Communications for Statistical Applications and Methods 2015, Vol. 22, No. 6, 655–664

# Nonresponse Adjusted Raking Ratio Estimation

Mingue Park<sup>1,a</sup>

<sup>a</sup>Department of Statistics, Korea University, Korea

#### Abstract

A nonresponse adjusted raking ratio estimator that consists of weighting adjustment using estimated response probability and raking procedure is often used to reduce the nonresponse bias and keep the calibration property of the estimator. We investigated asymptotic properties of nonresponse adjusted raking ratio estimator and proposed a variance estimator. A simulation study is used to examine the performance of suggested estimators.

Keywords: raking ratio estimator, logistic regression, propensity score, nonresponse, regression estimator

## 1. Introduction

Most surveys of human respondents entail a certain degree of nonresponse. Many studies on reducing or removing nonresponse bias have been done in survey statistics. A common way to handle unit nonresponse is weight adjustment in which the sampling weights for the respondents are adjusted so that the estimator based on the respondents only is (approximately) unbiased to the parameter of interest.

Using the framework of two-phase sampling design, weight adjustment is performed by deriving the response probability (often called propensity score). Nonresponse adjusted weights are then obtained by multiplying the inverse of estimated propensity score to the original sampling weight. Statistical models are often employed for the derivation of estimated response probability. One of the commonly used models, known as cell response model, assumes independent and identical response distribution for every elements in the same cell. The other popular model is logistic regression model that relates the binary response variable to the set of explanatory variables. The cell response model can be viewed as a logistic regression model with categorical explanatory variables that define cells. Results on the weight adjustment under the two-phase sampling framework were reviewed by Särndal and Lundström (2005) and Särndal *et al.* (1992). Kim and Kim (2007) gave the asymptotic results of weight adjusted estimator and proposed a possible variance estimator. Ekholm and Laaksonen (1991) is a practical example of this method.

Calibration estimators introduced by Deville and Särndal (1992) are often used to improve the efficiency of the estimator or to handle nonresponse when auxiliary information on the population (such as population total of auxiliary variables) is available. One typical example of such estimators is the raking ratio method in which the marginal distributions of several auxiliary categorical variables were used as auxiliary information. Raking ratio method, originally introduced by Deming and Stephan

This research was supported by Basic Science Research Program through the National Research Foundation of Korea

<sup>(</sup>NRF) funded by the Ministry of Education, Science and Technology (2013R1A1A2006363).

<sup>&</sup>lt;sup>1</sup> Department of Statistics, Korea University, 145 Anam-ro, Seongbuk-gu, Seoul 02841, Korea. E-mail: mpark2@korea.ac.kr

(1940) is often used to ensure that the estimator of the marginal distribution of auxiliary variables is equivalent to the known population distribution. The asymptotic results and small sample properties of the raking ratio estimator were investigated by Deville and Särndal (1992) and Deville *et al.* (1993).

In practice, a set of adjusted weights is often defined through two steps. In the first step, response probability, is estimated from a logistic regression model using appropriate set of explanatory variables. Then the final set of adjusted weights is obtained through raking ratio method or post stratification. Nonresponse adjusted raking ratio weights for both household and individuals were derived for the analysis of California Health Interviewing Survey (2011), that is a combined landline and cell telephone survey. Other examples in which the final weights were derived through two steps, nonresponse adjustment followed by raking, are National Health and Nutrition Examination Survey (NHANES) of USA (1996) and National Health and Nutrition Survey in South Korea (2010).

In our study, we derive the asymptotic properties of the nonresponse adjusted raking ratio estimator and suggest an explicit form of asymptotically unbiased variance estimator using the result of Deville and Särndal (1992) and Kim and Kim (2007). We also performed a small simulation study to investigate the properties of the nonresponse adjusted raking ratio estimator and variance estimator.

#### 2. Nonresponse Adjusted Raking Ratio Estimator

Consider the finite population  $U = \{1, 2, ..., N\}$  where the population size N is known. The parameter of interest is the population mean of the variable of interest, y, denoted by  $\bar{y}_N = N^{-1} \sum_{i \in U} y_i$ . If all sampled elements were observed, we consider the Horvitz-Thompson estimator (Horvitz and Thompson, 1952) of the form,  $\bar{y}_{HT} = N^{-1} \sum_{i \in A} \pi_i^{-1} y_i$ , as a design unbiased estimator of the population mean where  $A \in U$  is the set of indices in the sample and  $\pi_i = \Pr\{i \in A\}$ .

