

Estimating the distance distribution of subpopulations for a large-scale complex survey

Jianqiang Wang
Department of Statistics
Iowa State University

Jean Opsomer
Professor, Department of Statistics
Iowa State University

Abstract

Many finite populations targeted by sample surveys consist of homogenous subpopulations with respect to the variables being collected. In on-going survey operations, it is often of interest to be able to assess whether a new observation belongs to one of those subpopulations or should be flagged as not belonging to any of them. For this purpose, we propose a sample-based estimator for the subpopulation distribution functions of the distances between the elements and the subpopulation centers. We explore different ways to define the subpopulation centers and several distance metrics. We describe the theoretical properties of the estimator, and propose several approaches for design-based variance estimation. The practical properties of the procedures are evaluating in a simulation study.

Keywords: jackknife, kernel estimation, estimating equations, generalized median, elliptical distribution, classification.

1. Introduction

A common issue in large-scale complex surveys is the detection of outliers in the data. Such outliers can be caused by frame imperfections, which can lead to ineligible units being selected, or by errors during data collection. If these outliers remain in the survey dataset, they can cause inference based on the survey to be invalid for the population of interest. Most survey operations therefore incorporate data editing and validation as part of the post-data collection steps, where they attempt to identify suspicious observations and either remove or correct them. When outliers exhibit “extreme” values on one or several survey variables, they can be detected relatively easily. However, because surveys often collect large numbers of variables, there is the potential for other outliers which are not extreme on any single variable. Detecting such outliers is more difficult, and identifying unusual or suspicious patterns in the data often requires substantial subject-matter knowledge.

In practice, many finite populations targeted by surveys consist of a number of relatively homogeneous subpopulations, and this structure can be exploited to develop a statistical approach for flagging suspicious observations that does not require detailed knowledge of the relationships between the variables. The type of surveys we are targeting here is one in which recurring surveys are

made over time, and the characteristics and composition of the subpopulations remains relatively stable between surveys. The focus of the current article is to propose an estimator for the “outlyingness” of an observation relative to a subpopulation, based on the distance between an observation and the center of the subpopulation, and to derive its statistical properties in an asymptotic design-based context.

The remainder of the paper is as follows. Section 2 defines notation and establishes preliminary results for design-based inference. Section 3 lists our general design assumptions and assumptions on the sequence of finite populations. Section 4 presents our theoretical results showing asymptotic properties of distance distribution functions using either means or medians as subpopulation centers. Simulation results are detailed in Wang (2008) due to space limits.

2. Notation, assumption and preliminary results

2.1 Notation and definitions

Suppose we have an increasing sequence of finite populations $(U_\nu)_{\nu=1}^\infty$ of sizes N_ν with $N_\nu < N_{\nu+1}$. Associated with the i -th population element is a p -dimensional vector of observations

$$\mathbf{y}_i = (y_{i,1}, \dots, y_{i,p}),$$

and let \mathcal{F}_ν be the power set of ν -th finite population $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{N_\nu}\}$. The finite population U_ν can be partitioned into G subpopulations, with $U_\nu = \bigcup_{g=1}^G U_{\nu g}$. The knowledge of this structure of finite population is usually available for longitudinal surveys of stable populations, and we assume the partition into subpopulations is provided.

We take a sample $\mathcal{S}_\nu = \bigcup_{g=1}^G \mathcal{S}_{\nu g}$ of size n_ν from population U_ν , and the sampling design may be a complex design with stratification or multi-stage sampling. Here, $\mathcal{S}_{\nu g}$ are mutually exclusive subsets of \mathcal{S}_ν which contain elements from subpopulation g only. We assume the inclusion probability of the i -th population element, is known without error,

$$P(i \in \mathcal{S}_\nu) = \pi_i,$$

and the sample size

$$n_\nu = \sum_{g=1}^G n_{\nu g},$$

where $n_{\nu g}$ is the number of elements in sample \mathcal{S}_ν that come from subpopulation g . We use $n_\nu^* = E(n_\nu | \mathcal{F}_\nu)$ to denote the expected sample size conditioning on finite population and $n_{\nu g}^*$ to denote the expected sample size for subpopulation g .

We define the population level distribution of distances as

$$D_{\nu g, d}(\gamma_{\nu g}) = \frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \mathbf{I}(\|\mathbf{y}_i - \gamma_{\nu g}\| \leq d), \quad (1)$$

where $\gamma_{\nu g}$ is some measure of center of population $U_{\nu g}$. Given a sample \mathcal{S}_ν , $D_{\nu g, d}(\gamma_{\nu g})$ is estimated by,

$$\widehat{D}_{\nu g, d}(\hat{\gamma}_{\nu g}) = \frac{1}{\widehat{N}_{\nu g}} \sum_{\mathcal{S}_{\nu g}} \frac{1}{\pi_i} \mathbf{I}(\|\mathbf{y}_i - \hat{\gamma}_{\nu g}\| \leq d), \quad (2)$$

or

$$\widetilde{D}_{\nu g, d}(\hat{\gamma}_{\nu g}) = \frac{1}{N_{\nu g}} \sum_{\mathcal{S}_{\nu g}} \frac{1}{\pi_i} \mathbf{I}(\|\mathbf{y}_i - \hat{\gamma}_{\nu g}\| \leq d), \quad (3)$$

where $\hat{\gamma}_{\nu g}$ is an estimator of $\gamma_{\nu g}$, some measure of the center of subpopulation g formed from sample data $\mathcal{S}_{\nu g}$. Here, $\widehat{N}_{\nu g} = \sum_{i \in \mathcal{S}_{\nu g}} \frac{1}{\pi_i}$ is an estimator of size of subpopulation g which is generally not known a priori.

In equation (1), we assume $\gamma_{\nu g}$ is a nonrandom sequence of finite population centers but $\hat{\gamma}_{\nu g}$ is random due to sampling mechanism. Quantities $D_{\nu g, d}(\gamma_{\nu g})$ and $\widehat{D}_{\nu g, d}(\hat{\gamma}_{\nu g})$ are step functions of d and γ with jumps of size $O(\frac{1}{N_{\nu g}})$ and $O(\frac{1}{n_{\nu g}^*})$.

