

***L*-STATISTICS IN COMPLEX SURVEY PROBLEMS¹**

BY JUN SHAO

University of Wisconsin–Madison

We study linear combinations of order statistics (*L*-statistics) in survey problems with a stratified multistage sampling design. Two general types of *L*-statistics are considered: smooth *L*-statistics with weights generated by a smooth function and nonsmooth *L*-statistics (sample quantiles). The trimmed sample mean, the decile mean and variance, the sample Lorenz curve and the sample Gini's family parameters are examples of smooth *L*-statistics or functions of smooth *L*-statistics used in survey problems. It is shown that under weak conditions the smooth *L*-statistics are asymptotically normal and their asymptotic variances can be consistently estimated by jackknifing. For the sample quantiles, their asymptotic normality requires more conditions on the finite population distribution functions. Consistent estimators for the asymptotic variances of the sample quantiles are derived. Asymptotic validity of the Woodruff's confidence intervals for population quantiles is also proved.

1. Introduction. One of the most commonly used class of statistics in statistical analysis is the class of linear combinations of order statistics, termed *L*-statistics in the literature. Simple statistics such as the sample mean, the sample quantiles and the trimmed sample mean are special types of *L*-statistics. Asymptotic properties of *L*-statistics are well known in the case where the data are independent and identically distributed (i.i.d.) observations [e.g., Serfling (1980), Chapter 8]. Theoretical studies of the *L*-statistics in some non-i.i.d. cases have received considerable attention in recent years. For example, Bickel (1973), Koenker and Bassett (1978), Ruppert and Carroll (1980), Koenker and Portnoy (1987) and Welsh (1987) established many asymptotic properties of *L*-statistics in linear regression problems.

The purpose of this work is to study asymptotic properties of *L*-statistics with complex survey data. Sample surveys often use one or a combination of several of the following sampling methods: stratified sampling, cluster sampling, unequal probability sampling and multistage sampling [Kish and Frankel (1974)]. Simple random sampling is rarely used in practice because of both practical and theoretical considerations. Due to the complexity of the sampling design, the resulting data are heavily non-i.i.d.

Throughout the paper we consider the following commonly used stratified multistage sampling design. The population under consideration has been stratified into *L* strata. From each stratum, some first-stage units (clusters)

Received August 1991; revised March 1993.

¹The research was supported by the Natural Sciences and Engineering Research Council of Canada.

AMS 1991 subject classifications. Primary 62D05; secondary 62G20, 62G09.

Key words and phrases. Asymptotic normality, finite population, multistage sampling, trimmed mean, Lorenz curve, Gini's family, quantile, variance estimation, jackknife.

are selected using unequal probability sampling with replacement (or simple random sampling without replacement) and the samples are selected independently across the strata. Within each first-stage unit (cluster), some ultimate units are sampled according to some sampling methods. Note that we do not need to specify the number of stages and the sampling methods used after the first-stage sampling.

Under this stratified multistage sampling design, Krewski and Rao (1981) established some asymptotic results for the sample mean, a special L -statistic and functions of several sample means (e.g., ratio estimators). The use of other L -statistics in complex survey problems has become more and more popular in recent years. For example, the sample quantiles are important statistics for survey problems related to earnings. An asymptotic study is given by Francisco and Fuller (1991). Other examples, including the trimmed means, the Gini's family and Lorenz curve, are described in Section 2.

Section 2 contains an introduction of notation, a formal definition of L -statistics and some examples. The asymptotic normality of the L -statistics in stratified multistage survey problems is established in Sections 3 and 4. Consistent estimators of the asymptotic variances of the L -statistics are also derived. These results are useful for large-sample inferences.

2. Definitions and examples. Let $\{\mathcal{P}_k, k = 1, 2, \dots\}$ be a sequence of finite populations. Throughout the paper k is used as the index of the finite population, but it may be omitted frequently for simplicity. Each \mathcal{P}_k contains L strata, and the h th stratum contains N_h clusters. Associated with the j th ultimate unit in the i th cluster of stratum h is a characteristic $Y_{hij}, j = 1, \dots, N_{hi}; i = 1, \dots, N_h; h = 1, \dots, L$. Here N_{hi} is the number of ultimate units in the i th cluster of stratum h . Note that L, N_h, N_{hi}, Y_{hij} , and so on depend on k also but the subscript k is omitted. Let $I(A)$ be the indicator function of the set A . The k th finite population distribution function is

$$(2.1) \quad F_k(x) = \frac{1}{M} \sum_{h=1}^L \sum_{i=1}^{N_h} \sum_{j=1}^{N_{hi}} I(Y_{hij} \leq x),$$

where $M = \sum_{h=1}^L \sum_{i=1}^{N_h} N_{hi}$.

Suppose that in the first-stage sampling, $n_h \geq 2$ clusters are selected from stratum h (independently across the strata) with probabilities $p_{hi} > 0, i = 1, \dots, N_h; h = 1, \dots, L; \sum_{i=1}^{N_h} p_{hi} = 1$. For a fixed h , if $p_{hi} \equiv 1/N_h$, the sampling method is simple random sampling (equal probability sampling) and the clusters are selected either with replacement or without replacement. If unequal probability sampling is used, we assume that the clusters are selected independently (with replacement). Within each cluster, some second-stage, third-stage, ... units are sampled according to some methods (e.g., stratified sampling, cluster sampling, unequal probability sampling). Let n_{hi} be the number of ultimate units drawn from the i th cluster of stratum h , y_{hij} and \tilde{w}_{hij} be the observed characteristic and the survey weight associated with the (h, i, j) th ultimate unit,

$j = 1, \dots, n_{hi}; i = 1, \dots, n_h; h = 1, \dots, L$. Again, $n_h, n_{hi}, y_{hi,j}$ and $\tilde{w}_{hi,j}$ depend on k but the subscript k is omitted. Suppose that the survey weights are constructed so that

$$(2.2) \quad G_k(x) = \frac{1}{M} \sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{n_{hi}} \tilde{w}_{hi,j} I(y_{hi,j} \leq x)$$

is unbiased for $F_k(x)$, that is,

$$E[G_k(x)] = F_k(x) \quad \text{for any } x$$

[see, e.g., Krewski and Rao (1981)]. However, G_k may not be a distribution function since $G_k(\infty) = M^{-1} \sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{n_{hi}} \tilde{w}_{hi,j}$ is not necessarily equal to 1. Furthermore, in many cases M , the total number of ultimate units in the population, is unknown and therefore G_k is not an estimator. Hence we consider the normalized G_k :

$$(2.3) \quad \hat{F}_k(x) = G_k(x)/q_k,$$

where $q_k = G_k(\infty)$. Then \hat{F}_k is a distribution function for any k and can be used to estimate F_k since the unknown M is cancelled in the ratio in (2.3).

Let $n = \sum_{h=1}^L n_h$ be the number of sampled clusters in the first-stage sampling. We always assume that $n \rightarrow \infty$ as $k \rightarrow \infty$. Denote $\tilde{w}_{hi,j}/M$ by $w_{hi,j}$. Assuming

$$(A1) \quad \max_{j \leq N_{hi}, i \leq N_h, h \leq L} n_{hi} w_{hi,j} = O\left(\frac{1}{n}\right)$$

and using

$$\begin{aligned} \text{Var}[G_k(x)] &\leq \sum_{h=1}^L \sum_{i=1}^{n_h} E\left(\sum_{j=1}^{n_{hi}} w_{hi,j} I(y_{hi,j} \leq x)\right)^2 \\ &\leq E\left(\sum_{h=1}^L \sum_{i=1}^{n_h} n_{hi} \sum_{j=1}^{n_{hi}} w_{hi,j}^2\right) \leq \max_{j, i, h} n_{hi} w_{hi,j}, \end{aligned}$$

we obtain that, for any x , $G_k(x) - F_k(x) \rightarrow_p 0$. In particular, since $q_k = G_k(\infty)$ and $F_k(\infty) = 1$, $q_k - 1 \rightarrow_p 0$ and therefore

$$\hat{F}_k(x) - F_k(x) \rightarrow_p 0 \quad \text{for any } x.$$

Assumption (A1) means that no survey weight is disproportionately large. It is of interest to see what (A1) reduces to in some simple special cases. If the sampling design is only one-stage and simple random sampling is used, then $w_{hi,j} = N_h/Mn_h$ and (A1) is the same as

$$(2.4) \quad \max_{h \leq L} \frac{N_h}{Mn_h} = O\left(\frac{1}{n}\right),$$

an assumption given in Rao and Wu (1985). If the sampling design is a stratified two-stage sampling design and the second-stage sampling is also simple random sampling (within each cluster), then $w_{hij} = N_h N_{hi} / M n_h n_{hi}$ and (A1) reduces to

$$(2.5) \quad \max_{i \leq N_h, h \leq L} \frac{N_h N_{hi}}{M n_h} = O\left(\frac{1}{n}\right).$$

Note that N_{hi} is the number of units in the i th cluster of stratum h . If N_{hi} , $i = 1, \dots, N_h$, are relatively the same or they are bounded, then (2.5) is the same as (2.4).

