# Some Methods of Model-Based Sampling

# Survey & Health Policy Research Center
## Technical Report

## June 2023

dongguk UNIVERSITY SHPRC

# Some Methods of Model-Based Sampling

**Sun Woong Kim[1], Steven G. Heeringa[2], Sung Joon Hong[3],
So Hyung Park[4], Hong Yup Ahn[5]**

[1]*Survey & Health Policy Research Center and Department of Statistics, Dongguk
University, 26, 3Ga, Pil-Dong, Jung-Gu, Seoul, South Korea 100-715
E-mail: sunwk@dongguk.edu*
[2]*Institute for Social Research, University of Michigan, 426 Thompson Street, Ann Arbor,
MI 48106 USA
E-mail: sheering@isr.umich.edu*
[3]*Program in Survey & Data Science, Dongguk University, Seoul, South Korea
E-mail: hsj8129@naver.com*
[4]*Department of Statistics, Dongguk University, Seoul, South Korea
E-mail: 12astro@naver.com*
[5]*Department of Statistics, Dongguk University, Seoul, South Korea
E-mail: ahn@dongguk.edu*

## Summary

With respect to commonly used $\pi PS$ sampling techniques, samplers are often interested in reducing the design variance of the Horvitz and Thompson estimator. We first describe the differences between the mechanisms of conventional design-based $\pi PS$ sampling methods of Mizuno and Brewer and two model-based $\pi PS$ sampling methods developed by Kim *et al*. We also suggest two new model-based $\pi PS$ sampling methods, and empirically compare the efficiency of the new methods to the previous model-based sampling methods and design-based $\pi PS$ and non-$\pi PS$ sampling methods. The case where the sample size is two as in "two per stratum" designs is of particular interest for empirical comparison. In regard to the design variance, model-based $\pi PS$ sampling methods are preferable to design-based $\pi PS$ sampling methods. One of the new methods performs best. This new method is comparable to the method of Murthy, which is a design-based non-$\pi PS$ sampling procedure. Moreover, model-based sampling methods are preferable to design-based sampling methods due to the flexibility in the choice of sampling design that has better stability of the variance estimator.

*Key words*: superpopulation regression model; average variance; optimization; design variance; stability of variance estimator; maximum likelihood; restricted maximum likelihood

**Résumé**

Parmi techniques d'échantillonnage couramment utilisées, les échantillonneurs sont souvent intéressés par la réduction de la variance du plan d'échantillonnage de l'estimateur de Horvitz et Thompson. Nous décrivons d'abord les différences entre les mécanismes de conception classique fondée sur les méthodes d'échantillonnage de Mizuno et Brewer et deux méthodes d'échantillonnage fondées sur un modèle mis au point par Kim et autres. Nous suggérons également deux nouvellesaux approches d'échantillonnage fondées sur un modéle, et réalisons une comparaisons empirique de l'efficacité des nouvelles méthodes par rapport aux méthodes d'échantillonnage basées sur un modèle et les approches fondées sur un plan non-échantillonnale. La comparaison empirique revêt un intérêt particulier lorsque la taille de l'échantillon est de deux. S'agissant de la variance du plan d'échantillonnage, les méthodes d'échantillonnage basées sur un modèle sont préférables aux méthodes fondées sur un plan d'échantillonnage. Une des nouvelles méthodes se distingue par sa performance. Cette nouvelle méthode est comparable à celle méthode Murthy, fondée sur un plan non-échantillonnale. En outre, les méthodes d'échantillonnage fondées sur un modèle sont préférables aux approches fondées sur un plan en raison de la flexibilité dans le choix du plan d'échantillonnage pour assurer une meilleure stabilité de l'estimateur de la variance.

# 1. Introduction

Since Hansen and Hurwitz (1943) first suggested the selection of primary sampling units from each stratum with probabilities proportional to size (PPS), a large number of techniques for sampling without replacement with unequal probabilities have been developed.