We can use the usual mean vector as a measure of center in equations (1)-(3),

$$\boldsymbol{\mu}_{\nu g} = \frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \mathbf{y}_i, \quad (4)$$

which is estimated by,

$$\hat{\boldsymbol{\mu}}_{\nu g} = \frac{1}{\widehat{N}_{\nu g}} \sum_{\mathcal{S}_{\nu g}} \frac{\mathbf{y}_i}{\pi_i}. \quad (5)$$

We can also use a generalized version of median as a measure of center in multivariate case. The definition of spatial median was given in Brown (1983) and Small (1990). We generalize this idea and define the multivariate median of a finite population to be the location with smallest overall distance to all population units with respect to some norm $\|\cdot\|$. The norm is chosen at our choice, and commonly we use Minkoski distance or some shape-respecting quadratic distance.

$$\mathbf{q}_{\nu g} = \arg \inf_{\boldsymbol{\gamma}} \sum_{U_{\nu g}} \|\mathbf{y}_i - \boldsymbol{\gamma}\|. \quad (6)$$

The sample-based estimator of $\mathbf{q}_{\nu g}$ is given as follows,

$$\hat{\mathbf{q}}_{\nu g} = \arg \inf_{\boldsymbol{\gamma}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{1}{\pi_i} \|\mathbf{y}_i - \boldsymbol{\gamma}\|. \quad (7)$$

We use $\boldsymbol{\gamma}_{\nu g}$ to denote a general measure of center, which can be mean vector $\boldsymbol{\mu}_{\nu g}$ or median $\mathbf{q}_{\nu g}$, and $\hat{\boldsymbol{\gamma}}_{\nu g}$ denotes the estimator of $\boldsymbol{\gamma}_{\nu g}$.

For now, we assume that the subpopulation association for each \mathbf{y}_i in the sample is known.

3. Assumptions

3.1 General design assumptions

In this paper, we estimate the distribution of subpopulation distances in design-based framework. We assume the sequence of finite populations to be fixed and randomness only comes from the sampling mechanism. We do not want to restrict our attention to a specific sampling design but make rather general assumptions to cover various sampling schemes. Assumptions 3.1.1 and 3.1.2 ensure the design consistency and asymptotic normality of our estimator.

Assumption 3.1.1. *The following conditions hold for population size, inclusion probabilities π_i and design variance of Horvitz-Thompson estimator of the mean,*

1. $N_{\nu g} = O(N_\nu)$, $n_\nu = O_p(N_\nu^\beta)$ with $\beta \in (\frac{1}{2}, 1]$.
2. $K_L \leq \frac{N_\nu}{n_\nu} \pi_i \leq K_U$ for all i , where K_L and K_U are positive constants.
3. For any vector \mathbf{z} with finite $2 + \delta$ moments, define $\bar{\mathbf{z}}_{\nu, \pi} = \frac{1}{N_\nu} \sum_{\mathcal{S}_\nu} \frac{\mathbf{z}_i}{\pi_i}$ as the Horvitz-Thompson estimator of $\bar{\mathbf{z}}_\nu = \frac{1}{N_\nu} \sum_{U_\nu} \mathbf{z}_i$. We assume

$$\text{Var}(\bar{\mathbf{z}}_{\nu, \pi} | \mathcal{F}_\nu) \leq K_1 \text{Var}_{SRS}(\bar{\mathbf{z}}_{\nu, \pi} | \mathcal{F}_\nu),$$

for some constant K_1 , where $\text{Var}_{SRS}(\bar{\mathbf{z}}_{\nu, \pi} | \mathcal{F}_\nu)$ is the design variance-covariance matrix of $\bar{\mathbf{z}}_{\nu, \pi}$ under simple random sampling of size n_ν^* .

It can be shown that under Assumption 3.1.1(3), $\frac{n_{\nu g}}{n_\nu^*} \xrightarrow{P} 1$ by bounding its design variance.

We will make the following normality assumption on Horvitz-Thompson estimator for a general vector with moment conditions, similar to Fuller (2007).

Assumption 3.1.2. *For any \mathbf{z} with positive variance-covariance matrix and finite fourth moment*

$$n_\nu^{1/2}(\bar{\mathbf{z}}_{\nu, \pi} - \bar{\mathbf{z}}_\nu) | \mathcal{F}_\nu \xrightarrow{d} N(\mathbf{0}, \Sigma_{\nu, \mathbf{z}}), \quad (8)$$

and

$$[V(\bar{\mathbf{z}}_{\nu, \pi} | \mathcal{F}_\nu)]^{-1} \hat{V}_{HT}\{\bar{\mathbf{z}}_{\nu, \pi}\} - \mathbf{I}_{p \times p} = O_p(n_\nu^{*-1/2}), \quad (9)$$

where $\bar{\mathbf{z}}_\nu = \frac{1}{N_\nu} \sum_{i=1}^{N_\nu} \mathbf{z}_i$, $\bar{\mathbf{z}}_{\nu, \pi} = \frac{1}{N_\nu} \sum_{\mathcal{S}_\nu} \frac{\mathbf{z}_i}{\pi_i}$, and $\hat{V}_{HT}\{\bar{\mathbf{z}}_{\nu, \pi}\}$ is the Horvitz-Thompson estimator of the variance of $\bar{\mathbf{z}}_{\nu, \pi} | \mathcal{F}_\nu$.

3.2 Assumptions on finite population

To show the design properties of the distance distribution estimator, we need to assume a number of regularity conditions on the sequence of finite populations.

Assumption 3.2.1. *The sequence of population vectors \mathbf{y}_i 's in subpopulation g has bounded $4 + \delta$ moments,*

$$\lim_{N_{\nu g} \rightarrow \infty} N_{\nu g}^{-1} \sum_{i \in U_{\nu g}} |\mathbf{y}_i|^{4+\delta} < \infty,$$

for some $\delta > 0$.

Assumption 3.2.2. *1. The limit of population level distance distribution exists,*

$$\lim_{\nu \rightarrow \infty} D_{\nu g, d}(\boldsymbol{\gamma}) = \mathcal{D}_{g, d}(\boldsymbol{\gamma})$$

on $(d, \boldsymbol{\gamma}) \in [0, \infty) \times \mathbb{R}^p$.

2. The limiting function $\mathcal{D}_{g, d}(\boldsymbol{\gamma})$ is continuous in $d \in [0, \infty)$ and $\boldsymbol{\gamma} \in \mathbb{R}^p$. Additionally, the derivatives $\frac{\partial \mathcal{D}_{g, d}(\boldsymbol{\gamma})}{\partial d}$, $\frac{\partial \mathcal{D}_{g, d}(\boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}}$ and $\frac{\partial^2 \mathcal{D}_{g, d}(\boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}^2}$ all exist and are finite in $(d, \boldsymbol{\gamma}) \in [0, +\infty) \times \mathbb{R}^p$.