Let J be a function on the interval $[0, 1]$ and

$$(2.6) \quad T(F) = \int x J(F(x)) dF(x),$$

for any distribution function F . An L -statistic is generally of the form

$$(2.7) \quad T_k = T(\hat{F}_k) + \sum_{t=1}^m a_t \hat{F}_k^{-1}(p_t),$$

where m is a fixed integer, p_t and a_t are constants, $0 < p_t < 1$ and $F^{-1}(p) = \inf\{x: F(x) \geq p\}$ for any distribution function F . Note that $\hat{F}_k^{-1}(p)$ is simply the p -th quantile of the distribution function \hat{F}_k . T_k is also called the L -estimator of $T(F_k) + \sum_{t=1}^m a_t F_k^{-1}(p_t)$. If $a_t = 0$ for all t , then $T_k = T(\hat{F}_k)$ is called a smooth L -statistic (an L -statistic with smooth weights). If $J(t) \equiv 0$, then $T_k = \sum_{t=1}^m a_t \hat{F}_k^{-1}(p_t)$ is a linear combination of sample quantiles and will be called a non-smooth L -statistic.

Let $\{y_{(l)}, l = 1, \dots, \sum_{h=1}^L \sum_{i=1}^{n_h} n_{hi}\}$ be the order statistics of the sample $\{y_{hij}, j = 1, \dots, n_{hi}; i = 1, \dots, n_h; h = 1, \dots, L\}$ and $\omega_l = w_{hij}/q_k$ if $y_{(l)} = y_{hij}$. Then

$$(2.8) \quad T(\hat{F}_k) = \sum_l c_l y_{(l)}, \quad c_l = \omega_l J\left(\sum_{t=1}^l \omega_t\right)$$

and T_k is a linear combination of the order statistics.

We now discuss some examples of useful L -statistics in survey problems. Obviously the (weighted) sample mean ($J \equiv 1$ and $a_t \equiv 0$) and the sample quantile ($J \equiv 0$) are the most commonly used L -statistics. Some other examples are given as follows.

EXAMPLE 1 (Trimmed sample mean). Let $a_t \equiv 0$ and $J(t) = (\beta - \alpha)^{-1} I(\alpha \leq t \leq \beta)$, where $0 \leq \alpha < \beta \leq 1$. The resulting T_k is then a trimmed sample mean. T_k is not necessarily asymptotically normal if $\{F_k\}$ does not have a continuous

and differentiable limit [Stigler (1973)]. Hence we may smooth $J(t)$ and use

$$\tilde{J}(t) = \begin{cases} \frac{2(t - \alpha)}{(a - \alpha)(\beta - \alpha + b - a)}, & \alpha \leq t < a, \\ \frac{2}{(\beta - \alpha + b - a)}, & a \leq t < b, \\ \frac{2(\beta - t)}{(\beta - b)(\beta - \alpha + b - a)}, & b \leq t < \beta, \\ 0, & \text{otherwise,} \end{cases}$$

where a and b are constants satisfying $\alpha < a \leq b < \beta$.

EXAMPLE 2 (Weighted decile mean and variance). In the problems where we study income shares for different quantile groups we need to consider the (weighted) mean and variance of a given group. Suppose that we divide the sample into $D = 10$ ordered decile groups. For $d = 1, \dots, D - 1$, let $\hat{\mu}_d = \sum_l c_{dl} y_{(l)} / \lambda_d$ and $\hat{\rho}_d^2 = \sum_l c_{dl} [y_{(l)} - \hat{\mu}_d]^2 / \lambda_d$, where $\lambda_d = \sum_l c_{dl}$ and c_{dl} is generated according to (2.8) with the function

$$J(t) = J_d(t) = \begin{cases} 1, & \frac{d-1}{D} < t \leq \frac{d}{D}, \\ 0, & \text{otherwise,} \end{cases} \quad d = 1, \dots, D.$$

The decile mean $\hat{\mu}_d$ and variance $\hat{\rho}_d^2$ are functions of λ_d , $\xi_{1d} = \sum_l c_{dl} y_{(l)}$ and $\xi_{2d} = \sum_l c_{dl} y_{(l)}^2$; ξ_{1d} is an L -statistic; λ_d and ξ_{2d} are L -statistics of the more general form considered by Chernoff, Gastwirth and Johns (1967) with $J = J_d$ but with the characteristic $y_{(l)}$ replaced by 1 and $y_{(l)}^2$, respectively.

EXAMPLE 3 (Lorenz curve). The Lorenz curve (or curve of concentration) is one of the most frequently used devices to describe inequality in income or wealth distributions [Beach and Kaliski (1986)]. The sample Lorenz curve ordinates are expressed as $\hat{\gamma}_d / (D\hat{\gamma}_D)$, where $\hat{\gamma}_1 = \hat{\mu}_1$, $\hat{\gamma}_d = (d-1/d)\hat{\gamma}_{d-1} + (1/d)\hat{\mu}_d$, $d = 2, \dots, D$, and $\hat{\mu}_d$ are given in Example 2. The $\hat{\gamma}_d$ are functions of smooth L -statistics.

EXAMPLE 4 (Gini's family). The sample Gini's mean difference is a well-known smooth L -statistic with $J(t) = 4t - 2$ [Serfling (1980), Example 8.2.4B]. A related quantity is the sample Gini's coefficient (a measure of variability), which is equal to

$$\frac{[\text{sample Gini's mean difference} - \int x d\hat{F}_k(x)/M]}{\int x d\hat{F}_k(x)}.$$

This is a function of smooth L -statistics.

More generally, in survey problems we may consider a family of parameters called the Gini's family [Nygard and Sandström (1985)]:

$$(2.9) \quad GI_k = \frac{\int x J(F_k(x)) dF_k(x)}{\int x dF_k(x)}.$$

For example, when $J(t) = 1 - 3(1 - t)^2$, GI_k in (2.9) is the Mehran's measure; when $J(t) = (3t^2 - 1)/2$, GI_k is the Piesch's measure. The sampled version of GI_k is a function of smooth L -statistics.

3. Results for smooth L -statistics. We establish asymptotic normality of smooth L -statistics in this section and will treat nonsmooth L -statistics in the next section. The results can then be extended to general L -statistics of the form (2.7) in a straightforward manner. For smooth L -statistics $T(\widehat{F}_k)$ with T given by (2.6), we can write [see, e.g., Serfling (1980), Chapter 8]

$$T(\widehat{F}_k) - T(F_k) = \int \phi(x, F_k) d\widehat{F}_k(x) - \int Q_k(x) [\widehat{F}_k(x) - F_k(x)] dx,$$

where

$$Q_k(x) = \begin{cases} \frac{\psi(\widehat{F}_k(x)) - \psi(F_k(x))}{\widehat{F}_k(x) - F_k(x)} - J(F_k(x)), & \text{if } \widehat{F}_k(x) \neq F_k(x), \\ 0, & \text{if } \widehat{F}_k(x) = F_k(x), \end{cases}$$

$\psi(t) = \int_0^t J(s) ds$, and

$$\phi(x, F_k) = - \int [I(x \leq y) - F_k(y)] J(F_k(y)) dy$$

is the influence function. Let $z_{hij} = \phi(y_{hij}, F_k)$ and $Z_{hij} = \phi(Y_{hij}, F_k)$. Then

$$(3.1) \quad \int \phi(x, F_k) d\widehat{F}_k(x) = \sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{n_{hi}} \frac{w_{hij} z_{hij}}{q_k} = \frac{\bar{z}}{q_k}$$

and

$$E\bar{z} = \frac{1}{M} \sum_{h=1}^L \sum_{i=1}^{N_h} \sum_{j=1}^{N_{hi}} Z_{hij} = 0.$$

There are two types of J -functions, trimmed and untrimmed J , and we will treat them separately.

