^{-1}$	1.571	0.455	0.176	0.039	0.067	0.219

Table 4. Approximations to $F_S^{-1}(s)$ of the trimmed mean under the population \mathcal{U}_2

s	0.01	0.05	0.10	0.90	0.95	0.99
F_S^{-1}	-3.453	-2.186	-1.637	1.122	1.412	1.945
Φ^{-1}	-2.326	-1.645	-1.282	1.282	1.645	2.326
F_{Sz}^{-1}	-2.488	-1.716	-1.321	1.288	1.668	2.411
\tilde{F}_{SS}^{-1}	-3.078	-1.973	-1.474	1.178	1.490	2.057
$E_m F_{SBw}^{-1}$	-3.443	-2.212	-1.642	1.135	1.436	1.978
$E_m F_{SB}^{-1}$	-3.447	-2.211	-1.640	1.136	1.436	1.978
$E_m F_{SSBw}^{-1}$	-3.079	-1.971	-1.471	1.179	1.493	2.063
$E_m F_{SSB}^{-1}$	-3.072	-1.968	-1.469	1.181	1.495	2.068
$E_m G_{SBw}^{-1}$	-2.770	-1.973	-1.524	1.091	1.354	1.730
$R_m \Phi^{-1}$	1.127	0.541	0.356	0.160	0.233	0.381
$R_m F_{Sz}^{-1}$	0.965	0.469	0.316	0.166	0.256	0.465
$R_m \tilde{F}_{SS}^{-1}$	0.376	0.212	0.164	0.056	0.078	0.112
$R_m F_{SBw}^{-1}$	0.343	0.202	0.137	0.036	0.052	0.083
$R_m F_{SB}^{-1}$	0.385	0.222	0.149	0.039	0.056	0.090
$R_m F_{SSBw}^{-1}$	0.409	0.225	0.170	0.060	0.085	0.128
$R_m F_{SSB}^{-1}$	0.429	0.232	0.174	0.062	0.089	0.136
$R_m G_{SBw}^{-1}$	0.686	0.222	0.127	0.049	0.080	0.238

Table 5. Approximations to $F_S^{-1}(s)$ of the Gini mean difference under the population \mathcal{U}_3

s	0.01	0.05	0.10	0.90	0.95	0.99
F_S^{-1}	-4.243	-2.853	-2.231	1.012	1.202	1.514
Φ^{-1}	-2.326	-1.645	-1.282	1.282	1.645	2.326
F_{Sz}^{-1}	-4.726	-3.099	-2.381	1.039	1.276	1.668
\tilde{F}_{SS}^{-1}	-5.627	-3.542	-2.758	0.952	1.143	1.454
$E_m F_{SBw}^{-1}$	-4.524	-3.008	-2.303	1.121	1.396	1.918
$E_m F_{SB}^{-1}$	-4.361	-2.833	-2.123	1.181	1.480	2.060
$E_m F_{SSBw}^{-1}$	-5.826	-3.716	-2.872	1.037	1.275	1.674
$E_m F_{SSB}^{-1}$	-5.668	-3.472	-2.675	1.073	1.330	1.771
$E_m G_{SBw}^{-1}$	-2.904	-2.126	-1.665	1.053	1.275	1.571
$R_m \Phi^{-1}$	1.917	1.208	0.949	0.269	0.443	0.812
$R_m F_{Sz}^{-1}$	0.483	0.246	0.150	0.027	0.074	0.154
$R_m \tilde{F}_{SS}^{-1}$	1.384	0.689	0.527	0.060	0.059	0.060
$R_m F_{SBw}^{-1}$	1.567	1.234	1.054	0.159	0.260	0.498
$R_m F_{SB}^{-1}$	1.861	1.409	1.092	0.208	0.338	0.637
$R_m F_{SSBw}^{-1}$	2.325	1.556	1.353	0.065	0.118	0.230
$R_m F_{SSB}^{-1}$	2.857	1.851	1.623	0.085	0.158	0.305
$R_m G_{SBw}^{-1}$	1.358	0.764	0.605	0.126	0.189	0.320

Table 6. Approximations to $F_S^{-1}(s)$ of the trimmed mean under the population \mathcal{U}_3 .