Assumption 3.2.3. *The following population quantities converge to zero:*

1.

$$\sqrt{N_{\nu g}} \left\{ \frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \mathbf{I}_{(d < \|\mathbf{y}_i - \boldsymbol{\gamma}\| \leq d + h_{N_{\nu g}})} - \frac{\partial \mathcal{D}_{g, d}(\boldsymbol{\gamma})}{\partial d} h_{N_{\nu g}} \right\}$$

converges to zero, where $h_{N_{\nu g}} = O(N_{\nu g}^{-\alpha})$ and $\alpha \in (\frac{1}{4}, 1)$.

2.

$$\frac{n_{\nu g}^* 1/2}{N_{\nu g}} \sum_{U_{\nu g}} \left[\mathbf{I}_{(\|\mathbf{y}_i - \boldsymbol{\gamma} - n_{\nu g}^{-1/2} \mathbf{s}\| \leq d)} - \mathbf{I}_{(\|\mathbf{y}_i - \boldsymbol{\gamma}\| \leq d)} - \mathcal{D}_{g, d}(\boldsymbol{\gamma} + n_{\nu g}^* 1/2 \mathbf{s}) + \mathcal{D}_{g, d}(\boldsymbol{\gamma}) \right]$$

converges to zero uniformly for $\boldsymbol{\gamma} \in \mathbb{R}^p$ and $\mathbf{s} \in C_{\mathbf{s}}$, a large enough compact set in \mathbb{R}^p .

Assumptions 3.2.4-3.2.8 are necessary when deriving median-based results, where we need stronger conditions on the spreading of finite population elements and restrictions on the norm. Assumption 3.2.4 guarantees the uniqueness of population median $\mathbf{q}_{\nu g}$. Assumption 3.2.6 restricts our concern to an inner product space, and Assumption 3.2.7 will be used in showing asymptotic normality of $\hat{\mathbf{q}}_{\nu g}$.

Assumption 3.2.4. *The finite population $U_{\nu g}$ only puts finitely many points \mathbf{y}_i on a line in \mathbb{R}^p .*

Assumption 3.2.5. *The norm $\|\cdot\|$ is continuous on \mathbb{R}^p , with a continuous gradient vector $\boldsymbol{\psi}(\boldsymbol{\gamma})$, and bounded second derivative matrix $H_s(\boldsymbol{\gamma})$.*

Assumption 3.2.6. *The norm $\|\cdot\|$ has the following inner product representation*

$$\|\boldsymbol{\gamma}\| = \sqrt{\langle \boldsymbol{\gamma}, \boldsymbol{\gamma} \rangle}$$

for $\boldsymbol{\gamma} \in \mathbb{R}^p$ and some inner product $\langle \cdot, \cdot \rangle$.

Assumption 3.2.7. *For any $\boldsymbol{\gamma}$ in a neighborhood of $\mathbf{q}_{\nu g}$, $\frac{1}{N_{\nu g}} \sum_{U_{\nu g}} H_s(\mathbf{y}_i - \boldsymbol{\gamma})$ is a nonsingular matrix. Further, the sequence of $H_s(\mathbf{y}_i - \boldsymbol{\gamma}_{\nu g})$ has bounded first two moments at population level.*

Assumption 3.2.8. *Assume that the linearized term $\boldsymbol{\psi}(\mathbf{y}_i - \boldsymbol{\gamma}_{\nu g})$ has bounded fourth population moments,*

$$\frac{1}{N_{\nu g}} \sum_{U_{\nu g}} |\boldsymbol{\psi}(\mathbf{y}_i - \boldsymbol{\gamma}_{\nu g})|^4 < \infty.$$

We use a kernel smoothing estimator to estimate the derivative of limiting smooth function as needed in the expression of asymptotic variance. Assumptions 3.2.9 and 3.2.10 give conditions on the choice of kernel function and bandwidth.

Assumption 3.2.9. *The kernel function $K(t)$ is symmetric with $\int_{-\infty}^{\infty} K(t) dt = 1$, and $K(t)$ is an absolutely continuous function with finite derivatives $K'(t)$. Further, let $R(K) = \int_{-\infty}^{\infty} K^2(t) dt < \infty$ and $\sigma_K^2 = \int_{-\infty}^{\infty} t^2 K(t) dt < \infty$.*

Assumption 3.2.10. *The smoothing bandwidth $h \rightarrow 0$, and $N_{\nu g} h (\log N_{\nu g})^{-1} \rightarrow \infty$, as $N_{\nu g} \rightarrow \infty$.*

Assumption 3.2.11. *There exists a constant c , such that $|\frac{1}{h^2} K'(\frac{x}{h})| \leq c$, for any $x \neq 0$ and h arbitrarily small.*

4. Main results

4.1 Mean-based inference

Lemma 1. *Under Assumptions 3.1.1 and 3.2.2-3.2.3,*

$$n_{\nu g}^* 1/2 \left(\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) - \widehat{D}_{\nu g, d}(\boldsymbol{\mu}_{\nu g}) - \mathcal{D}_{g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) + \mathcal{D}_{g, d}(\boldsymbol{\mu}_{\nu g}) \right) \tag{10}$$

converges to zero in design.

For a proof of this result, see Wang (2008).

Theorem 1. *Under Assumptions 3.1.1 and 3.2.2-3.2.3, the sample-based quantity $\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g})$ is $\sqrt{n_{\nu g}^*}$ -consistent for the corresponding population quantity $D_{\nu g, d}(\boldsymbol{\mu}_{\nu g})$, namely,*

$$n_{\nu g}^* 1/2 \left(\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) - D_{\nu g, d}(\boldsymbol{\mu}_{\nu g}) \right) = O_p(1).$$

Proof. We use the following decomposition,

$$\begin{aligned} & n_{\nu g}^* 1/2 \left(\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) - D_{\nu g, d}(\boldsymbol{\mu}_{\nu g}) \right) \\ &= n_{\nu g}^* 1/2 \left(\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) - \widehat{D}_{\nu g, d}(\boldsymbol{\mu}_{\nu g}) - \mathcal{D}_{g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) + \mathcal{D}_{g, d}(\boldsymbol{\mu}_{\nu g}) \right) \\ &+ n_{\nu g}^* 1/2 \left(\widehat{D}_{\nu g, d}(\boldsymbol{\mu}_{\nu g}) - D_{\nu g, d}(\boldsymbol{\mu}_{\nu g}) \right) \\ &+ n_{\nu g}^* 1/2 \left(\mathcal{D}_{g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) - \mathcal{D}_{g, d}(\boldsymbol{\mu}_{\nu g}) \right), \end{aligned} \tag{11}$$

where the first term is $o_p(1)$ by Lemma 1, and the last two terms are both $O_p(1)$. \square