We define the response indicator variable of the unit *i* by

$$R_i = \begin{cases} 1, & \text{if unit } i \text{ responds,} \\ 0, & \text{if unit } i \text{ does not respond,} \end{cases}$$
(2.1)

for  $i \in A$  and define  $0 < p_{i|A} = \Pr\{R_i = 1 | i \in A\}$  as the response probability of sampled unit *i*. We assume  $R_i$ 's are independent with  $\operatorname{Var}(R_i|A) = p_{i|A}(1-p_{i|A})$ . By using the estimated response probability  $\hat{p}_{i|A}$  based on the logistic regression model,

$$logit(p_{i|A}|A) = \mathbf{z}_i' \alpha_A, \tag{2.2}$$

we define the nonresponse adjusted estimator

$$\bar{y}_{\scriptscriptstyle NA} = N^{-1} \sum_{i \in A} \pi_i^{-1} \, \hat{p}_{\scriptscriptstyle i|A}^{-1} \, R_i \, y_i.$$
(2.3)

The estimator of (2.3) is also called nonresponse weights adjustment (NWA) estimator, see Rosenbaum (1987).

Assume there exist *P*-variables of which population marginal distributions are known and they are used for defining raking ratio estimator. Let  $\mathbf{x}_i = (\mathbf{x}_{1i}, \dots, \mathbf{x}_{pi}, \dots, \mathbf{x}_{Pi})$  and  $\mathbf{x}_{pi} = (x_{1,pi}, \dots, x_{d_p,pi})$  be a set of indicators such that  $x_{j,pi} = 1$  if element *i* is in *j*<sup>th</sup> category of *p*<sup>th</sup> categorical variable and  $x_{j,pi} = 0$  otherwise, where  $d_p + 1$  is the number of categories defined by  $\mathbf{x}_p$ . Define the nonresponse adjusted raking ratio estimator by

$$\bar{y}_{_{NARR}} = N^{-1} \sum_{i \in A} \pi_i^{-1} \, \hat{p}_{_{i|A}}^{-1} \, \exp\left(\mathbf{x}_i' \hat{\lambda}\right) \, R_i \, y_i, \tag{2.4}$$

where  $\hat{\lambda}$  and  $\hat{\alpha}$  are the solution to

$$N^{-1} \sum_{i \in A} \pi_i^{-1} \hat{p}_{i|A}^{-1} \exp\left(\mathbf{x}_i' \hat{\boldsymbol{\lambda}}\right) R_i \mathbf{x}_i = \bar{\mathbf{x}}_N,$$
(2.5)

and

$$\sum_{i \in A} \pi_i^{-1} \{ R_i - [1 + \exp(-\mathbf{z}_i' \hat{\alpha}_A)] \} \mathbf{z}_i = \mathbf{0},$$
(2.6)

respectively, and  $\hat{p}_{i_A} = [1 + \exp(-\mathbf{z}'_i \hat{\alpha}_A)]^{-1}$ . We assume that there exists a unique solution to (2.5) and (2.6) and the solution  $\hat{\alpha}_A$  of (2.6) is a consistent estimator with respect to the response mechanism. Note that same variables could be used for both  $\mathbf{x}_i$  and  $\mathbf{z}_i$ . Equation (2.5) can be understood as a calibration equation so that the nonresponse adjusted raking ratio estimators of *x*-variables are equivalent to known population values. Equation (2.6) provides the (weighted) maximum likelihood estimator of  $\alpha_A$ .

To investigate the asymptotic properties of the  $\bar{y}_{NARR}$ , we consider a sequence of populations, samples and sampling designs assumed by Isaki and Fuller (1982). Assume population size,  $N_n(>n)$  increases as *n* increases. Assume  $\lim_{n\to\infty} \bar{\omega}_N$  exists, where  $\bar{\omega}_N = N^{-1} \sum_{i=1}^N \omega_i$  and  $\omega_i = (1, \mathbf{x}'_i, y_i)'$ . Also assume  $\omega_i$  has the finite fourth moments and its sample moments converge to the population moments such that

$$\left(N^{-1}\sum_{i\in A}\pi_i^{-1}\omega_i\omega_i'-N^{-1}\sum_{i=1}^N\omega_i\omega_i'\right)\bigg|\mathcal{F}_N=O_p\left(n^{-\frac{1}{2}}\right),\tag{2.7}$$

where  $\mathcal{F}_N = (\omega'_1, \dots, \omega'_N)'$  is the finite population. The assumption (2.7) means that mean and variance of Horvitz-Thompson estimator are well defined and converge to the corresponding population parameter with the order  $n^{-1/2}$  in probability. In most of sampling designs, Horvitz-Thompson estimator of the mean is unbiased and also has the finite variance of order  $n^{-1}$  and thus assumption, (2.7) is usually satisfied.