s	0.01	0.05	0.10	0.90	0.95	0.99
F_S^{-1}	-3.197	-2.012	-1.495	1.172	1.483	2.058
Φ^{-1}	-2.326	-1.645	-1.282	1.282	1.645	2.326
F_{Sz}^{-1}	-3.953	-2.241	-1.624	1.131	1.422	1.956
\tilde{F}_{SS}^{-1}	-2.901	-1.896	-1.428	1.198	1.519	2.117
$E_m F_{SBw}^{-1}$	-3.162	-1.999	-1.483	1.176	1.488	2.063
$E_m F_{SB}^{-1}$	-3.220	-2.059	-1.536	1.153	1.461	2.030
$E_m F_{SSBw}^{-1}$	-2.928	-1.907	-1.434	1.195	1.518	2.111
$E_m F_{SSB}^{-1}$	-2.946	-1.916	-1.440	1.194	1.515	2.106
$E_m G_{SBw}^{-1}$	-2.691	-1.895	-1.459	1.136	1.420	1.858
$R_m \Phi^{-1}$	0.871	0.367	0.213	0.110	0.162	0.269
$R_m F_{Sz}^{-1}$	0.755	0.229	0.129	0.041	0.062	0.102
$R_m \tilde{F}_{SS}^{-1}$	0.297	0.116	0.067	0.026	0.035	0.059
$R_m F_{SBw}^{-1}$	0.261	0.145	0.107	0.040	0.051	0.075
$R_m F_{SB}^{-1}$	0.313	0.162	0.108	0.048	0.065	0.098
$R_m F_{SSBw}^{-1}$	0.297	0.117	0.068	0.028	0.042	0.070
$R_m F_{SSB}^{-1}$	0.299	0.119	0.069	0.029	0.044	0.075
$R_m G_{SBw}^{-1}$	0.509	0.126	0.053	0.045	0.076	0.216

4 Conclusions

Our simulations suggest that the constructed calibrated bootstrap and saddlepoint approximations to the distribution function of the Studentized L -statistic improve the respective approximations based only on the sample data if the study and auxiliary variables are well correlated. There are also numerical evidences that the new approximations adapt better to estimate extreme quantiles of the distribution of interest than the empirical Edgeworth expansion with calibrated parameters. Moreover, the latter approximation exhibits larger biases.

The calibrated bootstrap approximation can be interpreted as a nonlinear combination of the bootstrap and synthetic approximations, which adapts to the quality of auxiliary information; that is, for significantly different distributions of the study and auxiliary variables, as in the second population of the simulations, the efficiency of the combination is smaller, but it is still greater than that of the bootstrap approximation, whereas the synthetic approximation fails. The latter approximation is not reliable in the third population, where the real data contain outliers.

The calibrated saddlepoint method evaluates the “true” saddlepoint approximation quite well. However, all presented saddlepoint approximations do not take into account the nonlinear part of the L -statistic, whereas the other competitive approximations do. More accurate saddlepoint approximations can be constructed by using higher-order terms of the Hoeffding decomposition, similarly as in [17]. This is a question for future research. According to the simulations, the calibrated saddlepoint approximation is slightly worse than the calibrated bootstrap approximation.

