Scientific
Research

# A New Estimator Using Auxiliary Information in Stratified Adaptive Cluster Sampling

**Nipaporn Chutiman[*], Monchaya Chiangpradit, Sujitta Suraphee**

Department of Mathematics, Faculty of Science, Mahasarakham University, Maha Sarakham, Thailand
Email: [*]j3832024@hotmail.com

## ABSTRACT

In this paper, we study the estimators of the population mean in stratified adaptive cluster sampling by using the information of the auxiliary variable. Simulations showed that if the variable of interest ($y$) and the auxiliary variables ($x$,$z$) have high positive correlation then the estimate of the mean square error of the ratio estimators is less than the estimate of the mean square error of the product estimator. The estimators which use only one auxiliary variable were better than the estimators which use two auxiliary variables.

**Keywords:** Stratified Adaptive Cluster Sampling; Auxiliary Variable; Ratio Estimator; Product Estimator

## 1. Introduction

Adaptive cluster sampling, proposed by Thompson [1], is an efficient method for sampling rare and hidden clustered populations. In adaptive cluster sampling, an initial sample of units is selected by simple random sampling. If the value of the variable of interest from a sampled unit satisfies a pre-specified condition $C$, that is $\{i, y_i \geq c\}$, then the unit's neighborhood will also be added to the sample. If any other units that are "adaptively" added also satisfy the condition $C$, then their neighborhoods are also added to the sample. This process is continued until no more units that satisfy the condition are found. The set of all units selected and all neighboring units that satisfy the condition is called a network. The adaptive sample units, which do not satisfy the condition are called edge units. A network and its associated edge units are called a cluster. If a unit is selected in the initial sample and does not satisfy the condition $C$, then there is only one unit in the network. A neighborhood must be defined such that if unit $i$ is in the neighborhood of unit $j$ then unit $j$ is in the neighborhood of unit $i$. In this paper, a neighborhood of a unit is defined as the four spatially adjacent units, that is to the left, right, top and bottom of that unit as shown in **Figure 1**.

**Figure 1** illustrates the example of a network. The unit with a star is the initial unit selected. The condition to adaptively added units is a value greater than or equal to 1.

Units that are to the left, right, top, and bottom of one another making up a neighborhood. The units in the gray shading form a single network. The units in bold numbers are edge units of the network. The network and its edge units make up a cluster.

Adaptive cluster sampling are applied in stratified random sampling. In adaptive cluster sampling, an initial stratified sample is selected from a population, and whenever the variable of interest for any unit is observed to satisfy the condition, the neighborhood of that unit is added in the sample. Sometimes other variables are related to the variable of interest $y$. We can obtain additional information for estimating the population mean. The use of an auxiliary variable is a common method to improve the precision of estimates of a population mean. In this paper, we will study the estimator of population mean in stratified adaptive cluster sampling using an auxiliary
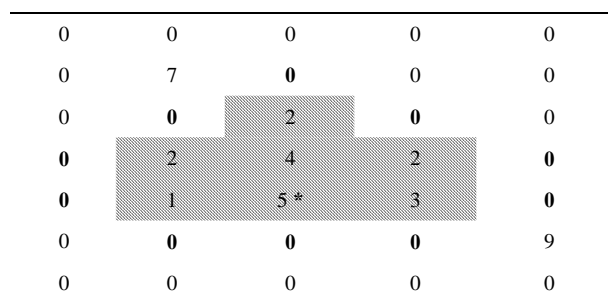


**Figure 1. The example of network where a unit neighborhood is defined as four spatially adjacent units.**

[*]Corresponding author.

variable. Some comparisons are made using a simulation.

## 2. Stratified Adaptive Cluster Sampling

For stratified adaptive cluster sampling, the population consists of $N$ units partitioned into $L$ strata based on prior information about units that are similar, and it is assumed that the population ignores crossover between strata. The population in each stratum consists of $N_h$ units $(h = 1, 2, \cdots, L)$. The population mean of the variable of interest in stratum $h$ is $\mu_{yh}$. An initial sample of unit size $n_h$ is selected by simple random sampling without replacement and for those units selected that satisfy the condition. Then the unit's neighborhood is added to the sample.