**Result 1.** Under the assumptions on response variable  $R_i$  of (2.1), logistic regression model of (2.2) and the assumption (2.7) on the sequence of populations and samples,

$$\bar{y}_{NARR} - \bar{y}_{N} = N^{-1} \sum_{i \in A} \pi_{i}^{-1} \left\{ p_{ijA} \mathbf{z}_{i}' \boldsymbol{\delta}_{N} + p_{ijA}^{-1} R_{i} \left( e_{i} - p_{ijA} \mathbf{z}_{i}' \boldsymbol{\delta}_{N} \right) \right\} + O_{p} \left( n^{-1} \right),$$
(2.8)

where  $e_i = y_i - \mathbf{x}'_i \boldsymbol{\beta}_N$ ,

$$\boldsymbol{\delta}_{\scriptscriptstyle N} = \left[\sum_{i \in U} p_{\scriptscriptstyle i\!i\!A} \left(1 - p_{\scriptscriptstyle i\!i\!A}\right) \mathbf{z}_i \mathbf{z}'_i\right]^{-1} \sum_{i \in U} \left(1 - p_{\scriptscriptstyle i\!i\!A}\right) \mathbf{z}_i e_i$$

and

$$\boldsymbol{\beta}_{N} = \left(\sum_{i=1}^{N} \mathbf{x}_{i} \mathbf{x}_{i}'\right)^{-1} \sum_{i=1}^{N} \mathbf{x}_{i} y_{i}$$

The brief proof of the result is given in Appendix. Using the result of (2.8), we could derive an asymptotic variance of the nonresponse adjusted raking ratio estimator as below.

$$V\left\{\bar{y}_{NARR} - \bar{y}_{N}|\mathcal{F}_{N}\right\} = E\left[V\left\{\bar{y}_{NARR} - \bar{y}_{N}|A_{N}, \mathcal{F}_{N}\right\}|\mathcal{F}_{N}\right] + V\left[E\left\{\bar{y}_{NARR} - \bar{y}_{N}|A_{N}, \mathcal{F}_{N}\right\}|\mathcal{F}_{N}\right],$$
(2.9)

where

$$E\left[V\left\{\bar{\mathbf{y}}_{\scriptscriptstyle NARR}-\bar{\mathbf{y}}_{\scriptscriptstyle N}|A_{\scriptscriptstyle N},\mathcal{F}_{\scriptscriptstyle N}\right\}|\mathcal{F}_{\scriptscriptstyle N}\right]\approx E\left[N^{-2}\sum_{i\in A}\pi_{i}^{-2}p_{_{\eta_{A}}}^{-1}\left(1-p_{_{\eta_{A}}}\right)\left(e_{i}-p_{_{\eta_{A}}}\mathbf{z}_{i}'\boldsymbol{\delta}_{\scriptscriptstyle N}\right)^{2}\Big|\mathcal{F}_{\scriptscriptstyle N}\right],$$
$$V\left[E\left\{\bar{\mathbf{y}}_{_{\scriptscriptstyle NARR}}-\bar{\mathbf{y}}_{\scriptscriptstyle N}|A_{\scriptscriptstyle N},\mathcal{F}_{\scriptscriptstyle N}\right\}|\mathcal{F}_{\scriptscriptstyle N}\right]\approx V\left(N^{-1}\sum_{i\in A}\pi_{i}^{-1}e_{i}\Big|\mathcal{F}_{\scriptscriptstyle N}\right),$$

 $A_N$  is the set of indices in the sample selected from the population  $U_N$ , and  $U_N$  is the  $N^{th}$  population of size N in the sequence.

The first term of (2.9) is the design expectation of conditional variance of the estimator conditioning on population and sample. The second term is the design variance of the Horvitz-Thompson estimator of the population mean of residuals.

For the definition of an asymptotically unbiased variance estimator of  $\bar{y}_{NARR}$ , assume that there exists a unbiased variance estimator under the full responses as a form of

$$\hat{V}(\bar{y}_{HT}) = \sum_{i \in A} \sum_{j \in A} \Omega_{ij} y_i y_j.$$
(2.10)

As an example of (2.10), we could consider Horvitz-Thompson variance estimator or Sen-Yates-Grundy variance estimator. For the definition of Horvitz-Thompson variance estimator and Sen-Yates-Grundy variance estimator, see Särndal *et al.* (1992).

One possible variance estimator of  $\bar{y}_{NARR}$  is

$$\hat{V} = \hat{V}_{res} + \hat{V}_{sam},\tag{2.11}$$

where

$$\begin{split} \hat{V}_{sam} &= \sum_{i \in A_R} \Omega_{ii} \hat{p}_i^{-1} (g_i \hat{e}_i)^2 + \sum_{i \neq j} \sum_{i,j \in A_R} \Omega_{ij} \hat{p}_i^{-1} \hat{p}_j^{-1} (g_i \hat{e}_i) (g_j \hat{e}_j), \\ \hat{V}_{res} &= N^{-2} \sum_{i \in A_R} \pi_i^{-2} \hat{p}_i^{-2} (1 - \hat{p}_i) (\hat{e}_i - \hat{p}_i \mathbf{z}_i' \hat{\boldsymbol{\gamma}})^2, \\ \hat{e}_i &= y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}_{NA}, \\ g_i &= \exp(\mathbf{x}_i' \hat{\boldsymbol{\lambda}}), \\ \hat{\boldsymbol{\gamma}} &= \left(\sum_{A_R} \pi_i^{-1} \hat{p}_i^{-1} (1 - \hat{p}_i) \mathbf{z}_i \mathbf{z}_i'\right)^{-1} \left(\sum_{A_R} \pi_i^{-1} (1 - \hat{p}_i) \hat{e}_i'\right), \end{split}$$