3.1. Trimmed L -statistics. A trimmed J satisfies

$$(3.2) \quad J(t) = 0 \quad \text{if } t < \alpha \text{ or } t > \beta,$$

where α and β are constants satisfying $0 < \alpha < \beta < 1$. The corresponding L -statistic is called a trimmed L -statistic. The trimmed sample mean in Example 1 and the weighted decile mean and variance in Example 2 are trimmed L -statistics. Let

$$(3.3) \quad \sigma_k^2 = \text{Var} \left(\sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{n_{hi}} w_{hij} z_{hij} \right).$$

THEOREM 1. *Suppose that assumption (A1) and the following assumptions hold:*

(A2) $\liminf_k \rightarrow \infty n \sigma_k^2 > 0$.

(A3) *The function J is trimmed [i.e., J satisfies (3.2)], bounded and $m(D) = 0$, where m is the Lebesgue measure and $D = \{x: J \text{ is discontinuous at a limit point of } \{F_k(x), k = 1, 2, \dots\}\}$.*

(A4) *There are constants c_α and c_β such that*

$$\sup_k F_k(c_\alpha) < \alpha \quad \text{and} \quad \inf_k F_k(c_\beta) > \beta.$$

Then

$$(3.4) \quad T(\widehat{F}_k) = T(F_k) + \bar{z}/q_k + o_p(n^{-1/2})$$

and

$$[T(\widehat{F}_k) - T(F_k)]/\sigma_k \rightarrow_{\mathcal{L}} N(0, 1).$$

REMARKS.

(i) As we discussed in Section 2, the condition $m(D) = 0$ in (A3) cannot be relaxed. For the trimmed mean with $J(t) = (\beta - \alpha)^{-1} I(\alpha \leq t \leq \beta)$, $D = \{x: \alpha \text{ or } \beta \text{ is a limit point of } \{F_k(x)\}\}$ and J satisfies (A3) if $m(D) = 0$. Assumption (A3) is not needed if we use the smoothly trimmed sample mean in Example 1.

(ii) Assumption (A4) is much weaker than the tightness of the sequence $\{F_k\}$.

PROOF OF THEOREM 1. Under assumptions (A3) and (A4), $\{z_{hij}\}$ is bounded. Under assumptions (A1)–(A4), Liapounov's condition holds for the "weighted average" \bar{z} in (3.1). Thus, by Bickel and Freedman [(1984), Theorem 3] or Krewski and Rao [(1981), Theorem 3.1] $(\bar{z} - E\bar{z})/\sigma_k \rightarrow_{\mathcal{L}} N(0, 1)$. Since $E\bar{z} = 0$ and $q_k \rightarrow_p 1$, $z/(q_k \sigma_k) \rightarrow_{\mathcal{L}} N(0, 1)$ and the second assertion of the theorem follows from the first assertion. Since

$$T(\widehat{F}_k) - T(F_k) = \bar{z}/q_k - \int Q_k(x) [\widehat{F}_k(x) - F_k(x)] dx,$$

the result follows if

$$(3.5) \quad n^{1/2} \int Q_k(x) [\widehat{F}_k(x) - F_k(x)] dx \rightarrow_p 0.$$

Define

$$\mathcal{A}_k = \{\widehat{F}_k(c_\beta) > \beta \text{ and } \widehat{F}_k(c_\alpha) < \alpha\}.$$

From assumption (A4) and the fact that $\widehat{F}_k(c_\alpha) - F_k(c_\alpha) \rightarrow_p 0$ and $\widehat{F}_k(c_\beta) - F_k(c_\beta) \rightarrow_p 0$, we have $P(\mathcal{A}_k) \rightarrow 1$. On the event \mathcal{A}_k ,

$$\begin{aligned} (3.6) \quad & \left| \int Q_k(x) [\widehat{F}_k(x) - F_k(x)] dx \right| \\ &= \left| \int_{c_\alpha}^{c_\beta} Q_k(x) [\widehat{F}_k(x) - F_k(x)] dx \right| \\ &\leq \left\{ \int_{c_\alpha}^{c_\beta} [Q_k(x)]^2 dx \int_{c_\alpha}^{c_\beta} [\widehat{F}_k(x) - F_k(x)]^2 dx \right\}^{1/2}. \end{aligned}$$

From the assumption (A3) and $\widehat{F}_k(x) - F_k(x) \rightarrow_p 0$, $Q_k(x) \rightarrow_p 0$ for $x \notin D$. Also, $\|Q_k\|_\infty \leq 2\|J\|_\infty$. Hence

$$(3.7) \quad \int_{c_\alpha}^{c_\beta} [Q_k(x)]^2 dx \rightarrow_p 0.$$

For any x , $E[G_k(x) - F_k(x)]^2 = \text{Var}[G_k(x)] \leq \max_{j,i,h} n_{hi} w_{hij} \leq cn^{-1}$ by assumption (A1), where c is a constant. Hence

$$(3.8) \quad E \int_{c_\alpha}^{c_\beta} [G_k(x) - F_k(x)]^2 dx \leq \frac{c(c_\beta - c_\alpha)}{n} = O(n^{-1}).$$

Since $\widehat{F}_k = G_k/q_k$, we have

$$\begin{aligned} \int_{c_\alpha}^{c_\beta} [\widehat{F}_k(x) - F_k(x)]^2 dx &\leq 2 \int_{c_\alpha}^{c_\beta} [G_k(x) - F_k(x)]^2 dx + 2 \int_{c_\alpha}^{c_\beta} [\widehat{F}_k(x) - G_k(x)]^2 dx \\ &\leq 2 \int_{c_\alpha}^{c_\beta} [G_k(x) - F_k(x)]^2 dx + 2(q_k - 1)^2(c_\beta - c_\alpha) \\ &= O_p(n^{-1}), \end{aligned}$$

by (3.8) and $q_k - 1 = O_p(n^{-1/2})$. This, together with (3.6) and (3.7), implies (3.5). The proof is complete. \square

3.2. Untrimmed L-statistics. For untrimmed L-statistics, some moment conditions are necessary. We only consider smooth functions J .

THEOREM 2. Suppose that assumptions (A1), (A2) and the following hold:

(A3') The function J is Lipschitz continuous on $[0, 1]$;

(A4') $\sup_k \int_{-\infty}^{\infty} |x| dF_k(x) < \infty$ and there is a $\delta > 0$ such that

$$(3.9) \quad n^{1+\delta} \sum_{h=1}^L \sum_{i=1}^{n_h} E|u_{hi} - Eu_{hi}|^{2(1+\delta)} \rightarrow 0,$$

where $u_{hi} = \sum_{j=1}^{n_{hi}} w_{hij} z_{hij}$.

Then the conclusions of Theorem 1 still hold.

REMARK. Note that (3.9) is a Liapounov-type condition similar to condition C1 in Krewski and Rao (1981).

PROOF OF THEOREM 2. Under (A2) and (3.9), $(\bar{z} - E\bar{z})/\sigma_k \rightarrow_{\mathcal{L}} N(0, 1)$. Thus, from the proof of Theorem 1, we only need to show that (3.5) holds. Since J is Lipschitz continuous, there is a constant $c_J > 0$ such that $|J(t) - J(s)| \leq c_J |t - s|$, $t, s \in [0, 1]$. Then

$$\left| \int Q_k(x) [\hat{F}_k(x) - F_k(x)] dx \right| \leq c_J \int [\hat{F}_k(x) - F_k(x)]^2 dx.$$

From $\hat{F}_k = G_k/q_k$,

$$\begin{aligned} & \int_{-\infty}^0 [\hat{F}_k(x) - F_k(x)]^2 dx \\ & \leq 2 \int_{-\infty}^0 [G_k(x) - F_k(x)]^2 dx + 2 \int_{-\infty}^0 [\hat{F}_k(x) - G_k(x)]^2 dx \\ & = 2 \int_{-\infty}^0 [G_k(x) - F_k(x)]^2 dx + 2(q_k^{-1} - 1)^2 \int_{-\infty}^0 [G_k(x)]^2 dx. \end{aligned}$$

Note that $(q_k^{-1} - 1)^2 = q_k^{-2}(q_k - 1)^2 = O_p(n^{-1})$,

$$E \int_{-\infty}^0 [G_k(x)]^2 dx \leq \int_{-\infty}^0 E[G_k(x)] dx = \int_{-\infty}^0 F_k(x) dx \leq \int |x| dF_k(x)$$

and

$$\begin{aligned} E \int_{-\infty}^0 [G_k(x) - F_k(x)]^2 dx &= \int_{-\infty}^0 \text{Var}[G_k(x)] dx \\ &\leq \int_{-\infty}^0 \sum_{h=1}^L \sum_{i=1}^{n_h} E \left[\sum_{j=1}^{n_{hi}} w_{hij} I(y_{hij} \leq x) \right]^2 dx \\ &\leq O(n^{-1}) \int_{-\infty}^0 F_k(x) dx = O(n^{-1}), \end{aligned}$$

where the last inequality follows from assumption (A1). Hence

$$\int_{-\infty}^0 [\hat{F}_k(x) - \hat{F}_k(x)]^2 dx = O_p(n^{-1}).$$

Similarly,

$$\int_0^\infty [\hat{F}_k(x) - \hat{F}_k(x)]^2 dx = O_p(n^{-1}),$$

and therefore (3.5) follows. This completes the proof. \square

The J -functions corresponding to the untrimmed L -statistics in Examples 1–4 are clearly Lipschitz continuous on $[0, 1]$.