Define

$$w_{yhi} = \frac{1}{m_{hi}} \sum_{j \in \psi_{hi}} y_{hj}$$

is the average of the $y$-values of the network to which $u_{hi}$ belongs. $\psi_{hi}$ is the network that include unit $i$ in stratum $h$ and $m_{hi}$ is the size of network that include unit $i$ in stratum $h$. The estimator of the population mean based on Hansen-Hurwitz estimator (Thompson and Seber [2]) is

$$\bar{y}_{st\_a} = \sum_{h=1}^{L} \frac{N_h}{n_h} \bar{w}_{yh} \qquad (1)$$

where

$$\bar{w}_{yh} = \frac{1}{n} \sum_{h=1}^{L} w_{yhi}$$

The variance of $\bar{y}_{st\_a}$ is

$$V(\bar{y}_{st\_a}) = \frac{1}{N} \sum_{h=1}^{L} N_h (N_h - n_h) \frac{S_{yhw}^2}{n_h} \qquad (2)$$

where

$$S_{yhw}^2 = \frac{1}{N_h - 1} \sum_{h=1}^{N_h} \left( w_{yhi} - \bar{\xi}_{yh} \right)^2$$

and

$$\bar{\xi}_{yh} = \sum_{h=1}^{N_h} w_{yhi} \Big/ N_h \; .$$

The estimate of $V(\bar{y}_{st\_a})$ is

$$\hat{V}(\bar{y}_{st\_a}) = \frac{1}{N} \sum_{h=1}^{L} N_h (N_h - n_h) \frac{s_{yhw}^2}{n_h} \qquad (3)$$

where

$$s_{yhw}^2 = \frac{1}{n_h - 1} \sum_{h=1}^{n_h} \left( w_{yhi} - \bar{w}_{yh} \right)^2 \; .$$

## 3. Propose Estimators

The estimator of the population mean in stratified adaptive cluster sampling using two auxiliary variables ($x,z$) is (Walid A. Abu-Dayyeh, M. S. Ahmed, R. A. Ahmed and Hassen A. Muttlak, [3]),

$$\bar{y}_{st\_a\_xz} = \bar{y}_{st\_a} \left( \frac{\bar{x}_{st\_a}}{\mu_x} \right)^{\alpha_1} \left( \frac{\bar{z}_{st\_a}}{\mu_z} \right)^{\alpha_2} \qquad (4)$$

$\alpha_1 = 0$ and $\alpha_2 = 0$ is called mean per unit, $\alpha_1 = -1$ and $\alpha_2 = -1$ is called multivariate ratio estimator, $\alpha_1 = -1$ and $\alpha_2 = -1$ is called multivariate ratio estimator, $\alpha_1 = 1$ and $\alpha_2 = 1$ is called multivariate product estimator, $\alpha_1 = -1$ and $\alpha_2 = 0$ is called ratio estimator using $x$, $\alpha_1 = 0$ and $\alpha_2 = -1$ is called ratio estimator using $z$, $\alpha_1 = 1$ and $\alpha_2 = 0$ is called product estimator using $x$ and $\alpha_1 = 0$ and $\alpha_2 = 1$ is called product estimator using $z$.

Let

$$e_0 = \frac{\bar{y}_{st\_a} - \mu_y}{\mu_y}, \quad e_1 = \frac{\bar{x}_{st\_a} - \mu_x}{\mu_x}$$

and

$$e_2 = \frac{\bar{z}_{st\_a} - \mu_z}{\mu_z}$$

So

$$E(e_0) = E(e_1) = E(e_2) = 0$$

$$E(e_0^2) = V(e_0) + \left[ E(e_0) \right]^2$$

$$= V\left[ \frac{\bar{y}_{st\_a} - \mu_y}{\mu_y} \right] = \frac{1}{\mu_y^2} V(\bar{y}_{st\_a})$$

$$E(e_0^2) = \frac{1}{\mu_y^2} \sum_{h=1}^{L} N_h (N_h - n_h) \frac{S_{yhw}^2}{n_h} = \frac{1}{\mu_y^2} A_y,$$

$$A_y = \sum_{h=1}^{L} N_h (N_h - n_h) \frac{S_{yhw}^2}{n_h}$$

Thus

$$E(e_1^2) = \frac{A_x}{\mu_x^2}, \quad A_x = \sum_{h=1}^{L} N_h (N_h - n_h) \frac{S_{xhw}^2}{n_h}$$