where  $\Omega_{ij}$  is the term that are used to define an unbiased variance estimator of  $\bar{y}_{HT}$  in (2.10). The variance estimator is obtained by estimating each term of (2.9) unbiasedly. For the definition of the variance estimator, we considered unbiased estimator with respect to design and response mechanism. We also use *g*-weight suggested by Särndal, Swensson and Wretman (1989) for the variance estimator of (2.11).

### 3. Simulation Study

A simulation study was performed to investigate the properties of the nonresponse adjusted raking ratio estimator and its variance estimator investigated at Section 2. For the simulation study, 18 different stratified populations of size 10,000 were generated from multivariate normal distribution

$$\begin{pmatrix} x_{hi} \\ z_{hi} \\ y_{hi} \end{pmatrix} \sim \text{MVN} \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix}, \begin{pmatrix} 1 & \rho_1 & \rho_2 \\ 1 & \rho_3 \\ & & 1 \end{bmatrix}, \quad h = 1, 2, 3, 4, \ i = 1, \dots, N_h,$$
(3.1)

where  $N_h = 1,000, 2,000, 3,000, 4,000$  for h = 1, 2, 3, 4 respectively. For  $\rho_1$ , we considered two levels 0.0 and 0.3, and three levels of (0.0, 0.3, 0.6) are considered for  $\rho_2$  and  $\rho_3$ . We also considered the variable  $w_{hi}$  generated from  $\chi^2$  distribution with 2 degrees of freedom.

The variable z generated the response probability using the logistic regression model,

$$R_i|A \sim \text{Bernoulli}(p_{i|A}), \quad \log(p_{i|A}) = -1 + z_i.$$
 (3.2)

Variable x and w were used to define membership variables that were used to calculate raking ratio estimator. The vector of indicator variables defined is

$$\mathbf{I}_{i} = \left(\mathbf{I}_{i1}^{\prime}, \mathbf{I}_{i2}^{\prime}\right)^{\prime}, \qquad (3.3)$$

where

$$\mathbf{I}_{i1}' = \begin{cases} (1,0,0,0), & \text{if } x_i < q_{0.25}^N, \\ (0,1,0,0), & \text{if } q_{0.25}^N \le x_i < q_{0.50}^N, \\ (0,0,1,0), & \text{if } q_{0.50}^N \le x_i < q_{0.75}^N, \\ (0,0,0,1), & \text{if } q_{0.75}^N \le x_i, \end{cases} \quad \mathbf{I}_{i2}' = \begin{cases} (1,0,0), & \text{if } w_i < q_{0.25}^{\chi^2}, \\ (0,1,0), & \text{if } q_{0.25}^{\chi^2} \le w_i < q_{0.50}^{\chi^2}, \\ (0,0,1), & \text{if } q_{0.75}^{\chi^2} \le x_i, \end{cases}$$

and  $q_p^N$  and  $q_p^{\chi^2}$  denote the  $p^{th}$  percentiles of the normal distribution with mean 2 and variance 1 and  $\chi^2$  distribution with 2 degree of freedom, respectively.

From each generated finite population, two sets of stratified random samples of size n = 400and n = 800 were selected and the stratum sample size were equal as  $n_h = n/4$  for h = 1, 2, 3, 4. We assumed the value of  $z_i$  is known for every element in the sample but the vector of membership indicators and variables of interest are known only for the responding element, that is,  $R_i = 1$ . We also assume that population total (or mean) of  $\mathbf{I}_i$  is known so that we could calculate the raking ratio estimator. The Monte Carlo sample sizes are all 20,000 and average response rate is about 0.70. As an estimator for the population mean, we considered  $\bar{y}_{NA}$  of (2.3) and  $\bar{y}_{NARR}$  of (2.4) with  $\mathbf{x}_i = \mathbf{I}_i$  of (3.3).

Table 1 shows the Monte Carlo relative bias in percent and variance of the two estimators. Monte Carlo relative bias of  $\hat{\theta}$ , that is an estimator of  $\theta$ , is defined as  $\theta^{-1}[E_{MC}(\hat{\theta}) - \theta]$ , where  $E_{MC}(\hat{\theta})$  is the Monte Carlo mean of  $\hat{\theta}$  and  $\theta$  is the parameter. With a sample of size 400 and 800, both estimators have negligible bias with less than 1% in absolute. Variance of  $\bar{y}_{NARR}$  decreases as the correlation between *z* and *y* gets stronger as expected. The variance of  $\bar{y}_{NARR}$  decreases as the correlation between *x* and *y* increases because the raking ratio estimator is asymptotically equivalent to the regression estimator efficient when the auxiliary variables are correlated with the variable of interest. Across all samples, the average reduction of the variance obtained by using  $\bar{y}_{NARR}$  is about 13% with maximum 36%.