In some cases we need to consider a statistic S_k which is a function of several L -statistics (Examples 2–4), say $T_k^{(j)}$, $j = 1, \dots, m$. The asymptotic distribution of S_k can be obtained by using the results in Theorems 1 and 2. Since each $T_k^{(j)}$ can be approximated by a weighted average according to (3.4), we can apply the central limit theorem and the δ -method. The details are omitted.

3.3. Variance estimation by jackknife. For various purposes in statistical analysis we need to estimate the unknown asymptotic variance of a given statistic. When the given statistic is a smooth L -statistic, a formula for its asymptotic variance is given by (3.3). When the given statistic is a known function of several L -statistics, a formula for its asymptotic variance can be obtained by using the δ -method and calculating partial derivatives of the known function. In some cases, the asymptotic variance can be estimated by substituting the unknown quantities in the formula of the asymptotic variance by some estimators. This is called the substitution method. Another very popular method in sample surveys is the jackknife [Tukey (1958)]. Unlike the substitution method, the jackknife does not require knowing a formula of the asymptotic variance, which is one of the reasons why the jackknife is so attractive. Note that for statistics such as the sample Lorenz curve ordinates (Example 3), the derivation of their asymptotic variance is quite complicated and tedious. The jackknife replaces the theoretical derivation by repeated computations of the given statistic. Many agencies (e.g., Statistics Canada) have computer software to implement the computation of the jackknife variance estimators.

Variance estimation by jackknife is based on the variability among a number of replicate statistics computed from overlapping subsamples of the total sample. Let

$$G_k^{(gl)}(x) = \sum_{h \neq g} \sum_{i=1}^{n_h} \sum_{j=1}^{n_{hi}} w_{hij} I(y_{hij} \leq x) + \frac{n_g}{n_g - 1} \sum_{i \neq l} \sum_{j=1}^{n_{gi}} w_{gij} I(y_{gij} \leq x),$$

$$q_k^{(gl)} = G_k^{(gl)}(\infty), \quad \widehat{F}_k^{(gl)} = \frac{\widehat{G}_k^{(gl)}}{q_k^{(gl)}}$$

$$T_k^{(gl)} = T(\widehat{F}_k^{(gl)}), \quad l = 1; \dots, n_g; g = 1, \dots, L.$$

The jackknife estimator of the asymptotic variance of $T_k = T(\widehat{F}_k)$ is

$$(3.10) \quad \widehat{\sigma}_k^2 = \sum_{h=1}^L \frac{(1-f_h)(n_h-1)}{n_h} \sum_{i=1}^{n_h} (T_k^{(hi)} - T_k)^2,$$

where $f_h = n_h/N_h$ if the first-stage sampling is without replacement and $f_h = 0$ if the first-stage sampling is with replacement. When the statistic T_k is a function

of several L -statistics, its jackknife estimator of the asymptotic variance is still given by (3.10) with $T_k^{(hi)}$ being defined as the statistic based on the subsample obtained by removing the i th cluster of the h th stratum.

We now establish the consistency of the jackknife variance estimator $\hat{\sigma}_k^2$ in (3.10). First, consider the case where $T_k = T(\hat{F}_k)$.

THEOREM 3 (Trimmed L -statistics). *Suppose that assumptions (A1)–(A4) and the following hold:*

(A5) *There are two sequence of sets $\{N_{1k}\}$ and $\{N_{2k}\}$ such that, for each k ,*

$$\{n_h, h = 1, \dots, L\} = N_{1k} \cup N_{2k},$$

$$\sup_k \max\{n_h \in N_{1k}\} < \infty \quad \text{and} \quad \lim_{k \rightarrow \infty} \min\{n_h \in N_{2k}\} = \infty.$$

Then the jackknife estimator $\hat{\sigma}_k^2$ in (3.10) is consistent for σ_k^2 in (3.3), that is,

$$\hat{\sigma}_k^2 / \sigma_k^2 \rightarrow_p 1.$$

REMARK. Assumption (A5) is satisfied in the following two common situations in surveys: (i) all the n_h are small (bounded by a constant); (ii) all the n_h are large. Assumption (A5) allows a mixture of some small n_h 's and some large n_h 's.

PROOF OF THEOREM 3. Let a_{hij} denote the function $a_{hij}(x) = w_{hij}I(y_{hij} \leq x)$. Then

$$\begin{aligned} \hat{F}_k^{(hi)} - \hat{F}_k &= \frac{G_k^{(hi)}}{q_k^{(hi)}} - \frac{G_k}{q_k} = G_k^{(hi)} \left(\frac{1}{q_k^{(hi)}} - \frac{1}{q_k} \right) + \frac{1}{q_k} (G_k^{(hi)} - G_k) \\ &= G_k^{(hi)} \left(\frac{1}{q_k^{(hi)}} - \frac{1}{q_k} \right) + \frac{n_h}{q_k(n_h - 1)} \left(\frac{1}{n_h} \sum_{i=1}^{n_h} \sum_{j=1}^{n_{hi}} a_{hij} - \sum_{j=1}^{n_{hi}} a_{hij} \right). \end{aligned}$$

Let u_{hi} be defined in (3.9), $\bar{u}_h = n_h^{-1} \sum_{i=1}^{n_h} u_{hi}$ and $\bar{z}^{(hi)} = \int \phi(x, F_k) dG_k^{(hi)}(x)$. Then

$$\begin{aligned} T_k^{(hi)} - T_k &= \frac{n_h}{q_k(n_h - 1)} (\bar{u}_h - u_{hi}) + \bar{z}^{(hi)} \left(\frac{1}{q_k^{(hi)}} - \frac{1}{q_k} \right) \\ &\quad - \int Q_k^{(hi)}(x) [\hat{F}_k^{(hi)}(x) - \hat{F}_k(x)] dx, \end{aligned}$$

where

$$Q_k^{(hi)}(x) = \begin{cases} \frac{\psi(\hat{F}_k^{(hi)}(x)) - \psi(\hat{F}_k(x))}{\hat{F}_k^{(hi)}(x) - \hat{F}_k(x)} - J(F_k(x)), & \text{if } \hat{F}_k^{(hi)}(x) \neq \hat{F}_k(x), \\ 0, & \text{otherwise,} \end{cases}$$

and $\psi(t) = \int_0^t J(s) ds$. It follows from (3.10) that

$$(3.11) \quad \hat{\sigma}_k^2 = \alpha_{1k} + \alpha_{2k} + \alpha_{3k} + 2(\text{cross-product terms}),$$

where

$$(3.12) \quad \alpha_{1k} = \sum_{h=1}^L \frac{(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} \frac{(u_{hi} - \bar{u}_h)^2}{q_k^2},$$

$$(3.13) \quad \alpha_{2k} = \sum_{h=1}^L \frac{(1-f_h)(n_h-1)}{n_h} \sum_{i=1}^{n_h} (\bar{z}_h^{(hi)})^2 \left(\frac{1}{q_k^{(hi)}} - \frac{1}{q_k} \right)^2,$$

$$(3.14) \quad \alpha_{3k} = \sum_{h=1}^L \frac{(1-f_h)(n_h-1)}{n_h} \sum_{i=1}^{n_h} (r_k^{(hi)})^2$$

and

$$(3.15) \quad r_k^{(hi)} = - \int Q_k^{(hi)}(x) [\hat{F}_k^{(hi)}(x) - \hat{F}_k(x)] dx.$$

To establish the result, it suffices to show that

$$(3.16) \quad n(\alpha_{1k} - \sigma_k^2) \rightarrow_p 0,$$

$$(3.17) \quad n\alpha_{2k} \rightarrow_p 0 \quad \text{and} \quad n\alpha_{3k} \rightarrow_p 0,$$

since under (3.16) and (3.17), the cross-product terms are of the order $o_p(n^{-1})$.