$$E(e_2^2) = \frac{A_z}{\mu_z^2}, \quad A_z = \sum_{h=1}^{L} N_h (N_h - n_h) \frac{S_{zhw}^2}{n_h},$$

$$E(e_0 e_1) = \frac{1}{\mu_x \mu_y} \sum_{h=1}^{L} N_h (N_h - n_h) \frac{S_{xhw,yhw}}{n_h} = \frac{1}{\mu_x \mu_y} A_{xy}$$

$$E(e_0 e_2) = \frac{1}{\mu_y \mu_z} \sum_{h=1}^{L} N_h (N_h - n_h) \frac{S_{yhw,zhw}}{n_h} = \frac{1}{\mu_y \mu_z} A_{yz}$$

and

$$E(e_1 e_2) = \frac{1}{\mu_x \mu_z} \sum_{h=1}^{L} N_h (N_h - n_h) \frac{S_{xhw,zhw}}{n_h} = \frac{1}{\mu_x \mu_z} A_{xz}.$$

So,

$$\bar{y}_{st\_a\_xz} = \mu_y (1+e_0)(1+e_1)^{\alpha_1}(1+e_2)^{\alpha_2}$$

$$= \mu_y \left\{ 1 + e_0 + \alpha_1 e_1 + \alpha_2 e_2 + \alpha_1 e_0 e_1 + \alpha_2 e_0 e_2 \right.$$

$$\left. + \alpha_1 \alpha_2 e_1 e_2 + \frac{\alpha_1(\alpha_1 - 1)}{2} e_1^2 + \frac{\alpha_2(\alpha_2 - 1)}{2} e_2^2 + \cdots \right\}$$

$$MSE(\bar{y}_{st\_a\_xz}) = E(\bar{y}_{st\_a\_xz} - \mu_y)^2$$

$$= \mu_y^2 E\left[ e_0^2 + \alpha_1^2 e_1^2 + \alpha_2^2 e_2^2 + 2\alpha_1 e_0 e_1 \right.$$

$$\left. + 2\alpha_2 e_0 e_2 + 2\alpha_1 \alpha_2 e_1 e_2 + \cdots \right]$$

$$= \mu_y^2 \left\{ \frac{A_y}{\mu_y^2} + \alpha_1^2 \frac{A_x}{\mu_x^2} + \alpha_2^2 \frac{A_z}{\mu_z^2} + 2\alpha_1 \frac{A_{xy}}{\mu_x \mu_y} \right.$$

$$\left. + 2\alpha_2 \frac{A_{yz}}{\mu_y \mu_z} + 2\alpha_1 \alpha_2 \frac{A_{xz}}{\mu_x \mu_z} \right\}$$

To find $\alpha_1$ and $\alpha_2$ which minimizes $MSE(\bar{y}_{st\_a\_xz})$ take partial derivative of $MSE(\bar{y}_{st\_a\_xz})$ with respect to $\alpha_1$, $\alpha_2$ and set it equal to zero.

$$\frac{\partial MSE(\bar{y}_{st\_a\_xz})}{\partial \alpha_1} = 0$$

and

$$\frac{\partial MSE(\bar{y}_{st\_a\_xz})}{\partial \alpha_2} = 0$$

So the optimum values of $\alpha_1$ and $\alpha_2$ are

$$\alpha_1 = -\left( \frac{\mu_x A_{xy}}{A_x \mu_y} + \alpha_2 \frac{\mu_x A_{xz}}{A_x \mu_z} \right)$$

and

$$\alpha_2 = \frac{\mu_z}{\mu_y} \frac{(A_{xy} A_{xz} - A_x A_{yz})}{(A_x A_z - A_{xz}^2)}$$