Expression (8) of Assumption 3.1.2 implies the following multivariate normality,

$$\frac{n_{\nu g}^*}{N_{\nu g}} \sum_{U_{\nu g}} \underbrace{\begin{bmatrix} \mathbb{I}(\|\mathbf{y}_i - \boldsymbol{\mu}_{\nu g}\| \leq d) \\ 1 \\ \mathbf{y}_i \end{bmatrix}}_{\mathbf{b}_{\boldsymbol{\mu}, gi}} \left[\frac{\mathbb{I}(i \in \mathcal{S}_{\nu g})}{\pi_i} - 1 \right] \Bigg| \mathcal{F}_{\nu} \xrightarrow{d} N(0, \Sigma_{\boldsymbol{\mu}, d}), \quad (12)$$

where

$$\Sigma_{\boldsymbol{\mu}, d} = \frac{n_{\nu g}^*}{N_{\nu g}^2} \sum_{U_{\nu g}} \sum_{U_{\nu g}} (\pi_{ij} - \pi_i \pi_j) \frac{\mathbf{b}_{\boldsymbol{\mu}, gi} \mathbf{b}'_{\boldsymbol{\mu}, gj}}{\pi_i \pi_j}, \quad (13)$$

as we are estimating a domain quantity and $\mathbf{b}_{\boldsymbol{\mu}, gi}$ has bounded second moments under Assumption 3.2.1.

Theorem 2. *Under Assumptions 3.1.1-3.1.2 and 3.2.1-3.2.3,*

$$n_{\nu g}^{1/2} \left[V \left(\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) \right) \right]^{-1/2} \left(\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) - D_{\nu g, d}(\boldsymbol{\mu}_{\nu g}) \right) \xrightarrow{d} N(0, 1),$$

where

$$V \left(\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g}) \right) = \boldsymbol{\alpha}'_{\boldsymbol{\mu}} \Sigma_{\boldsymbol{\mu}, d} \boldsymbol{\alpha}_{\boldsymbol{\mu}}, \quad (14)$$

$$\boldsymbol{\alpha}_{\boldsymbol{\mu}} = \left(1, -D_{\nu g, d}(\boldsymbol{\mu}_{\nu g}) - \left(\frac{\partial \mathcal{D}_{g, d}(\boldsymbol{\mu}_{\nu g})}{\partial \boldsymbol{\mu}_{\nu g}} \right)' \boldsymbol{\mu}_{\nu g}, \left(\frac{\partial \mathcal{D}_{g, d}(\boldsymbol{\mu}_{\nu g})}{\partial \boldsymbol{\mu}_{\nu g}} \right)' \right)' \quad (15)$$

and $\Sigma_{\boldsymbol{\mu}, d}$ is defined in (13).

Proof. We still use decomposition (11), then it follows from Lemma 1, expression (4.1), Taylor Linearization and Slutsky's Theorem. \square

The asymptotic variance of $\widehat{D}_{\nu g, d}(\hat{\boldsymbol{\mu}}_{\nu g})$ consists of two pieces, the piece from estimating the distribution function with known center and the piece due to the uncertainty of estimating population center. The first piece of variance can be easily estimated using plug-in estimator or replication procedure, but the second piece involves an unknown derivative of limiting smooth function. The derivative can be estimated by kernel smoothing and incorporated into either plug-in or replication estimator.

Remarks: In Lemma 1, we use the limiting function $\mathcal{D}_{g, d}(\cdot)$ instead of population quantity $D_{\nu g, d}(\cdot)$, because we want the third term in (11) to be a smooth function so we can use linearization to quantify the randomness due to estimating subpopulation center.

4.2 Median-based inference

Now let us look at the case when we use median as a measure of subpopulation center. We introduce the following estimating equations at population and sample level, respectively,

$$\sum_{U_{\nu g}} \boldsymbol{\psi}(\mathbf{y}_i - \boldsymbol{\gamma}) = \mathbf{0}, \quad (16)$$

and

$$\sum_{i \in \mathcal{S}_{\nu g}} \frac{\boldsymbol{\psi}(\mathbf{y}_i - \boldsymbol{\gamma})}{\pi_i} = \mathbf{0}. \quad (17)$$

The medians defined by (6)-(7) are related to the roots of the estimating equations (16)-(17). More details are held until Appendix B.

Lemma 2. *Under Assumptions 3.2.4 and 3.2.6, for a large enough population, $\sum_{U_{\nu g}} \|\mathbf{y}_i - \boldsymbol{\gamma}\|$ has only one local minimum, which is also its global minimum.*

Remark 1: Lemma 2 states a stronger result than $\mathbf{q}_{\nu g}$ being unique, and it also says there are no other local minimums for $\sum_{U_{\nu g}} \|\mathbf{y}_i - \boldsymbol{\gamma}\|$.

Remark 2: If we make a similar assumption on the sample $\mathcal{S}_{\nu g}$, then it is obvious that $\sum_{\mathcal{S}_{\nu g}} \frac{\|\mathbf{y}_i - \boldsymbol{\gamma}\|}{\pi_i}$ has a unique global minimizer and no other local minimizers. But under a complex sampling design, we may have nonzero probability of selecting a sample where all points are on the same line. But considering the increasing sequence of finite populations and sequence of sampling designs, the probability of selecting a finite-sized sample will go to zero.

In establishing the weak convergence of $\hat{\mathbf{q}}_{\nu g}$ for $\mathbf{q}_{\nu g}$, we adopt the definition of weak convergence from P.24 of Billingsley (1968) and use general norm $\|\cdot\|$ as a discrepancy measure.

Theorem 3. *Under Assumptions 3.1.1, 3.2.4 and 3.2.6, any sequence $\hat{\mathbf{q}}_{\nu g}$ that satisfies (7) is design consistent for $\mathbf{q}_{\nu g}$.*

Proof. Negation. Suppose $\hat{\mathbf{q}}_{\nu g}$ does not converge to $\mathbf{q}_{\nu g}$ in probability, then there exists $\epsilon_1, \delta > 0$, such that

$$P(\|\hat{\mathbf{q}}_{\nu g} - \mathbf{q}_{\nu g}\| > \epsilon_1) > \delta, \quad (18)$$

for all ν . Equation (18) together with Assumption 3.2.4 would imply $\exists \epsilon_2 > 0$, such that

$$P\left(\frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \|\mathbf{y}_i - \hat{\mathbf{q}}_{\nu g}\| - \frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \|\mathbf{y}_i - \mathbf{q}_{\nu g}\| > \epsilon_2\right) > \delta, \quad (19)$$

for all ν .