Sample size	01	00	02	Relative	Bias in %	Variance	$e \times 10^{3}$
Sample Size	$\rho_1$	$p_2$	$p_3$	$\bar{y}_{NA}$	$\bar{y}_{NARR}$	$\overline{y}_{NA}$	$\bar{y}_{NARR}$
			0.0	-0.293	-0.284	4.940	4.738
		0.0	0.3	0.058	0.407	4.616	4.665
			0.6	-0.522	0.077	4.263	4.512
			0.0	0.502	0.435	5.058	4.268
	0.0	0.3	0.3	0.121	0.010	4.538	4.236
			0.6	0.591	0.826	3.922	3.847
			0.0	0.336	0.498	5.363	3.272
		0.6	0.3	-0.371	-0.031	4.610	3.249
n - 400			0.6	-0.190	0.154	4.096	2.804
n = 400			0.0	0.784	0.799	4.822	4.565
		0.0	0.3	-0.419	-0.381	4.429	4.409
			0.6	0.516	0.648	4.161	4.213
			0.0	-0.534	-0.589	4.911	4.266
	0.3	0.3	0.3	0.085	0.226	4.486	4.211
			0.6	-0.581	-0.115	4.137	4.097
			0.0	0.107	0.192	5.021	3.227
		0.6	0.3	-0.168	0.041	4.493	3.187
			0.6	-0.162	0.257	4.131	3.154
			0.0	-0.352	-0.341	2.323	2.256
		0.0	0.3	0.691	0.592	2.252	2.261
			0.6	0.136	0.330	1.966	2.064
			0.0	0.521	0.434	2.419	2.050
	0.0	0.3	0.3	0.399	0.559	2.039	1.878
			0.6	0.247	0.261	1.965	1.979
			0.0	0.361	0.526	2.537	1.588
		0.6	0.3	-0.048	0.035	2.193	1.483
n = 800			0.6	0.268	0.402	1.978	1.409
n = 800			0.0	0.784	0.802	2.323	2.197
		0.0	0.3	-0.436	-0.592	2.221	2.214
			0.6	0.356	0.400	1.884	1.990
		-	0.0	-0.493	-0.541	2.300	2.006
	0.3	0.3	0.3	0.272	0.159	2.300	2.065
			0.6	0.514	0.537	1.830	1.814
			0.0	0.130	0.244	2.315	1.515
		0.6	0.3	-0.087	-0.094	2.239	1.543
			0.6	0.449	0.581	1.965	1.532

Table 1: Monte Carlo relative bias and variance of estimators

We also calculate the variance estimator of  $\bar{y}_{NARR}$  given in (2.11). Table 2 shows the Monte Carlo properties of  $\hat{V}_{sam}$ ,  $\hat{V}_{res}$  and  $\hat{V}$ . The relative bias of the variance estimator in percent is

$$\left[V_{MC}(\bar{y}_{NARR})\right]^{-1}\left[E_{MC}\left(\hat{V}\right)-V_{MC}\left(\bar{y}_{NARR}\right)\right].$$

For the samples of size 400, suggested variance estimator underestimates the true variance in all population. Relative bias of the variance estimator is about -4% to -7%. Underestimation of the variance estimator with sample size 400 is mainly due to the Taylor approximation in which the second and higher order term of the approximation are not considered for the variance estimation. The absolute relative bias of the suggested variance estimator with a sample of size 800 is less than 2.2%. For both sample sizes, the component of the variance estimator,  $\hat{V}_{res}$  makes up about 34% to 36% of the  $\hat{V}$ . Thus if only  $\hat{V}_{sam}$  is used to estimate the variance of  $\bar{y}_{NARR}$ , which is common in practice, it is expected to underestimate the true variance severely.