The estimate of $\alpha_1$ is

$$\hat{\alpha}_1 = -\left( \frac{\mu_x \hat{A}_{xy}}{\hat{A}_x \mu_y} + \hat{\alpha}_2 \frac{\mu_x \hat{A}_{xz}}{\hat{A}_x \mu_z} \right)$$

and the estimate of $\alpha_2$ is

$$\hat{\alpha}_2 = \frac{\mu_z}{\mu_y} \frac{(\hat{A}_{xy} \hat{A}_{xz} - \hat{A}_x \hat{A}_{yz})}{(\hat{A}_x \hat{A}_z - \hat{A}_{xz}^2)}$$

where

$$\hat{A}_x = \sum_{h=1}^{L} N_h (N_h - n_h) \frac{s_{xhw}^2}{n_h},$$

$$\hat{A}_z = \sum_{h=1}^{L} N_h (N_h - n_h) \frac{s_{zhw}^2}{n_h},$$

$$\hat{A}_{xy} = \sum_{h=1}^{L} N_h (N_h - n_h) \frac{s_{xhw,yhw}}{n_h},$$

$$\hat{A}_{xz} = \sum_{h=1}^{L} N_h (N_h - n_h) \frac{s_{xhw,zhw}}{n_h}$$

and

$$\hat{A}_{xy} = \sum_{h=1}^{L} N_h (N_h - n_h) \frac{s_{yhw,zhw}}{n_h}.$$

## 4. Simulation Study

This section, the simulation $x$-values $z$-values and $y$-values from Chutiman, N. and Kumphon, B. [4] were studied. The data partition into 4 stratum. The stratum size is $20 \times 5 = 100$ units. The populations were shown in **Figures 2-4**. Sample of units is selected by simple random sampling without replacement. The $y$-values are obtained for keeping the sample network. In each the sample network, the $x$-values and $z$-values are obtained. The condition for added units in the sample is defined by $C = \{y : y > 0\}$.

For each estimator 5000 iterations were performed to obtain an accuracy estimate. Initial SRS sizes were varied $n_h$ = 5, 10, 15, 20 and 30 were used. The estimated mean square error of the estimate mean is

$$M\hat{S}E(\bar{y}) = \frac{1}{5000} \sum_{i=1}^{5000} (\bar{y}_i - \mu_y)^2,$$

where $\bar{y}_i$ is the value for the relevant estimator for sample $i$.

The estimate of the mean square of the estimators $M\hat{S}E(\bar{y}_{\alpha_1 \alpha_2})$ are shown in **Table 1**, where $\alpha_1 = 0$ and $\alpha_2 = 0$ is called $M\hat{S}E$ of mean per unit, $\alpha_1 = -1$ and $\alpha_2 = -1$ is called $M\hat{S}E$ of multivariate ratio estimator, $\alpha_1 = -1$ and $\alpha_2 = -1$ is called $M\hat{S}E$ of multivariate ratio estimator, $\alpha_1 = 1$ and $\alpha_2 = 1$ is called $M\hat{S}E$ of multivariate product estimator, $\alpha_1 = -1$ and $\alpha_2 = 0$ is called $M\hat{S}E$ of ratio estimator using $x$, $\alpha_1 = 0$ and $\alpha_2 = -1$ is called $M\hat{S}E$ of ratio estimator using $z$, $\alpha_1 = 1$ and $\alpha_2 = 0$ is called $M\hat{S}E$ of product estimator using $x$ and $\alpha_1 = 0$ and $\alpha_2 = 1$ is called $M\hat{S}E$ of product estimator using $z$.

## 5. Conclusion

Stratified adaptive cluster sampling is an efficient method for sampling rare and hidden clustered populations. The numerical study showed that if the variable of in-

| Stratum1 | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 2 | 24 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 22 | 5 | 4 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 2 | 0 | 4 | 8 | 0 | 0 | 0 | 0 | 0 | 33 | 0 | 0 | 0 | 27 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 6 | 7 | 1 | 0 | 5 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 21 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 5 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 5 | 7 | 0 | 7 | 7 | 6 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 5 | 4 | 3 | 0 | 5 | 8 | 4 | 5 | 1 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 7 | 65 | 0 | 4 | 5 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 3 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 4 | 5 | 0 | 7 | 3 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 27 | 0 | 0 | 21 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 29 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

*(Header row spanning: Stratum1 | Stratum 2 | Stratum 3 | Stratum 4)*

**Figure 2. *Y* values.**

| Stratum1 | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 11 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 11 | 2 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 2 | 4 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 0 | 15 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 2 | 3 | 0 | 0 | 2 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 2 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 2 | 0 | 3 | 3 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 2 | 2 | 1 | 0 | 2 | 3 | 2 | 2 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 3 | 18 | 0 | 2 | 2 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 2 | 2 | 0 | 3 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 12 | 0 | 0 | 12 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 27 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