Now, let us show that

$$\frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\|\mathbf{y}_i - \hat{\mathbf{q}}_{\nu g}\|}{\pi_i} - \frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \|\mathbf{y}_i - \hat{\mathbf{q}}_{\nu g}\| \xrightarrow{p} 0. \quad (20)$$

Define

$$Q_n(\gamma) = \frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\|\mathbf{y}_i - \gamma\|}{\pi_i} - \frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \|\mathbf{y}_i - \gamma\|,$$

note that,

$$P(|Q_n(\hat{\mathbf{q}}_{\nu g})| > \epsilon) \leq P[\sup_{\gamma \in C} |Q_n(\gamma)| > \epsilon'] + P(\hat{\mathbf{q}}_{\nu g} \notin C),$$

where C is a large enough compact set. Now it is left to show that

$$\sup_{\gamma \in C} |Q_n(\gamma)| \xrightarrow{p} 0. \quad (21)$$

The equation above can be shown by covering technique. Equations (19) and (20) imply that, there exists $\epsilon_3 > 0$ and $\delta_1 > 0$, such that

$$P\left(\frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\|\mathbf{y}_i - \hat{\mathbf{q}}_{\nu g}\|}{\pi_i} - \frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \|\mathbf{y}_i - \mathbf{q}_{\nu g}\| > \epsilon_3\right) > \delta_1 \quad (22)$$

As $\frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\|\mathbf{y}_i - \mathbf{q}_{\nu g}\|}{\pi_i} \geq \frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\|\mathbf{y}_i - \hat{\mathbf{q}}_{\nu g}\|}{\pi_i}$, for the same ϵ_3 and δ_1 ,

$$P\left(\frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\|\mathbf{y}_i - \mathbf{q}_{\nu g}\|}{\pi_i} - \frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \|\mathbf{y}_i - \mathbf{q}_{\nu g}\| > \epsilon_3\right) > \delta_{\Sigma_{\nu g, \psi}} \quad (23)$$

as $\nu \rightarrow \infty$. Contradicting the fact that $\frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\|\mathbf{y}_i - \mathbf{q}_{\nu g}\|}{\pi_i}$ is design consistent for $\frac{1}{N_{\nu g}} \sum_{U_{\nu g}} \|\mathbf{y}_i - \mathbf{q}_{\nu g}\|$. \square

Theorem 4. Under Assumptions 3.1.1-3.1.2, 3.2.5 and 3.2.7, and let $\hat{\mathbf{q}}_{\nu g}$ be a design consistent sequence, then we have the following asymptotic normality for sample median,

$$n_{\nu g}^{*1/2}(\hat{\mathbf{q}}_{\nu g} - \mathbf{q}_{\nu g}) \xrightarrow{d} N_p(\mathbf{0}, \Sigma_{\nu g, \mathbf{q}}), \quad (23)$$

where

$$\Sigma_{\nu g, \mathbf{q}} = A \Sigma_{\nu g, \psi} A', \quad A = \left[\frac{1}{N_{\nu g}} \sum_{i=1}^{N_{\nu g}} H_s(\mathbf{y}_i - \mathbf{q}_{\nu g}) \right]^{-1}$$

and

$$\Sigma_{\nu g, \psi} = \frac{n_{\nu g}}{N_{\nu g}^2} \sum_{U_{\nu g}} \sum_{U_{\nu g}} (\pi_{ij} - \pi_i \pi_j) \frac{\psi(\mathbf{y}_i - \mathbf{q}_{\nu g}) \psi(\mathbf{y}_{jj} - \mathbf{q}_{\nu g})}{\pi_i \pi_j}.$$

And further,

$$n_{\nu g}^{*1/2}(\hat{\mathbf{q}}_{\nu g} - \mathbf{q}_{\nu g}) \xrightarrow{d} N_p(\mathbf{0}, \Sigma_{\nu g, \mathbf{q}}). \quad (24)$$

Proof. The consistency of $\hat{\mathbf{q}}_{\nu g}$ for $\mathbf{q}_{\nu g}$ gives,

$$\sum_{i \in \mathcal{S}_{\nu g}} \frac{1}{\pi_i} \psi(\mathbf{y}_i - \hat{\mathbf{q}}_{\nu g}) = \mathbf{0}$$

$$\Leftrightarrow \sum_{i \in \mathcal{S}_{\nu g}} \frac{1}{\pi_i} \psi(\mathbf{y}_i - \mathbf{q}_{\nu g} - (\hat{\mathbf{q}}_{\nu g} - \mathbf{q}_{\nu g})) = \mathbf{0}$$

$$\Leftrightarrow \sum_{i \in \mathcal{S}_{\nu g}} \frac{1}{\pi_i} \left\{ \psi(\mathbf{y}_i - \mathbf{q}_{\nu g}) - H_s(\mathbf{y}_i - \mathbf{q}_{\nu g})(\hat{\mathbf{q}}_{\nu g} - \mathbf{q}_{\nu g}) + o_p(\hat{\mathbf{q}}_{\nu g} - \mathbf{q}_{\nu g}) \right\} = \mathbf{0},$$

which implies

$$\begin{aligned} \hat{\mathbf{q}}_{\nu g} &= \mathbf{q}_{\nu g} \\ &+ \left[\frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{H_s(\mathbf{y}_i - \mathbf{q}_{\nu g})}{\pi_i} \right]^{-1} \frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\psi(\mathbf{y}_i - \mathbf{q}_{\nu g})}{\pi_i} \\ &+ o_p(\hat{\mathbf{q}}_{\nu g} - \mathbf{q}_{\nu g}) \end{aligned} \quad (25)$$

after using non-singularity condition in Assumption 3.2.7

to take the inverse of $\frac{1}{N_{\nu g}} \sum_{i=1}^{N_{\nu g}} H_s(\mathbf{y}_i - \mathbf{q}_{\nu g})$.