PROOF OF (3.16). Let

$$\sigma_{1k}^2 = \sum_{n_h \in \mathcal{N}_{1k}} \text{Var} \left(\sum_{i=1}^{n_h} u_{hi} \right) \quad \text{and} \quad \sigma_{2k}^2 = \sum_{n_h \in \mathcal{N}_{2k}} \text{Var} \left(\sum_{i=1}^{n_h} u_{hi} \right).$$

Then $\sigma_k^2 = \sigma_{1k}^2 + \sigma_{2k}^2$ and

$$E \left[\sum_{n_h \in \mathcal{N}_{1k}} \frac{n(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} (u_{hi} - \bar{u}_h)^2 \right] = n\sigma_{1k}^2.$$

Under assumptions (A1) and (A5), there is a constant $c > 0$ such that

$$n^{1+\delta} \sum_{n_h \in \mathcal{N}_{1k}} E \left| \frac{(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} (u_{hi} - \bar{u}_h)^2 \right|^{1+\delta} \leq cn^{1+\delta} \sum_{n_h \in \mathcal{N}_{1k}} \sum_{i=1}^{n_h} E u_{hi}^{2(1+\delta)} \rightarrow 0.$$

Hence

$$(3.18) \quad \sum_{n_h \in \mathcal{N}_{1k}} \frac{n(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} (u_{hi} - \bar{u}_h)^2 - n\sigma_{1k}^2 \rightarrow_p 0.$$

Let $\mu_h = Eu_{hi}$. Since

$$\begin{aligned} & E \left[\sum_{n_h \in \mathcal{N}_{2k}} \frac{n(1-f_h)n_h^2}{n_h-1} (\mu_h - \bar{u}_h)^2 \right] \\ &= \sum_{n_h \in \mathcal{N}_{2k}} \frac{n(1-f_h)n_h^2}{n_h-1} \text{Var}(\bar{u}_h) \\ &\leq \frac{n}{\min\{n_h \in \mathcal{N}_{2k}\} - 1} \sum_{n_h \in \mathcal{N}_{2k}} \text{Var} \left(\sum_{i=1}^{n_h} u_{hi} \right) \rightarrow 0 \end{aligned}$$

[by assumptions (A1)–(A5)] and

$$\begin{aligned} & E \left[\sum_{n_h \in \mathcal{N}_{2k}} \frac{n(1-f_h)}{n_h-1} \sum_{i=1}^{n_h} (u_{hi} - \mu_h)^2 \right] \\ &\leq \frac{n}{\min\{n_h \in \mathcal{N}_{2k}\} - 1} \sum_{n_h \in \mathcal{N}_{2k}} \text{Var} \left(\sum_{i=1}^{n_h} u_{hi} \right) \rightarrow 0, \end{aligned}$$

we have

$$\begin{aligned} \sum_{n_h \in \mathcal{N}_{2k}} \frac{n(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} (u_{hi} - \bar{u}_h)^2 &= \sum_{n_h \in \mathcal{N}_{2k}} \frac{n(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} (u_{hi} - \mu_h)^2 \\ &\quad - \sum_{n_h \in \mathcal{N}_{2k}} \frac{n(1-f_h)n_h^2}{n_h-1} \sum_{i=1}^{n_h} (\mu_h - \bar{u}_h)^2 \\ &= \sum_{n_h \in \mathcal{N}_{2k}} n(1-f_h) \sum_{i=1}^{n_h} (u_{hi} - \mu_h)^2 + o_p(1). \end{aligned}$$

Since

$$E \left[\sum_{n_h \in \mathcal{N}_{2k}} n(1-f_h) \sum_{i=1}^{n_h} (u_{hi} - \mu_h)^2 \right] = n\sigma_{2k}^2 + o(1)$$

and

$$n^{1+\delta} \sum_{n_h \in \mathcal{N}_{2k}} \sum_{i=1}^{n_h} (1-f_h)^{1+\delta} E|u_{hi} - \mu_h|^{2(1+\delta)} \rightarrow 0,$$

we have

$$(3.19) \quad \sum_{n_h \in \mathcal{N}_{2k}} \frac{n(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} (u_{hi} - \bar{u}_h)^2 - n\sigma_{2k}^2 \rightarrow 0.$$

Since $q_k^2 \rightarrow_p 1$, (3.16) follows from (3.12), (3.18) and (3.19). \square

PROOF OF (3.17). Let

$$\zeta_k = \sum_{h=1}^L \frac{(1-f_h)(n_h-1)}{n_h} \sum_{i=1}^{n_h} \left(\frac{1}{q_k^{(hi)}} - \frac{1}{q_k} \right)^2.$$

Then ζ_k is the jackknife estimator of the asymptotic variance of q_k^{-1} . Thus $n\zeta_k = O_p(1)$. Since $\bar{z}^{(hi)}$ are weighted averages, it can be shown that $\max_{h,i} |\bar{z}^{(hi)} - \bar{z}| \rightarrow_p 0$; but $E\bar{z} = 0$ and $\bar{z} - E\bar{z} \rightarrow_p 0$. Hence, from (3.13),

$$n\alpha_{2k} \leq \max_{h,i} (\bar{z}^{(hi)})^2 n\zeta_k = o_p(1).$$

Note that $\max_{h,i} \|\widehat{F}_k^{(hi)} - \widehat{F}_k\|_\infty \rightarrow 0$. Hence, similar to the proof of Theorem 1,

$$(r_k^{(hi)})^2 \leq \int_{c_\alpha}^{c_\beta} [Q_k^{(hi)}(x)]^2 dx \int_{c_\alpha}^{c_\beta} [\widehat{F}_k^{(hi)}(x) - \widehat{F}_k(x)]^2 dx$$

and

$$\max_{h,i} \int_{c_\alpha}^{c_\beta} [Q_k^{(hi)}(x)]^2 dx \rightarrow_p 0.$$

From (3.14) and (3.15),

$$\begin{aligned} n\alpha_{3k} &\leq \left(\max_{h,i} \int_{c_\alpha}^{c_\beta} [Q_k^{(hi)}(x)]^2 dx \right) n \sum_{h=1}^L \frac{(1-f_h)(n_h-1)}{n_h} \\ &\quad \times \sum_{i=1}^{n_h} \int_{c_\alpha}^{c_\beta} [\widehat{F}_k^{(hi)}(x) - \widehat{F}_k(x)]^2 dx \\ &\leq o_p(1) \left[n\zeta_k + n \sum_{h=1}^L \frac{(1-f_h)n_h}{n_h-1} \frac{1}{q_k^2} \sum_{i=1}^{n_h} \left(\sum_{j=1}^{n_{hi}} w_{hij} \right)^2 \right] = o_p(1), \end{aligned}$$

since, by (A1),

$$n \sum_{h=1}^L \sum_{i=1}^{n_h} \left(\sum_{j=1}^{n_{hi}} w_{hij} \right)^2 \leq n \max_{h,i,j} (n_{hi} w_{hij}) \sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{n_{hi}} w_{hij} = O_p(1).$$

This proves (3.17). The proof of Theorem 3 now is complete. \square

THEOREM 4 (Untrimmed L-statistics). *The conclusion of Theorem 3 holds also for untrimmed L-statistics if assumptions (A3) and (A4) are replaced by (A3') and (A4'), respectively.*

PROOF. Note that (3.11)–(3.15) still hold and we need to show (3.16) and (3.17). Using (A4') and a proof similar to that of Theorem 3, we can show (3.16)

and the first assertion of (3.17). It remains to show the second assertion of (3.17). From the Lipschitz continuity of J ,

$$(r_k^{(hi)})^2 \leq c_J^2 \left\{ \int [\widehat{F}_k^{(hi)}(x) - \widehat{F}_k(x)]^2 dx \right\}^2.$$