Somela size				М	onte Carlo me	ean	Relative bias	$V_{MC}(\hat{V})$
Sample size	$\rho_1$	$\rho_2$	$\rho_3$	$\hat{V}_{sam}$	$\hat{V}_{res}$	$\hat{V}$	of Ŷ (%)	$\times 10^{6}$
			0.0	2.863	1.669	4.531	-4.358	0.624
		0.0	0.3	2.943	1.509	4.452	-4.570	0.438
			0.6	2.920	1.374	4.294	-4.835	0.700
			0.0	2.576	1.487	4.062	-4.821	0.398
	0.0	0.3	0.3	2.676	1.359	4.035	-4.747	0.341
			0.6	2.590	1.111	3.701	-3.788	0.433
			0.0	1.948	1.149	3.097	-5.343	0.328
		0.6	0.3	2.019	1.067	3.087	-4.999	0.216
n = 400			0.6	1.889	0.732	2.621	-6.532	0.308
n = 400			0.0	2.821	1.574	4.394	-3.748	0.520
		0.0	0.3	2.731	1.493	4.223	-4.217	0.448
			0.6	2.768	1.255	4.024	-4.480	0.476
			0.0	2.615	1.426	4.042	-5.253	0.322
	0.3	0.3	0.3	2.606	1.459	4.066	-3.444	0.345
			0.6	2.709	1.163	3.871	-5.510	0.423
			0.0	1.971	1.098	3.069	-4.897	0.219
		0.6	0.3	1.941	1.052	2.992	-6.108	0.286
			0.6	1.958	0.982	2.940	-6.764	0.606
			0.0	1.412	0.878	2.290	1.497	0.090
		0.0	0.3	1.402	0.890	2.293	1.428	0.111
			0.6	1.402	0.663	2.065	0.024	0.062
			0.0	1.265	0.782	2.047	-0.149	0.053
	0.0	0.3	0.3	1.225	0.677	1.902	1.250	0.036
		0.3	0.6	1.306	0.683	1.989	0.493	0.360
			0.0	0.962	0.613	1.574	-0.864	0.052
		0.6	0.3	0.966	0.527	1.494	0.706	0.029
n = 800			0.6	0.960	0.452	1.413	0.295	0.076
			0.0	1.388	0.824	2.211	0.646	0.076
		0.0	0.3	1.407	0.807	2.214	0.002	0.076
			0.6	1.403	0.629	2.032	2.125	0.084
			0.0	1.287	0.744	2.032	1.267	0.041
	0.3	0.3	0.3	1.278	0.772	2.049	-0.737	0.054
			0.6	1.229	0.595	1.824	0.540	0.039
			0.0	0.971	0.570	1.541	1.692	0.027
		0.6	0.3	0.982	0.588	1.570	1.740	0.031
			0.6	0.948	0.554	1.502	-1.985	0.073

Table 2: Monte Carlo properties of variance estimators

# Appendix A: Proof of result 1

Step 1. Show

$$N^{-1} \sum_{i \in A} \pi_i^{-1} \hat{p}_{i|A}^{-1} R_i \omega_i \omega_i' - N^{-1} \sum_{i=1}^N \omega_i \omega_i' = O_p\left(n^{-\frac{1}{2}}\right).$$
(A.1)

**Proof**: Let  $v_i$  be an any element of  $\omega_i \omega'_i$ . Then, by Kim and Kim (2007),

$$\left(\bar{v}_{\scriptscriptstyle NA}-\bar{v}_{\scriptscriptstyle NAL}\right)\left|\mathcal{F}_{N}=O_{p}\left(n^{-1}\right),\right.$$

where

$$\bar{\nu}_{\scriptscriptstyle NAL} = N^{-1} \sum_{i \in A} \pi_i^{-1} \left\{ p_{\scriptscriptstyle i|A} \mathbf{z}'_i \boldsymbol{\gamma}_n + p_{\scriptscriptstyle i|A}^{-1} R_i \left( \nu_i - p_{\scriptscriptstyle i|A} \mathbf{z}'_i \boldsymbol{\gamma}_n \right) \right\},\,$$

and

$$\boldsymbol{\gamma}_n = \left[\sum_{i \in A} \pi_i^{-1} p_{i|A} \left(1 - p_{i|A}\right) \mathbf{z}_i \mathbf{z}'_i\right]^{-1} \sum_{i \in A} \pi_i^{-1} \left(1 - p_{i|A}\right) \mathbf{z}_i \boldsymbol{\nu}_i.$$

Note that

$$E\left(\bar{v}_{\scriptscriptstyle NAL}\big|\mathcal{F}_{\scriptscriptstyle N}\right) = E\left\{E\left(\bar{v}_{\scriptscriptstyle NAL}\big|A_{\scriptscriptstyle N},\mathcal{F}_{\scriptscriptstyle N}\right)\big|\mathcal{F}_{\scriptscriptstyle N}\right\} = \bar{v}_{\scriptscriptstyle N}$$

and

$$V\left(\bar{v}_{\scriptscriptstyle NAL}\middle|\mathcal{F}_{\scriptscriptstyle N}\right) = V\left(N^{-1}\sum_{i\in A}\pi_{i}^{-1}v_{i}\middle|\mathcal{F}_{\scriptscriptstyle N}\right) + N^{-2}E\left\{\sum_{i\in A}\pi_{i}^{-2}\left[p_{\scriptscriptstyle i|A}^{-1}(1-p_{\scriptscriptstyle i|A})\left(v_{i}-p_{\scriptscriptstyle i|A}\mathbf{z}_{i}'\boldsymbol{\gamma}_{n}\right)^{2}\right]\middle|\mathcal{F}_{\scriptscriptstyle N}\right\}$$
$$= O\left(n^{-1}\right).$$