References

1. K. Agho, W. Dai, and J. Robinson, Empirical saddlepoint approximations of the Studentized ratio and regression estimates for finite populations, *Stat. Probab. Lett.*, **71**(3):237–247, 2005.
2. A. Barbiero, G. Manzi, and F. Mecatti, Bootstrapping probability-proportional-to-size samples via calibrated empirical population, *J. Stat. Comput. Simulation*, **85**(3):608–620, 2015.
3. M. Bloznelis, Empirical Edgeworth expansion for finite population statistics. I, *Lith. Math. J.*, **41**(2):120–134, 2001.
4. M. Bloznelis, Edgeworth expansions for Studentized versions of symmetric finite population statistics, *Lith. Math. J.*, **43**(3):221–240, 2003.
5. M. Bloznelis, Bootstrap approximation to distributions of finite population U -statistics, *Acta Appl. Math.*, **96**:71–86, 2007.
6. M. Bloznelis and F. Götze, Orthogonal decomposition of finite population statistics and its applications to distributional asymptotics, *Ann. Stat.*, **29**(3):899–917, 2001.
7. M. Bogdan, Asymptotic distributions of linear combinations of order statistics, *Appl. Math.*, **22**(2):201–225, 1994.
8. J.G. Booth, R.W. Butler, and P. Hall, Bootstrap methods for finite populations, *J. Am. Stat. Assoc.*, **89**(428):1282–1289, 1994.
9. J. Chen and J. Qin, Empirical likelihood estimation for finite populations and the effective usage of auxiliary information, *Biometrika*, **80**(1):107–116, 1993.
10. A. Čiginas, *Approximations to Distributions of Linear Combinations of Order Statistics in Finite Populations*, PhD dissertation, Vilnius University, Lithuania, 2011.
11. A. Čiginas, An Edgeworth expansion for finite-population L -statistics, *Lith. Math. J.*, **52**(1):40–52, 2012.
12. A. Čiginas, Second-order approximations of finite population L -statistics, *Statistics*, **47**(5):954–965, 2013.
13. A. Čiginas, On the asymptotic normality of finite population L -statistics, *Stat. Pap.*, **55**(4):1047–1058, 2014.
14. A. Čiginas and D. Pumputis, Calibrated Edgeworth expansions of finite population L -statistics, *Math. Popul. Stud.*, 2019, available from: <https://doi.org/10.1080/08898480.2018.1553408>.
15. W. Dai and J. Robinson, Empirical saddlepoint approximations of the Studentized mean under simple random sampling, *Stat. Probab. Lett.*, **53**(3):331–337, 2001.
16. J.C. Deville and C.-E. Särndal, Calibration estimators in survey sampling, *J. Am. Stat. Assoc.*, **87**(418):376–382, 1992.

17. G.S. Easton and E. Ronchetti, General saddlepoint approximations with applications to L statistics, *J. Am. Stat. Assoc.*, **81**(394):420–430, 1986.
18. C. Field, Small sample asymptotic expansions for multivariate M -estimates, *Ann. Stat.*, **10**(3):672–689, 1982.
19. N.V. Gribkova and R. Helmers, The empirical Edgeworth expansion for a Studentized trimmed mean, *Math. Methods Stat.*, **15**(1):61–87, 2006.
20. N.V. Gribkova and R. Helmers, On the Edgeworth expansion and the M out of N bootstrap accuracy for a Studentized trimmed mean, *Math. Methods Stat.*, **16**(2):142–176, 2007.
21. P. Hall and A.R. Padmanabhan, On the bootstrap and the trimmed mean, *J. Multivariate Anal.*, **41**(1):132–153, 1992.
22. P. Hall and B. La Scala, Methodology and algorithms of empirical likelihood, *Int. Stat. Rev.*, **58**(2):109–127, 1990.
23. R. Helmers, On the Edgeworth expansion and the bootstrap approximation for a Studentized U -statistic, *Ann. Stat.*, **19**(1):470–484, 1991.
24. R. Helmers, B.-Y. Jing, G. Qin, and W. Zhou, Saddlepoint approximations to the trimmed mean, *Bernoulli*, **10**(3):465–501, 2004.
25. D. Hinkley and E. Schechtman, Conditional bootstrap methods in the mean-shift model, *Biometrika*, **74**(1):85–93, 1987.
26. Y. Maesono, An Edgeworth expansion and a normalizing transformation for L -statistics, *Bull. Inf. Cybern.*, **39**:25–43, 2007.
27. A. Owen, Empirical likelihood ratio confidence regions, *Ann. Stat.*, **18**(1):90–120, 1990.
28. D. Pumputis and A. Čiginas, Estimation of parameters of finite population L -statistics, *Nonlinear Anal. Model. Control*, **18**(3):327–343, 2013.
29. H. Putter and W.R. van Zwet, Empirical Edgeworth expansions for symmetric statistics, *Ann. Stat.*, **26**(4):1540–1569, 1998.
30. R.R. Sitter and C. Wu, Efficient estimation of quadratic finite population functions in the presence of auxiliary information, *J. Am. Stat. Assoc.*, **97**(458):535–543, 2002.
31. B. Zhang, Bootstrapping with auxiliary information, *Can. J. Stat.*, **27**(2):237–249, 1999.