As discussed by Brewer and Hanif (1983) and Särndal (1996), much research on sample selection has been focused on design-based inclusion probability proportional to size ($\pi PS$) sampling procedures in which the second-order inclusion probabilities (or joint probabilities), which indicate the probabilities that any two units in a population are both included in a sample, have a key role in the variance reduction. For example, the methods of Mizuno (1952) and Brewer (1963) that are draw-by-draw procedures and Sampford's (1967) method, a rejective procedure, are well known $\pi PS$ sampling and commonly employed by samplers. The method of Mizuno (1952) is available in the R package, and the methods of Brewer (1963) and Sampford (1967) are available in the R package and software such as SAS or SPSS. See "sampling" package of R on the Web site of the Comprehensive R Archive Network (CRAN) (https://cran.r-project.org/web/packages/sampling/index.html) and "Complex Samples" in PASW Statistics (formerly SPSS Statistics) and SAS/STAT (2023). Murthy's (1957) method, a non-$\pi PS$ draw-by-draw procedure, was noted by Rao and Bayless (1969) and Cochran (1977). Murthy's method is available in SAS or SPSS for selecting samples in "two per stratum" designs.

The comparative efficiency of these techniques in actual population sampling applications is an open question. As seen in many empirical studies, this may be due to the fact that the variances of the estimates of interest calculated from a sample selected by any sampling method are sensitive to population characteristics, and hence the user

may not be sure that the efficiency of a chosen method would be significantly better than other procedures. This is especially true when a small sample is selected from a population or population stratum. In many national surveys, deep stratification with a substantial number of strata is used, and only a small number of cluster units are sampled from each stratum. For example, two per stratum designs are common in stratified cluster sampling (e.g., see Wu and Thompson, 2020, pp. 71-73). Accordingly, a sampling method whose efficiency is robust in the case of small samples would be preferred.

Although a sample is selected from a finite population, considering the concept of an infinite superpopulation may be useful in the sample selection stage. In fact, an infinite superpopulation model has been often used in the estimation procedures, such as model-assisted estimation and model-dependent estimation. But with regard to sample selection the model has been used by many writers mainly for the theoretical comparisons among sampling procedures, not for the actual selection of a sample.

The model may be used to ensure that the second-order inclusion probabilities involving sampling design implemented by a $\pi PS$ sampling procedure would result in reasonable efficiency. Kim $et\ al$. (2006) developed a theory of model-based $\pi PS$ sampling procedures as a specification of the selection method using the model. Their procedures to yield optimal sampling designs that reduce the variance of the Horvitz and Thompson (1952) estimator were based on fairly practical linear superpopulation models and optimization theory.

In this paper, we first describe the mechanism of design-based $\pi PS$ sampling of Mizuno (1952) and Brewer (1963) and model-based $\pi PS$ sampling developed by Kim $et\ al$. (2006). Next, we describe new model-based $\pi PS$ sampling methods, and empirically compare their efficiency to that for the previous model-based sampling methods and the conventional design-based sampling methods of Mizuno (1952), Brewer (1963) and Murthy (1957). For estimation of the parameters of the superpopulation model in model-

based approaches, maximum likelihood estimation and restricted maximum likelihood estimation are used. The case where the sample size is two is of particular interest for empirical comparison, both for simplicity and because it is the most important situation in practice. The model in model-based sampling is not be central to the selection problem and is just a means to the end of achieving higher efficiency.

## 2. Mechanism of Design-Based $\pi PS$ Sampling

Before presenting model-based $\pi PS$ sampling procedures in the next section, we first describe design-based $\pi PS$ sampling procedures of Mizuno (1952), Brewer (1963), and Sampford (1967).

Consider a finite population of $N$ units, denoted by $U = \{u_1, \cdots, u_i, \cdots, u_N\}$. $y_i$ is the value of the variable of interest, $y$, for the $i$ th unit $u_i$. In order to estimate the total $Y = \sum_{i=1}^{N} y_i$, a sample $s$ of size $n$ is selected from the finite population. Let $S$ be the set of all possible samples. $p(s)$ ($s \in S$), simply called the sampling design (or sampling plan), is the probability of selecting $s$. In particular, let $p_d(s)$ denote the sampling design in design-based $\pi PS$ sampling. Let the $\pi_i$ be the first-order inclusion probabilities, denoted by $\pi_i = \sum_{i \in s} p_d(s)$, and let the $\pi_{ij}$ be the second-order inclusion probabilities given by $\pi_{ij} = \sum_{i,j \in s} p_d(s)$. When $n = 2$, simply $\pi_{ij} = p_d(s)$.

The method of Mizuno (1952) for $n \geq 2$ uses the selection procedure:

   i. Select the first unit $u_i$ with unequal probabilities, $p_i$, where $p_i = x_i / X$, $X = \sum_{i=1}^{N} x_i$,

   and $x_i$ is the value of the auxiliary variable, $x$, correlated with $y$.

   ii. Select the remaining units with equal probabilities.