It is easy to argue that

$$\left[\frac{1}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{H_s(\mathbf{y}_i - \mathbf{q}_{\nu g})}{\pi_i} \right]^{-1} \xrightarrow{p} \left[\frac{1}{N_{\nu g}} \sum_{U_{\nu g}} H_s(\mathbf{y}_i - \mathbf{q}_{\nu g}) \right]^{-1}, \quad (26)$$

and

$$\frac{n_{\nu g}^{*1/2}}{N_{\nu g}} \sum_{i \in \mathcal{S}_{\nu g}} \frac{\psi(\mathbf{y}_i - \mathbf{q}_{\nu g})}{\pi_i} \xrightarrow{d} N_p(\mathbf{0}, \Sigma_{\nu g, \psi}), \quad (27)$$

where

$$\Sigma_{\nu g, \psi} = \frac{n_{\nu g}^*}{N_{\nu g}^2} \sum_{U_{\nu g}} \sum_{U_{\nu g}} (\pi_{ij} - \pi_i \pi_j) \frac{\psi(\mathbf{y}_i - \mathbf{q}_{\nu g}) \psi(\mathbf{y}_{jj} - \mathbf{q}_{\nu g})}{\pi_i \pi_j}.$$

The proof is then completed by applying Slutsky's Theorem to obtain (23). Then we can use Slutsky's Theorem again and the fact that $\frac{n_{\nu g}}{n_{\nu g}^*} \xrightarrow{p} 1$ to obtain (24). \square

Theorem 5. Suppose Assumptions 3.1.1 and 3.2.4-3.2.7 are satisfied, then for any sequence $\hat{\mathbf{q}}_{\nu g}$ that satisfies (7), the estimated distance distribution $\hat{D}_{\nu g, d}(\hat{\mathbf{q}}_{\nu g})$ is $\sqrt{n_{\nu g}^*}$ -consistent for the corresponding population quantity $D_{\nu g, d}(\mathbf{q}_{\nu g})$, namely,

$$n_{\nu g}^{*1/2} \left(\hat{D}_{\nu g, d}(\hat{\mathbf{q}}_{\nu g}) - D_{\nu g, d}(\mathbf{q}_{\nu g}) \right) = O_p(1).$$

Proof. Similar to the proof of Theorem 2, we use the following decomposition,

$$\begin{aligned} &n_{\nu g}^{*1/2} \left(\hat{D}_{\nu g, d}(\hat{\mathbf{q}}_{\nu g}) - D_{\nu g, d}(\mathbf{q}_{\nu g}) \right) \\ &= n_{\nu g}^{*1/2} \left(\hat{D}_{\nu g, d}(\hat{\mathbf{q}}_{\nu g}) - \hat{D}_{\nu g, d}(\mathbf{q}_{\nu g}) - \mathcal{D}_{g, d}(\hat{\mathbf{q}}_{\nu g}) + \mathcal{D}_{g, d}(\mathbf{q}_{\nu g}) \right) \\ &+ n_{\nu g}^{*1/2} \left(\hat{D}_{\nu g, d}(\mathbf{q}_{\nu g}) - D_{\nu g, d}(\mathbf{q}_{\nu g}) \right) \\ &+ n_{\nu g}^{*1/2} \left(\mathcal{D}_{g, d}(\hat{\mathbf{q}}_{\nu g}) - \mathcal{D}_{g, d}(\mathbf{q}_{\nu g}) \right). \end{aligned} \quad (28)$$

We can show the first term converges to zero in design in a similar fashion to the proof of Lemma 1. The remainder of the proof follows as $\sqrt{n_{\nu g}^*}(\hat{\mathbf{q}}_{\nu g} - \mathbf{q}_{\nu g}) = O_p(1)$. \square

Assumption 3.2.8 together with Assumption 3.1.2 gives

$$\frac{n_{\nu g}^{*1/2}}{N_{\nu g}} \sum_{U_{\nu g}} \underbrace{\left[\frac{\mathbf{I}(\|\mathbf{y}_i - \mathbf{q}_{\nu g}\| \leq d)}{1} \right]}_{\mathbf{b}_{\mathbf{q}, g, i}} \left[\frac{\mathbf{I}(i \in \mathcal{S}_{\nu g})}{\pi_i} - 1 \right] \Bigg| \mathcal{F}_{\nu} \xrightarrow{d} N(0, \Sigma_{\mathbf{q}, d}), \quad (29)$$

where

$$\Sigma_{\mathbf{q},d} = \frac{n_{\nu g}^*}{N_{\nu g}^2} \sum_i \sum_j (\pi_{ij} - \pi_i \pi_j) \frac{\mathbf{b}_{\mathbf{q},gi} \mathbf{b}'_{\mathbf{q},gj}}{\pi_i \pi_j}. \quad (30)$$

Theorem 6. Under Assumptions 3.1.1-3.1.2 and 3.2.4-3.2.8, for any sequence $\hat{\mathbf{q}}_{\nu g}$ satisfying (7),

$$n_{\nu g}^{1/2} \left[V \left(\widehat{D}_{\nu g,d}(\hat{\mathbf{q}}_{\nu g}) \right) \right]^{-1/2} \left(\widehat{D}_{\nu g,d}(\hat{\mathbf{q}}_{\nu g}) - D_{\nu g,d}(\mathbf{q}_{\nu g}) \right) \Big| \mathcal{F}_{\nu} \xrightarrow{d} N(0,1),$$

where

$$V \left(\widehat{D}_{\nu g,d}(\hat{\mathbf{q}}_{\nu g}) \right) = \mathbf{a}'_{\mathbf{q}} \Sigma_{\mathbf{q},d} \mathbf{a}_{\mathbf{q}}, \quad (31)$$

and

$$\mathbf{a}_{\mathbf{q}} = \begin{pmatrix} 1 \\ -D_{\nu g,d}(\mathbf{q}_{\nu g}) - \left(\frac{\partial \mathcal{D}_{g,d}(\mathbf{q}_{\nu g})}{\partial \mathbf{q}_{\nu g}} \right)' H_{s,N_{\nu g}}^{-1} \mathbf{q}_{\nu g} \\ H_{s,N_{\nu g}}^{-1} \left(\frac{\partial \mathcal{D}_{g,d}(\mathbf{q}_{\nu g})}{\partial \mathbf{q}_{\nu g}} \right) \end{pmatrix}, \quad (32)$$

where $H_{s,N_{\nu g}} = \frac{1}{N_{\nu g}} \sum_{i=1}^{N_{\nu g}} H_s(\mathbf{y}_i - \mathbf{q}_{\nu g})$ and $\Sigma_{\mathbf{q},d}$ is defined in (30).

Proof. We can use decomposition (28), then the leading term in decomposition is asymptotically normally distributed follows by using Slutsky's Theorem. \square

4.3 Variance Estimation

This section deals with estimating the variances of $\widehat{D}_{\nu g,d}(\hat{\boldsymbol{\mu}}_{\nu g})$ and $\widehat{D}_{\nu g,d}(\hat{\mathbf{q}}_{\nu g})$. We will introduce a naive estimator, a kernel estimator estimating the effect of error in $\hat{\gamma}$ and a jackknife estimator. The naive estimator ignores the error in estimating population center, and the extra piece of variance can be estimated by kernel smoothing and included in analytic variance estimator. An alternative is to incorporate a smoothing term in jackknife to get a consistent variance estimator. The smoothing term is only estimated once using the whole sample, so the benefit of replication procedure is not much affected.