*(Header row spanning: Stratum1 | Stratum 2 | Stratum 3 | Stratum 4)*

**Figure 3. *X* values.**

| Stratum1 | | | | | Stratum 2 | | | | | Stratum 3 | | | | | Stratum 4 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 0 | 12 | 77 | 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 57 | 10 | 8 | 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 10 | 0 | 10 | 8 | 0 | 0 | 0 | 0 | 0 | 55 | 1 | 0 | 0 | 97 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 17 | 6 | 6 | 1 | 0 | 9 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 74 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 11 | 0 | 10 | 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 12 | 12 | 0 | 10 | 10 | 12 | 12 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 6 | 8 | 12 | 0 | 12 | 8 | 12 | 10 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 8 | 53 | 0 | 8 | 6 | 0 | 8 | 0 | 0 | 0 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 14 | 10 | 0 | 8 | 12 | 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 59 | 0 | 0 | 54 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 63 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**Figure 4. *Z* values.**

**Table 1. The estimate mean square error of the estimators.**

| $n_h$ | $n$ | $M\hat{S}E(\bar{y}_{(00)})$ | $M\hat{S}E(\bar{y}_{(10)})$ | $M\hat{S}E(\bar{y}_{(-10)})$ | $M\hat{S}E(\bar{y}_{(01)})$ | $M\hat{S}E(\bar{y}_{(0-1)})$ | $M\hat{S}E(\bar{y}_{(11)})$ | $M\hat{S}E(\bar{y}_{(-1-1)})$ | $M\hat{S}E(\bar{y}_{(\alpha_1^*\alpha_2^*)})$ |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 20 | 0.8105 | 8.2450 | 0.1341 | 7.3278 | 0.1589 | 79.6045 | 5.7501 | 1.9343 |
| 10 | 40 | 0.3212 | 2.1277 | 0.0397 | 1.8104 | 0.0427 | 10.0623 | 2.7507 | 0.9560 |
| 15 | 60 | 0.1935 | 1.3305 | 0.0239 | 0.9954 | 0.0189 | 4.7206 | 0.6095 | 0.6242 |
| 20 | 80 | 0.1104 | 0.5740 | 0.0204 | 0.4883 | 0.0146 | 1.6832 | 0.4179 | 0.3891 |
| 25 | 100 | 0.0744 | 0.4013 | 0.0190 | 0.3362 | 0.0130 | 1.0909 | 0.2390 | 0.2796 |
| 30 | 120 | 0.0619 | 0.3207 | 0.1154 | 0.2825 | 0.0081 | 0.8702 | 0.1537 | 0.2221 |
| 40 | 160 | 0.0442 | 0.2242 | 0.0105 | 0.1927 | 0.0049 | 0.5527 | 0.1194 | 0.1595 |

terest $(y)$ and the auxiliary variables $(x, z)$ have high positive correlation then the estimate of the mean square error of the ratio estimators is less than the estimate of the mean square error of the product estimator. The estimators which use only one auxiliary variable were better than the estimators which use two auxiliary variables.

## 6. Acknowledgements

## REFERENCES

[1]   S. K. Thompson, "Adaptive Cluster Sampling," *Journal of the American Statistical Association*, Vol. 85, No. 412, 1990, pp. 1050-1059. doi:10.1080/01621459.1990.10474975

[2]   S. K. Thompson and G. A. F. Seber, "Adaptive Sampling," Wiley, New York, 1996.

[3]   W. A. Abu-Dayyeh, M. S. Ahmed, R. A. Ahmed and H. A. Muttlak, "Some Estimators of a Finite Population Mean Using Auxiliary Information," *Applied Mathematics and Computation*, Vol. 139, No. 2-3, 2003, pp. 287-298. doi:10.1016/S0096-3003(02)00180-7

[4]   N. Chutiman and B. Kumphon, "Ratio Estimator Using Two Auxiliary Variables for Adaptive Cluster Sampling," *Journal of the Thai Statistical Association*, Vol. 6, No. 2, 2008, pp. 241-256.

*OJS*