Thus, for all element  $v_i$  of  $\omega_i \omega'_i$ ,

$$\left(\bar{\nu}_{\scriptscriptstyle NA}-\bar{\nu}_{\scriptscriptstyle N}\right)\left|\mathcal{F}_{\scriptscriptstyle N}=O_p\left(n^{-\frac{1}{2}}\right).\right.$$

L		
L		
L		_

Step 2. Show

$$\bar{y}_{\scriptscriptstyle NARR} = \bar{y}_{\scriptscriptstyle NA} + (\bar{\mathbf{x}}_{\scriptscriptstyle N} - \bar{\mathbf{x}}_{\scriptscriptstyle NA})\hat{\boldsymbol{\beta}}_{\scriptscriptstyle NA} + O_p\left(n^{-1}\right),$$

where

$$\hat{\boldsymbol{\beta}}_{\scriptscriptstyle NA} = \left(\sum_{A} \pi_i^{-1} \hat{p}_{\scriptscriptstyle \partial A}^{-1} R_i \mathbf{x}_i' \mathbf{x}_i\right)^{-1} \left(\sum_{A} \pi_i^{-1} \hat{p}_{\scriptscriptstyle \partial A}^{-1} R_i \mathbf{x}_i' y_i\right).$$

**Proof**: For any sample and a set of respondents, there exists a solution of  $\hat{\lambda}$  that satisfies the equation (2.5) as explained by Deville and Särndal (1992). Define the function  $\phi(\cdot)$  as

$$\phi\left(\hat{\lambda}\right) = N^{-1} \sum_{i \in A} \pi_i^{-1} \hat{p}_{i|\lambda}^{-1} \left[ \exp\left(\mathbf{x}_i'\hat{\lambda}\right) - 1 \right] \mathbf{x}_i R_i.$$
(A.2)

Note that  $\phi(\mathbf{0}) = 0$ . By using the Taylor expansion of  $\phi(\hat{\lambda})$  at  $\hat{\lambda} = \mathbf{0}$ , we obtain

$$\phi\left(\hat{\lambda}\right) = \phi(\mathbf{0}) + \phi'(\boldsymbol{\xi})\left(\hat{\lambda} - \mathbf{0}\right) = \phi'(\boldsymbol{\xi})\hat{\lambda},$$

where  $\boldsymbol{\xi}$  is a vector in the interval **0** and  $\hat{\boldsymbol{\lambda}}$ , called **B**. Now let  $\boldsymbol{\eta}$  be a vector such that

$$\boldsymbol{\eta}=\bar{\mathbf{x}}_{\scriptscriptstyle N}-\bar{\mathbf{x}}_{\scriptscriptstyle NA}.$$

Then, by (2.5),

$$\hat{\boldsymbol{\lambda}} = [\boldsymbol{\phi}'(\boldsymbol{\xi})]^{-1}(\boldsymbol{\eta})$$

662

and thus, by the characteristics of exponential function,

$$\left\|\hat{\lambda}\right\| \leq K'_{\phi} \left\|\eta\right\|,$$

where  $K'_{\phi}$  is a positive constant such that  $\|[\phi'(\xi)]^{-1}\| \le K_{\phi}$  for all  $\xi$  in **B**. By (2.8),  $\|\eta\| = O_p(n^{-1/2})$  and thus

$$\hat{\lambda} = O_p\left(n^{-\frac{1}{2}}\right).$$

Now using  $\exp(\mathbf{x}'_i \hat{\boldsymbol{\lambda}}) = 1 + \mathbf{x}'_i \hat{\boldsymbol{\lambda}} + \theta(\mathbf{x}'_i \boldsymbol{\xi})$  and (2.7), we have

$$\boldsymbol{\lambda} = \left[ N^{-1} \sum_{i \in A} \pi_i^{-1} \hat{p}_{iA}^{-1} R_i \mathbf{x}_i \mathbf{x}_i' \right]^{-1} (\bar{\mathbf{x}}_N - \bar{\mathbf{x}}_{NA}) + O_p(n^{-1}),$$

because  $\max_{\boldsymbol{\xi}} \theta(\mathbf{x}'_{i}\boldsymbol{\xi}) = O(n^{-1})$ . Thus,

$$\bar{y}_{\scriptscriptstyle NARR} = N^{-1} \sum_{i \in A} \pi_i^{-1} \hat{p}_{i|A} \exp\left(\mathbf{x}_i' \hat{\boldsymbol{\lambda}}\right) R_i y_i$$
  
=  $\bar{y}_{\scriptscriptstyle NA} + (\bar{\mathbf{x}}_{\scriptscriptstyle N} - \bar{\mathbf{x}}_{\scriptscriptstyle NA}) \hat{\boldsymbol{\beta}}_{\scriptscriptstyle NA} + O_p\left(n^{-1}\right).$  (A.3)

_		

Step 3. Show the final result.