Note that $\max_{h,i} |1/q_k^{(hi)} - 1/q_k| = o_p(1)$ and $\max_{h,i} 1/q_k^{(hi)} = O_p(1)$. Then

$$\begin{aligned} & \sum_{h=1}^L \frac{n(1-f_h)(n_h-1)}{n_h} \sum_{i=1}^{n_h} \left\{ \int_{-\infty}^0 [G_k(x)]^2 \left(\frac{1}{q_k^{(hi)}} - \frac{1}{q_k} \right)^2 dx \right\}^2 \\ & \leq \max_{h,i} \left(\frac{1}{q_k^{(hi)}} - \frac{1}{q_k} \right)^2 \left\{ \int_{-\infty}^0 [G_k(x)]^2 dx \right\}^2 n \zeta_k \\ & = o_p(1) O_p(1) O_p(1) = o_p(1), \end{aligned}$$

and by (A1) and (A4'),

$$\begin{aligned} & \sum_{h=1}^L \frac{n(1-f_h)(n_h-1)}{n_h} \sum_{i=1}^{n_h} \left(\frac{1}{q_k^{(hi)}} \right)^2 \left\{ \int_{-\infty}^0 [G_k^{(hi)}(x) - G_k(x)]^2 dx \right\}^2 \\ & \leq O_p(n^{-1}) \sum_{h=1}^L \sum_{i=1}^{n_h} \left[\int_{-\infty}^0 \sum_{j=1}^{n_{hi}} a_{hij}(x) dx \right]^2 \\ & \leq O_p(n^{-1}) \left[\int_{-\infty}^0 G_k(x) dx \right]^2 = O_p(n^{-1}). \end{aligned}$$

From $\widehat{F}_k^{(hi)} - \widehat{F}_k = (G_k^{(hi)} - G_k)/q_k^{(hi)} + G_k(1/q_k^{(hi)} - 1/q_k)$, we have

$$\sum_{h=1}^L \frac{n(1-f_h)(n_h-1)}{n_h} \sum_{i=1}^{n_h} \left\{ \int_{-\infty}^0 [\widehat{F}_k^{(hi)}(x) - \widehat{F}_k(x)]^2 dx \right\}^2 \rightarrow_p 0.$$

Similarly,

$$\sum_{h=1}^L \frac{n(1-f_h)(n_h-1)}{n_h} \sum_{i=1}^{n_h} \left\{ \int_0^{\infty} [\widehat{F}_k^{(hi)}(x) - \widehat{F}_k(x)]^2 dx \right\}^2 \rightarrow_p 0.$$

This proves (3.17) and completes the proof. \square

When the statistic of interest is a function of several smooth L -statistics, the consistency of its jackknife variance estimator can be easily established by using the results in Theorems 3 and 4 and the argument used in the proof of Theorem 3.4 of Krewski and Rao (1981), who proved the consistency of the jackknife variance estimator for a function of several sample means. Therefore, our results can be applied to Examples 1–4 for variance estimation.

4. Results for sample quantiles. For nonsmooth L -statistics, we focus on the sample p th quantile with a given $p \in (0, 1)$. The case of functions of several sample quantiles can be treated by using the δ -method.

Let $\theta_k = F_k^{-1}(p)$ be the population p th quantile, and let $\hat{\theta}_k = \hat{F}_k^{-1}(p)$ be the sample p th quantile. Francisco and Fuller (1991) established a Bahadur-type representation for $\hat{\theta}_k$ in complex survey problems. A similar result was obtained in Shao and Wu (1992) for the case of one-stage stratified simple random sampling. Francisco and Fuller (1991) also proposed an estimator [see (4.15)] of the asymptotic variance of $\hat{\theta}_k$ and proved its consistency under their conditions 1–7. However, some of these conditions (conditions 5–7) are unnecessary. In this section we derive an estimator of the asymptotic variance of $\hat{\theta}_k$ and establish its consistency under weaker conditions.

Suppose that assumption (A1) and the following assumptions hold.

(A6) The sequence $\{\theta_k\}$ is bounded.

(A7) There is a sequence of functions $\{f_k(\cdot)\}$ such that

$$(4.1) \quad \lim_{k \rightarrow \infty} \left[\frac{F_k(\theta_k + \delta_k) - F_k(\theta_k)}{\delta_k} - f_k(\theta_k) \right] = 0,$$

for any sequence $\{\delta_k\}$ of the order $O(n^{-1/2})$ and

$$0 < \inf_k f_k(\theta_k) \leq \sup_k f_k(\theta_k) < \infty.$$

Then, by using a proof similar to those in Francisco and Fuller (1991) and Shao and Wu (1992), we can show that

$$(4.2) \quad \hat{\theta}_k = \theta_k + \frac{F_k(\theta_k) - \hat{F}_k(\theta_k)}{f_k(\theta_k)} + o_p(n^{-1/2}).$$

Furthermore, it can be shown that

$$(4.3) \quad n^{1/2} [\hat{F}_k^{-1}(p_k) - \hat{F}_k^{-1}(p)] = \frac{c}{f_k(\theta_k)} + o_p(1),$$

where $p_k = p + cn^{-1/2}$ and c is a constant, under the following assumption:

$$(A8) \quad \lim_{k \rightarrow \infty} \left[\frac{F_k^{-1}(p_k) - F_k^{-1}(p)}{p_k - p} - \frac{1}{f_k(\theta_k)} \right] = 0 \text{ and } \lim_{k \rightarrow \infty} |f_k(F_k^{-1}(p)) - f_k(\theta_k)| = 0.$$

Assumptions (A7) and (A8) are essentially the same as Condition 4 in Francisco and Fuller (1991). Condition (4.1) is a type of smoothness condition requiring that F_k is nearly differentiable at θ_k when k is large, although F_k is not differentiable for each fixed k . When the units in the finite populations \mathcal{P}_k are samples from continuous and differentiable superpopulations, (4.1) is satisfied.

As a direct consequence of (4.2),

$$(4.4) \quad \frac{\hat{\theta}_k - \theta_k}{\sigma_k} \rightarrow_{\mathcal{L}} N(0, 1),$$

where

$$(4.5) \quad \sigma_k^2 = \frac{v_k(\theta_k)}{[f_k(\theta_k)]^2}$$

is the asymptotic variance of $\hat{\theta}_k$,

$$(4.6) \quad v_k(\theta_k) = F_k^2(\theta_k) \text{Var}(q_k) - 2F_k(\theta_k) \text{Cov}[q_k, G_k(\theta_k)] + \text{Var}[G_k(\theta_k)]$$

is the asymptotic variance of $\hat{F}_k(\theta_k)$, and $\liminf_k n v_k(\theta_k) > 0$ is assumed.

To estimate σ_k^2 in (4.5), however, is much more difficult (or requires much stronger conditions) than to estimate the asymptotic variance of a smooth L -statistic. Unlike the case of smooth L -statistics, the jackknife variance estimator may not work for sample quantiles. More complicated data-resampling methods such as the delete- d jackknife [Shao and Wu (1989)], the balanced repeated replication [McCarthy (1969) and Shao and Wu (1992)] and the bootstrap [Efron (1979)] need to be considered. The estimator of σ_k^2 considered here is of the form

$$(4.7) \quad \hat{\sigma}_k^2 = \hat{v}_k(\hat{\theta}_k) \hat{u}_k^2,$$

where

$$(4.8) \quad \hat{u}_k = n^{1/2} [\hat{F}_k^{-1}(p + n^{-1/2}) - \hat{F}_k^{-1}(p - n^{-1/2})] / 2$$

is a consistent estimator of $1/f_k(\theta_k)$ by (4.3),

$$\begin{aligned} \hat{v}_k(x) &= \frac{\hat{F}_k^2(x) \widehat{\text{var}}(\infty) - 2\hat{F}_k(x) \widehat{\text{cov}}(x) + \widehat{\text{var}}(x)}{q_k^2}, \\ \widehat{\text{var}}(x) &= \sum_{h=1}^L \frac{(1-f_h)n_h}{n_h-1} \sum_{l=1}^{n_h} \left[a_{hi}(x) - \frac{1}{n_h} \sum_{l=1}^{n_h} a_{hl}(x) \right]^2, \\ \widehat{\text{cov}}(x) &= \sum_{h=1}^L \frac{(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} \left[a_{hi}(x) - \frac{1}{n_h} \sum_{l=1}^{n_h} a_{hl}(x) \right] \left[a_{hi}(\infty) - \frac{1}{n_h} \sum_{l=1}^{n_h} a_{hl}(\infty) \right] \end{aligned}$$

and $a_{hi}(x) = \sum_{j=1}^{n_{hi}} a_{hij}(x) = \sum_{j=1}^{n_{hi}} w_{hij} I(y_{hij} \leq x)$.