4.3.1 Naive estimator

To estimate $V \left(\widehat{D}_{\nu g,d}(\hat{\boldsymbol{\mu}}_{\nu g}) \right)$ and $V \left(\widehat{D}_{\nu g,d}(\hat{\mathbf{q}}_{\nu g}) \right)$ as defined in (14) and (31), we need to estimate gradient vectors $\frac{\partial \mathcal{D}_{g,d}(\boldsymbol{\mu}_{\nu g})}{\partial \boldsymbol{\mu}_{\nu g}}$ and $\frac{\partial \mathcal{D}_{g,d}(\mathbf{q}_{\nu g})}{\partial \mathbf{q}_{\nu g}}$. We can either use some nonparametric smoothing method or propose a parametric model for the subpopulation distribution. First, we will show that if the subpopulation is a realization of an elliptical distribution, the variance due to estimating subpopulation center can be ignored.

Lemma 3. Assume random variable $\mathbf{Y}_g \sim EC_p(\boldsymbol{\mu}_g, \Lambda_g, \phi)$ with mean vector $\boldsymbol{\mu}_g$, and Λ_g is a non-negative definite matrix. We define the norm as,

$$\| \mathbf{u} \| = \sqrt{\mathbf{u}' \mathbf{B} \mathbf{u}}, \quad (33)$$

where \mathbf{B} is a non-negative definite matrix. Then

1. the partial derivative evaluated at superpopulation mean is $\mathbf{0}$, $\frac{\partial \mathcal{D}_{g,d}(\boldsymbol{\mu}_g)}{\partial \boldsymbol{\mu}_g} = \mathbf{0}$.
2. mean $\boldsymbol{\mu}_g$ coincide with median \mathbf{q}_g , $\boldsymbol{\mu}_g = \mathbf{q}_g$.

Lemma 4. Assume 3.2.2, $\frac{\partial \mathcal{D}_{g,d}(\boldsymbol{\gamma}_g)}{\partial \boldsymbol{\gamma}_g} = \mathbf{0}$ for some constant vector $\boldsymbol{\gamma}_g$, and the sequence of subpopulation centers $\boldsymbol{\gamma}_{\nu g}$ converges to $\boldsymbol{\gamma}_g$, $\lim_{\nu \rightarrow \infty} \boldsymbol{\gamma}_{\nu g} = \boldsymbol{\gamma}_g$. Then $\frac{\partial \mathcal{D}_{g,d}(\boldsymbol{\gamma}_{\nu g})}{\partial \boldsymbol{\gamma}_{\nu g}} = o(1)$.

Proof. The proof follows from Taylor expansion,

$$\frac{\partial \mathcal{D}_{g,d}(\boldsymbol{\gamma}_{\nu g})}{\partial \boldsymbol{\gamma}_{\nu g}} = \frac{\partial \mathcal{D}_{g,d}(\boldsymbol{\gamma}_g)}{\partial \boldsymbol{\gamma}_g} + \frac{\partial^2 \mathcal{D}_{g,d}(\boldsymbol{\gamma}_{\nu g}^*)}{\partial \boldsymbol{\gamma}_{\nu g}^2} (\boldsymbol{\gamma}_{\nu g} - \boldsymbol{\gamma}_g) = o(1). \quad \square$$

Lemma 3 and 4 imply that the extra variance due to estimating subpopulation centers can be ignored in elliptical distributions with a norm specified by (33). This special case is similar to case A of Randles (1982), where we can pretend that we are using true population center $\boldsymbol{\gamma}_{\nu g}$ without affecting the leading variance. So we can propose the following naive plug-in variance estimator for the leading term,

$$\widehat{V}_{NV} \left(\widehat{D}_{\nu g,d}(\hat{\boldsymbol{\gamma}}_{\nu g}) \right) = \left(1, -\widehat{D}_{\nu g,d}(\hat{\boldsymbol{\gamma}}_{\nu g}) \right) \widehat{\Sigma}_{\boldsymbol{\gamma},d} \left(1, -\widehat{D}_{\nu g,d}(\hat{\boldsymbol{\gamma}}_{\nu g}) \right)', \quad (34)$$

with

$$\widehat{\Sigma}_{\boldsymbol{\gamma},d,NV} = \frac{n_{\nu g}}{\widehat{N}_{\nu g}^2} \sum_i \sum_j \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}} \frac{\left[\frac{\mathbf{I}(\|\mathbf{y}_i - \hat{\boldsymbol{\gamma}}_{\nu g}\| \leq d)}{1} \right]}{\pi_i \pi_j} \left[\frac{\mathbf{I}(\|\mathbf{y}_i - \hat{\boldsymbol{\gamma}}_{\nu g}\| \leq d), 1 \right]}{\pi_i \pi_j}.$$

4.3.2 Estimating the effect of error in $\hat{\gamma}$

The naive estimator (34) ignores the piece of variance due to estimating population center and tends to underestimate the true variance for a general subpopulation. So we need to estimate $\frac{\partial \mathcal{D}_{g,d}(\boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}}$ and incorporate the extra piece of variance. Let $\boldsymbol{\zeta}_{g,d}(\boldsymbol{\gamma}) = \frac{\partial \mathcal{D}_{g,d}(\boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}}$. We have proposed an estimator of $\boldsymbol{\zeta}_{g,d}(\boldsymbol{\gamma})$ using kernel smoothing,

$$\hat{\boldsymbol{\zeta}}_{\nu g,d}(\boldsymbol{\gamma}) = \frac{1}{\widehat{N}_{\nu g} h} \sum_{\mathcal{S}_{\nu g}} K \left(\frac{d - \|\mathbf{y}_i - \boldsymbol{\gamma}\|}{h} \right) \boldsymbol{\psi}(\mathbf{y}_i - \boldsymbol{\gamma}) \frac{1}{\pi_i}. \quad (35)$$

The idea of estimator $\hat{\boldsymbol{\zeta}}_{\nu g,d}(\boldsymbol{\gamma})$ is to estimate $\mathcal{D}_{g,d}(\boldsymbol{\gamma})$ using primitive function of kernel $K(\cdot)$, and then take derivative with respect to $\boldsymbol{\gamma}$.

Lemma 5. Under Assumptions 3.1.1, 3.2.5 and 3.2.8-3.2.10, the estimator $\hat{\zeta}_{\nu g,d}(\gamma)$ is design consistent for $\zeta_{\nu g,d}(\gamma)$.