Proof: By (A.1),

$$\hat{\boldsymbol{\beta}}_{\scriptscriptstyle NA} - \boldsymbol{\beta}_{\scriptscriptstyle N} = O_p\left(n^{-\frac{1}{2}}\right),$$

and thus, due to the result of Kim and Kim (2007),

$$\begin{split} \bar{\mathbf{y}}_{\scriptscriptstyle NARR} &= \bar{\mathbf{y}}_{\scriptscriptstyle NA} + (\bar{\mathbf{x}}_{\scriptscriptstyle N} - \bar{\mathbf{x}}_{\scriptscriptstyle NA}) \boldsymbol{\beta}_{\scriptscriptstyle N} + O_p \left( n^{-1} \right) \\ &= \bar{\mathbf{x}}_{\scriptscriptstyle N} \boldsymbol{\beta}_{\scriptscriptstyle N} + N^{-1} \sum_{A} \pi_i^{-1} \hat{p}_{\scriptscriptstyle i|A}^{-1} R_i \left( y_i - \mathbf{x}_i \boldsymbol{\beta}_{\scriptscriptstyle N} \right) + O_p \left( n^{-1} \right) \\ &= \bar{\mathbf{x}}_{\scriptscriptstyle N} \boldsymbol{\beta}_{\scriptscriptstyle N} + N^{-1} \sum_{i \in A} \pi_i^{-1} \left\{ p_{\scriptscriptstyle i|A} \mathbf{z}_i' \boldsymbol{\delta}_{\scriptscriptstyle N} + p_{\scriptscriptstyle i|A}^{-1} R_i \left( e_i - p_{\scriptscriptstyle i|A} \mathbf{z}_i' \boldsymbol{\delta}_{\scriptscriptstyle N} \right) \right\} + O_p \left( n^{-1} \right). \end{split}$$
(A.4)

#### Appendix B: Concluding remark

In many survey practices, both weight adjustment using the estimated response probability and raking ratio procedure are employed to reduce nonresponse bias and obtain the known population auxiliary information if the estimator were applied to auxiliary variables. We proposed a consistent variance estimator using the asymptotic result of nonresponse adjusted raking ratio estimator. A simulation study indicated that a variance estimator that does not consider variability due to nonresponse and adjustment could severely underestimate true variance.

#### References

- Center for Health Policy Research (2011). *Report 5: Weighting and Variance Estimation (CHIS 2009 Methodology Series)*, UCLA Center for Health Policy Research, Los Angeles, CA.
- Deming, W. E. and Stephan, F. F. (1940). On a least squares adjustment of a sampled frequency table when the expected marginal totals are known, *Annals of Mathematical Statistics*, **11**, 427–444.
- Deville, J. C. and Särndal, C. E. (1992). Calibration estimators in survey sampling, *Journal of the American Statistical Association*, 87, 376–382.
- Deville, J. C., Sarndal, C. E. and Sautory, O. (1993). Generalized raking procedures in survey sampling, *Journal of the American Statistical Association*, 88, 1013–1020.
- Ekholm, A. and Laaksonen, S. (1991). Weighting via response modeling in the Finnish household budget survey, *Journal of Official Statistics*, 7, 325–337.
- Horvitz, D. G. and Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe, *Journal of the American Statistical Association*, 47, 663–685.
- Isaki, C. T. and Fuller, W. A. (1982). Survey design under the regression superpopulation model, *Journal of the American Statistical Association*, **77**, 89–96.
- Kim, J. K. and Kim, J. J. (2007). Nonresponse weighting adjustment using estimated response probability, *Canadian Journal of Statistics*, 35, 501–514.
- National Health and Nutrition Examination Survey (1996). *Analytic and Reporting Guidelines: The Third National Health and Nutrition Examination Survey, NHANE III (1988–1994)*, National Center for Health Statistics, Centers for Disease Control and Prevention, Hyattsville, MD.
- National Health and Nutrition Survey (2010). *Analytic Guidelines: The fourth NHNS*, Ministry of Health and Welfare, Seoul, Korea.
- Rosenbaum, P. R. (1987). Model-based direct adjustment, *Journal of the American Statistical Asso*ciation, 82, 387–394.

Särndal, C. E. and Lundström, S. (2005). Estimation in Surveys with Nonresponse, Wiley, New York.

- Särndal, C. E., Swensson, B. and Wretman, J. H. (1989). The weighted residual technique for estimating the variance of the general regression estimator of the finite population total, *Biometrika*, 76, 527–537.
- Särndal, C. E., Swensson, B. and Wretman, J. H. (1992). Model Assisted Survey Sampling, Springer-Verlag, New York.

Received September 3, 2015; Revised October 6, 2015; Accepted October 6, 2015