Note that, for each fixed x , $\widehat{\text{var}}(x)$ and $\widehat{\text{cov}}(x)$ are unbiased and consistent estimators of $\text{Var}[G_k(x)]$ and $\text{Cov}[q_k, G_k(x)]$, respectively. Naturally, it is expected that $\hat{v}_k(\theta_k)$ is consistent for $v_k(\theta_k)$ in (4.6); but $\hat{v}_k(\theta_k)$ is not an estimator since θ_k is unknown. We have to replace θ_k by $\hat{\theta}_k$ and use $\hat{v}_k = \hat{v}_k(\hat{\theta}_k)$ as an estimator of $v_k(\theta_k)$. The estimator \hat{v}_k was also used by Francisco and Fuller (1991). The following result shows the consistency of \hat{v}_k and $\hat{\sigma}_k^2$.

THEOREM 5. Assume (A1), (A6) and the following:

(A9) For any $x \in \Theta$,

$$\lim_{t \rightarrow 0} \lim_{k \rightarrow \infty} [F_k(x+t) - F_k(x)] = 0,$$

where Θ is the set of limit points of the sequence $\{\theta_k\}$. Then

$$(4.9) \quad n[\hat{v}_k - v_k(\theta_k)] \rightarrow_p 0.$$

If (A7) and (A8) also hold, then $\hat{\sigma}_k^2$ in (4.7) is consistent for σ_k^2 in (4.5), that is,

$$(4.10) \quad \hat{\sigma}_k^2 / \sigma_k^2 \rightarrow_p 1.$$

REMARK. Similar to (4.1), assumption (A9) requires that F_k is nearly continuous at each $x \in \Theta$ when k is large; (A9) holds if $\{F_k\}$ has a continuous limit. If $\theta_k = \theta$ for all k , then (A9) is implied by (4.1).

PROOF OF THEOREM 5. Since (4.10) follows from (4.3), (4.7) and (4.9), we only need to show (4.9). Using a subsequence argument, we can assume that there is a $\theta \in \Theta$ such that

$$\theta_k \rightarrow \theta \quad \text{and} \quad \hat{\theta}_k \rightarrow \theta \quad \text{a.s.}$$

Define

$$\begin{aligned} \hat{\tau}_{1k}(x) &= \sum_{h=1}^L \frac{(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} a_{hi}^2(x), \\ \hat{\tau}_{2k}(x) &= \sum_{h=1}^L \frac{1-f_h}{n_h-1} \left[\sum_{i=1}^{n_h} a_{hi}(x) \right]^2, \\ \hat{\tau}_{3k}(x) &= \sum_{h=1}^L \frac{(1-f_h)n_h}{n_h-1} \sum_{i=1}^{n_h} a_{hi}(x) a_{hi}(\infty), \\ \hat{\tau}_{4k}(x) &= \sum_{h=1}^L \frac{1-f_h}{n_h-1} \left[\sum_{i=1}^{n_h} a_{hi}(x) \right] \left[\sum_{i=1}^{n_h} a_{hi}(\infty) \right] \end{aligned}$$

and

$$\tau_{1k}(x) = E\hat{\tau}_{1k}(x), \quad l = 1, \dots, 4.$$

Let R_0 be the collection of all rational numbers. Under (A1), $n[\hat{\tau}_{1k}(x) - \tau_{1k}(x)] \rightarrow 0$ a.s., for any $x \in R_0$. Hence, almost surely,

$$(4.11) \quad n[\hat{\tau}_{1k}(x) - \tau_{1k}(x)] \rightarrow 0 \quad \text{for all } x \in R_0.$$

For fixed sequences $\{\hat{\tau}_{1k}(x)\}$ and $\{\hat{\theta}_k\}$ satisfying (4.11) and $\hat{\theta}_k - \theta \rightarrow 0$, select $\varepsilon > 0$ such that $\theta \pm \varepsilon \in R_0$ and $\theta - \varepsilon \leq \hat{\theta}_k \leq \theta + \varepsilon$ for large k . Since $\hat{\tau}_{1k}(x)$ is nondecreasing in x , $\hat{\tau}_{1k}(\theta - \varepsilon) \leq \hat{\tau}_{1k}(\hat{\theta}_k) \leq \hat{\tau}_{1k}(\theta + \varepsilon)$ and

$$\begin{aligned} & \hat{\tau}_{1k}(\theta - \varepsilon) - \tau_{1k}(\theta - \varepsilon) + \tau_{1k}(\theta - \varepsilon) - \tau_{1k}(\theta) \\ & \leq \hat{\tau}_{1k}(\hat{\theta}_k) - \tau_{1k}(\theta) \\ & \leq \hat{\tau}_{1k}(\theta + \varepsilon) - \tau_{1k}(\theta + \varepsilon) + \tau_{1k}(\theta + \varepsilon) - \tau_{1k}(\theta). \end{aligned}$$

Letting $k \rightarrow \infty$, we have

$$\begin{aligned} \lim_{k \rightarrow \infty} n[\tau_{1k}(\theta - \varepsilon) - \tau_{1k}(\theta)] & \leq \liminf_{k \rightarrow \infty} n[\hat{\tau}_{1k}(\hat{\theta}_k) - \tau_{1k}(\theta)] \\ (4.12) \quad & \leq \limsup_{k \rightarrow \infty} n[\hat{\tau}_{1k}(\hat{\theta}_k) - \tau_{1k}(\theta)] \\ & \leq \lim_{k \rightarrow \infty} n[\tau_{1k}(\theta + \varepsilon) - \tau_{1k}(\theta)]. \end{aligned}$$

From the definition of $\tau_{1k}(x)$ and (A1),

$$\begin{aligned} & n[\tau_{1k}(\theta + \varepsilon) - \tau_{1k}(\theta)] \\ & = E \sum_{h=1}^L \frac{n(1-f_h)n_h}{n_h - 1} \sum_{i=1}^{n_h} [a_{hi}(\theta + \varepsilon) - a_{hi}(\theta)] [a_{hi}(\theta + \varepsilon) + a_{hi}(\theta)] \\ & \leq 4n \max_{j \leq N_{hi}, i \leq N_h, h \leq L} n_{hi} w_{hij} E \sum_{h=1}^L \sum_{i=1}^{n_h} \sum_{j=1}^{n_{hi}} w_{hij} I(\theta < y_{hij} \leq \theta + \varepsilon) \\ & = O(1)E[G_k(\theta + \varepsilon) - G_k(\theta)] = O(1)[F_k(\theta + \varepsilon) - F_k(\theta)]. \end{aligned}$$

Hence by (A9), $\lim_{\varepsilon \rightarrow 0} \lim_{k \rightarrow \infty} n[\tau_{1k}(\theta + \varepsilon) - \tau_{1k}(\theta)] = 0$. This result still holds if ε is replaced by $-\varepsilon$. Letting $\varepsilon \rightarrow 0$ in (4.12), we have $\lim_{k \rightarrow \infty} n[\hat{\tau}_{1k}(\hat{\theta}_k) - \tau_{1k}(\theta_k)] = 0$. This (and the subsequence argument) proves that

$$n[\hat{\tau}_{1k}(\hat{\theta}_k) - \tau_{1k}(\theta_k)] \rightarrow_p 0.$$

Similarly, we can show that, for $l = 2, 3$ and 4 ,

$$n[\hat{\tau}_{lk}(\hat{\theta}_k) - \tau_{lk}(\theta_k)] \rightarrow_p 0 \quad \text{and} \quad \hat{F}_k(\hat{\theta}_k) - F_k(\theta_k) \rightarrow_p 0.$$

Hence (4.9) follows from

$$\hat{v}_k = \left[\hat{F}_k^2(\hat{\theta}_k) \widehat{\text{var}}(\infty) - 2\hat{F}_k(\hat{\theta}_k)[\hat{\tau}_{3k}(\hat{\theta}_k) - \hat{\tau}_{4k}(\hat{\theta}_k)] + \hat{\tau}_{1k}(\hat{\theta}_k) - \hat{\tau}_{2k}(\hat{\theta}_k) \right] / q_k^2$$

and

$$v_k(\theta_k) = F_k^2(\theta_k) \text{Var}(q_k) - 2F_k(\theta_k)[\tau_{3k}(\theta_k) - \tau_{4k}(\theta_k)] + \tau_{1k}(\theta_k) - \tau_{2k}(\theta_k).$$