Lemma 6. Under Assumptions 3.1.1, 3.2.5 and 3.2.7-3.2.11, and assume the sequence of populations is such that

$$\sup_{\gamma \in \mathbb{R}^p} |\zeta_{\nu g,d}(\gamma) - \zeta_{g,d}(\gamma)| \rightarrow \mathbf{0}, \quad (36)$$

the kernel estimator $\hat{\zeta}_{\nu g,d}(\hat{\gamma}_{\nu g})$ is design consistent for $\zeta_{g,d}(\gamma_{\nu g})$ for every d .

Remarks: We have established uniform strong consistency of $\zeta_{g,d}(\gamma)$ for $\frac{\partial \mathcal{D}_{g,d}(\gamma)}{\partial \gamma}$, under appropriate superpopulation assumptions. Assumption (36) assumes that we are not working with the populations where $\zeta_{\nu g,d}(\gamma)$ does not converge to $\zeta_{g,d}(\gamma)$, which is on a zero-probability set.

We can directly plug the estimator (35) into estimators (14) and (31) to get a design-consistent variance estimator, and its performance will be evaluated in simulation study as shown in Wang (2008).

4.3.3 Jackknife variance estimator

This section only applies formally to mean-based estimator, but can be modified to include median-based estimator. To introduce the jackknife variance estimator for our application, we borrowed some strength from Da Silva and Opsomer (2006) and start by assuming there already exists a design consistent jackknife variance estimator for linear estimators. Then we define jackknife replicates in our case and show its consistency.

Theorem 7. Let $\hat{\theta}$ be a linear estimator for subpopulation g with

$$\hat{\theta}_g = \sum_{S_{\nu g}} w_i z_i,$$

where w_i is the survey weight and z_i has bounded $4 + \delta$ moments. Assume there is a jackknife replication procedure that generates L replicated estimates

$$\hat{\theta}_g^{(l)} = \sum_{S_{\nu g}} w_i^{(l)} z_i,$$

with $l = 1, 2, \dots, L$ and $w_i^{(l)}$ is replication weight for unite i in the l -th replicate. The replication variance estimator is defined as

$$\hat{V}_{JK}(\hat{\theta}_g) = \sum_{l=1}^L c_l \left(\hat{\theta}_g^{(l)} - \hat{\theta}_g \right)^2, \quad (37)$$

where c_l is a set of constants for the l -th replicate. Assumptions similar to (D1)-(D4) and (D6) in Da Silva and Opsomer (2006) are assumed.

We define the l -th jackknife replicate as

$$\begin{aligned} \hat{D}^{(l)}(\hat{\mu}_{\nu g}) &= \hat{D}_{\nu g,d}^{(l)}(\hat{\mu}_{\nu g}) \\ &+ \frac{1}{\hat{N}_{\nu g} h} \sum_{S_{\nu g}} \frac{1}{\pi_i} K\left(\frac{d - \|\mathbf{y}_i - \hat{\mu}_{\nu g}\|}{h}\right) \boldsymbol{\psi}^T(\mathbf{y}_i - \hat{\mu}_{\nu g})(\hat{\mu}_{\nu g}^{(l)} - \hat{\mu}_{\nu g}) \end{aligned} \quad (38)$$

$$\text{where } \hat{D}_{\nu g,d}^{(l)}(\hat{\mu}_{\nu g}) = \frac{1}{\hat{N}_{\nu g}^{(l)}} \sum_{i \in S_{\nu g}} w_i^{(l)} \mathbb{I}_{(\|\mathbf{y}_i - \hat{\mu}_{\nu g}\| \leq d)}, \hat{N}_{\nu g}^{(l)} = \sum_{i \in S_{\nu g}} w_i^{(l)} \text{ and } \hat{\mu}_{\nu g}^{(l)} = \frac{1}{\hat{N}_{\nu g}^{(l)}} \sum_{i \in S_{\nu g}} w_i^{(l)} \mathbf{y}_i.$$

Then the jackknife variance estimator

$$\hat{V}_{JK} \left(\hat{D}_{\nu g,d}(\hat{\mu}_{\nu g}) \right) = \sum_{l=1}^L c_l \left(\hat{D}^{(l)}(\hat{\mu}_{\nu g}) - \hat{D}(\hat{\mu}_{\nu g}) \right)^2 \quad (39)$$

is design consistent for $V \left(\hat{D}_{\nu g,d}(\hat{\mu}_{\nu g}) \right)$.

Proof. The use of this jackknife variance estimator is suggested by expression (11). The first term in (38) is used to approximate the variance of $\hat{D}_{\nu g,d}(\mu_{\nu g})$, and the second term approximates the variance due to estimating the center. \square

To examine the connection between the jackknife estimator and kernel estimator in section 4.3.2, we first ignore the second piece in (38) and compare it with the naive estimator (34). The naive estimator only approximates the linearized variance for the ratio, but the jackknife estimator usually overestimates this linearized variance. If we compare the whole jackknife replicate (38) with kernel estimator, the difference exists because of the nonlinearity of $\hat{D}_{\nu g,d}(\mu_{\nu g})$ and $\hat{\mu}_{\nu g}$.

A great advantage of jackknife variance estimator over the plug-in estimator is that we do not need to estimate the covariance matrix (13), which can be complicated in a large-scale complex survey. In jackknife variance estimation, we estimate the gradient vector based on the whole sample only once, but $\hat{D}_{\nu g,d}^{(l)}(\hat{\mu}_{\nu g})$ and $\hat{\mu}_{\nu g}^{(l)}$ will change with replicate. Variations of delete-1 jackknife like delete-d or delete-a-group jackknife can be used in complex surveys or in case of using median as measures of center.

Acknowledgement

This research was supported in part by the USDA Natural Resources Conservation Service cooperative agreement NRCS-68-3A75-4-122.

References

- Billingsley, P. (1968). *Convergence of Probability Measures*. John Wiley & Sons.
- Brown, B.M. (2007). Statistical uses of the spatial median. *Journal of the Royal Statistical Society, Series B: Methodology* 45. 25-30.

- Da Silva, D. and Opsomer, J.(2001). A kernel smoothing method to adjust for unit nonresponse in sample surveys. *the Canadian Journal of Statistics* 34. ??? - ???.
- Fuller, W. (2007). *Sampling statistics*.
- Randles, R.H.(1982). On the asymptotic normality of statistics with estimated parameters *The Annals of Statistics* 10. 462-474.
- Small, C.G. (2007). A survey of multidimensional medians *International Statistical Review* 58. 263-277.
- Wang, J. (2008). *Estimating the distance distribution of subpopulations for a large-scale complex survey*. Ph.D. thesis, Iowa State University.