This completes the proof. \square

From results (4.4) and (4.10), an approximate $(1 - 2\alpha)$ confidence interval for θ_k is

$$CI_k = [\hat{\theta}_k - z_\alpha \hat{\sigma}_k, \hat{\theta}_k + z_\alpha \hat{\sigma}_k],$$

where z_α is the $(1 - \alpha)$ quantile of the standard normal distribution (or the t -distribution with $n - 1$ degrees of freedom). Another confidence interval is Woodruff's (1952) interval obtained by inverting the sample distribution function. Let $s_k = [v_k(\theta_k)]^{1/2}$. If s_k were known, then by (4.3), an approximate $(1 - 2\alpha)$ confidence interval for $\hat{\theta}_k$ is

$$CI_k^* = [\hat{F}_k^{-1}(p - z_\alpha s_k), \hat{F}_k^{-1}(p + z_\alpha s_k)].$$

The "confidence interval" CI_k^* is better than CI_k since it avoids the estimation of $1/f_k(\theta_k)$. Since s_k is unknown, it has to be replaced by $\hat{s}_k = (\hat{v}_k)^{1/2}$, and the resulting confidence interval is Woodruff's interval:

$$CI_k^W = [\hat{F}_k^{-1}(p - z_\alpha \hat{s}_k), \hat{F}_k^{-1}(p + z_\alpha \hat{s}_k)].$$

Asymptotic validity of CI_k^W was first shown by Francisco and Fuller (1991), but they assumed some unnecessary conditions [Conditions 5–7 in Francisco and Fuller (1991)]. We now prove the same result with the conditions in Theorem 5.

THEOREM 6. *Suppose that assumptions (A1) and (A6)–(A9) hold. Then*

$$(4.13) \quad P\{\theta_k \in CI_k^W\} \rightarrow 1 - 2\alpha.$$

PROOF. From (4.9), $\hat{s}_k/s_k \rightarrow_p 1$. Let $\varepsilon > 0$ be arbitrarily given. Then

$$(4.14) \quad P\{(1 - \varepsilon)s_k \leq \hat{s}_k \leq (1 + \varepsilon)s_k\} \rightarrow 1.$$

Because (4.14), we may assume $(1 - \varepsilon)s_k \leq \hat{s}_k \leq (1 + \varepsilon)s_k$ in the following proof. Note that

$$\begin{aligned} & P\left\{\theta_k \in \left[\hat{F}_k^{-1}(p - z_\alpha(1 - \varepsilon)s_k), \hat{F}_k^{-1}(p + z_\alpha(1 - \varepsilon)s_k)\right]\right\} \\ & \leq P\{\theta_k \in CI_k^W\} \\ & \leq P\left\{\theta_k \in \left[\hat{F}_k^{-1}(p - z_\alpha(1 + \varepsilon)s_k), \hat{F}_k^{-1}(p + z_\alpha(1 + \varepsilon)s_k)\right]\right\} \end{aligned}$$

and, by (4.3),

$$\begin{aligned} & P\left\{\theta_k \in \left[\hat{F}_k^{-1}(p - z_\alpha(1 \pm \varepsilon)s_k), \hat{F}_k^{-1}(p + z_\alpha(1 \pm \varepsilon)s_k)\right]\right\} \\ & = P\left\{\theta_k \leq \hat{F}_k^{-1}(p + z_\alpha(1 \pm \varepsilon)s_k)\right\} - P\left\{\theta_k < \hat{F}_k^{-1}(p - z_\alpha(1 \pm \varepsilon)s_k)\right\} \\ & = P\left\{\theta_k \leq \hat{\theta}_k + z_\alpha(1 \pm \varepsilon)s_k / [f_k(\theta_k)n^{1/2}] + o_p(n^{-1/2})\right\} \\ & \quad - P\left\{\theta_k < \hat{\theta}_k - z_\alpha(1 \pm \varepsilon)s_k / [f_k(\theta_k)n^{1/2}] + o_p(n^{-1/2})\right\} \\ & \rightarrow 2\Phi(z_\alpha(1 \pm \varepsilon)) - 1, \end{aligned}$$

where Φ is the standard normal distribution function. Hence

$$\begin{aligned} 2\Phi(z_\alpha(1-\varepsilon)) - 1 &\leq \liminf_{k \rightarrow \infty} P\{\theta_k \in CI_k^W\} \\ &\leq \limsup_{k \rightarrow \infty} P\{\theta_k \in CI_k^W\} \\ &\leq 2\Phi(z_\alpha(1+\varepsilon)) - 1. \end{aligned}$$

The result follows since ε is arbitrary. \square

The variance estimator proposed by Francisco and Fuller (1991) is

$$(4.15) \quad \tilde{\sigma}_k^2 = \hat{v}_k \tilde{u}_k^2,$$

where

$$(4.16) \quad \tilde{u}_k = (2z_\alpha \hat{s}_k)^{-1} [F_k^{-1}(p + z_\alpha \hat{s}_k) - F_k^{-1}(p - z_\alpha \hat{s}_k)].$$

This estimator is motivated by Woodruff's interval CI_k^W . Comparing (4.8) with (4.16), and using (4.3), we can show that $\tilde{\sigma}_k^2/\hat{\sigma}_k^2 \rightarrow_p 1$.

Acknowledgments. I would like to thank Professor J. N. K. Rao at Carleton University, two referees and an Associate Editor for their helpful comments.

REFERENCES

- BEACH, C. M. and KALISKI, S. F. (1986). Lorenz curve inference with sample weights: an application to the distribution of unemployment experience. *J. Roy. Statist. Soc. Ser. C* **35** 38–45.
- BICKEL, P. J. (1973). On some analogues to linear combinations of order statistics in the linear model. *Ann. Statist.* **1** 597–616.
- BICKEL, P. J. and FREEDMAN, D. A. (1984). Asymptotic normality and the bootstrap in stratified sampling. *Ann. Statist.* **12** 470–482.
- CHERNOFF, H., GASTWIRTH, J. L. and JOHNS, M. V., JR. (1967). Asymptotic distribution of linear combinations of order statistics, with applications to estimation. *Ann. Math. Statist.* **38** 52–72.
- EFRON, B. (1979). Bootstrap methods: Another look at the jackknife. *Ann. Statist.* **7** 1–26.
- FRANCISCO, C. A. and FULLER, W. A. (1991). Quantile estimation with a complex survey design. *Ann. Statist.* **19** 454–469.
- KISH, L. and FRANKEL, M. R. (1974). Inference from complex samples (with discussion). *J. Roy. Statist. Soc. Ser. B* **36** 1–37.
- KOENKER, R. and BASSETT, G. (1978). Regression quantiles. *Econometrica* **46** 33–50.
- KOENKER, R. and PORTNOY, S. (1987). L-estimation for linear models. *J. Amer. Statist. Assoc.* **82** 851–857.
- KREWSKI, D. and RAO, J. N. K. (1981). Inference from stratified samples: Properties of the linearization, jackknife and balanced repeated replication methods. *Ann. Statist.* **9** 1010–1019.
- MCCARTHY, P. J. (1969). Pseudo-replication: half samples. *Internat. Statist. Rev.* **37** 239–264.
- NYGARD, F. and SANDSTRÖM, A. (1985). The estimation of Gini and the entropy inequality parameters in finite populations. *Journal of Official Statistics* **1** 399–412.
- RAO, J. N. K. and WU, C. F. J. (1985). Inference from stratified samples: second-order analysis of three methods for nonlinear statistics. *J. Amer. Statist. Assoc.* **80** 620–630.

- RUPPERT, D. and CARROLL, R. J. (1980). Trimmed least squares estimation in the linear model. *J. Amer. Statist. Assoc.* **75** 828–838.
- SERFLING, R. J. (1980). *Approximation Theorems of Mathematical Statistics*. Wiley, New York.
- SHAO, J. and WU, C. F. J. (1989). A general theory for jackknife variance estimation. *Ann. Statist.* **17** 1176–1197.
- SHAO, J. and WU, C. F. J. (1992). Asymptotic properties of the balanced repeated replication method for sample quantiles. *Ann. Statist.* **20** 1571–1593.
- STIGLER, S. M. (1973). The asymptotic distribution of the trimmed mean. *Ann. Statist.* **1** 472–477.
- TUKEY, J. (1958). Bias and confidence in not quite large samples (abstract). *Ann. Math. Statist.* **29** 614.
- WELSH, A. H. (1987). The trimmed mean in the linear model (with discussion). *Ann. Statist.* **15** 20–45.
- WOODRUFF, R. S. (1952). Confidence intervals for medians and other position measures. *J. Amer. Statist. Assoc.* **47** 635–646.

DEPARTMENT OF STATISTICS
UNIVERSITY OF WISCONSIN
MADISON, WISCONSIN 53706-1885