The $\pi_i$ and $\pi_{ij}$ are respectively:

$$\pi_i = p_i + (1 - p_i)\frac{n-1}{N-1} \tag{1}$$

$$\pi_{ij} = \left( (p_i + p_j)\frac{(n-1)(N-n)}{(N-1)(N-2)} \right) + \frac{(n-1)(n-2)}{(N-1)(N-2)} \tag{2}$$

For $n = 2$,

$$\pi_{ij} = p_d(s) = \frac{1}{N-1}(p_i + p_j) \tag{3}$$

In this case the sampling design is a simple function of $p_i$ and $p_j$.

The method of Brewer (1963), which is only for $n = 2$ and every $p_i < 1/2$, has the more complicated procedure, as follows:

    i. Select the first unit $u_i$ with probabilities $\frac{p_i(1-p_i)}{Q(1-2p_i)}$, where $Q = \frac{1}{2}\left(1 + \sum_{i=1}^{N}\frac{p_i}{1-2p_i}\right)$.

    ii. Select the second unit $u_j$ with probabilities $\frac{p_j}{1-p_i}$.

The procedure gives

$$\pi_i = 2p_i \tag{4}$$

$$\pi_{ij} = p_d(s) = \frac{2p_ip_j}{Q}\frac{(1-p_i-p_j)}{(1-2p_i)(1-2p_j)} \tag{5}$$

The method of Sampford (1967) is an extension of Brewer's (1963) method to samples of any size. Design-based $\pi PS$ sampling procedures of Brewer or Sampford have the properties:

a) The sampling method is a draw-by-draw procedure, where one unit is selected at each successive draw.

b) The sampling design $p_d(s)$, which keeps $\pi_i = \sum_{i \in s} p_d(s) = np_i$ called the $\pi PS$ requirement, is obtained according to the selection probability of each unit defined for each draw, and is a function of the relative sizes $p_i$ of the units.

A sampler may prefer a $\pi PS$ sampling yielding a smaller design variance. But as described above, the $p_d(s)$ in design-based $\pi PS$ sampling is a certain function of the relative sizes $p_i$ depending on only the values of the auxiliary variable $x$, and there is no definite indication of the strength and direction of a linear relationship between the variables $x$ and $y$. Thus, although $p_d(s)$ plays a central role in the reduction of the design variance, it is not clear whether $p_d(s)$ in any design-based $\pi PS$ sampling procedure would yield a low variance for any population of interest.

## 3. Mechanism of Model-Based $\pi PS$ Sampling

A generalized regression (GREG) estimator may be one of the useful estimators for the population total. But it is well-known that it might be appreciably biased for a small sample, although the bias is in modest for large samples. As an alternative, the Horvitz-Thompson (H-T) estimator (1952) in (6), which is unbiased for the population total and highly efficient under a good $\pi PS$ sampling method, can be used.

$$\hat{Y}_{HT} = \sum_{i=1}^{n} \frac{y_i}{\pi_i} \tag{6}$$

The H-T estimator is the only unbiased estimator in the subclass of linear estimators denoted by

$$\hat{Y} = \sum_{i=1}^{n} k_i y_i \qquad (7)$$

where $k_i$ is a constant to be used as a weight for the $i$ th unit whenever it selected for the sample, and hence the best linear estimator of the subclass (Horvitz and Thompson (1952), Godambe (1955)). Also, note that best linear estimate does not exist for the entire class of linear estimators (Godambe (1955)).

The variance of the H-T estimator is

$$Var\left(\hat{Y}_{HT}\right) = \sum_{i=1}^{N} \frac{y_i^2}{\pi_i} + 2\sum_{i=1}^{N}\sum_{j>i}^{N} \frac{\pi_{ij}}{\pi_i \pi_j} y_i y_j - Y^2 \qquad (8)$$

A model-based $\pi PS$ sampling method was first suggested by Des Raj (1956). Let $p_m(s)$ be the sampling design in his model-based $\pi PS$ sampling, based on the model $m$ denoted by $y_i = \alpha + \beta x_i$ reflecting a linear relationship between the variables $x$ and $y$. His method for $n = 2$ is a variance minimization sampling procedure, which first constructs an optimization problem consisting of an objective function and constraints in terms of sampling design $p_m(s)$ for minimizing $Var\left(\hat{Y}_{HT}\right)$ in (8) under the model. The method then attempts to obtain an optimal set of $p_m(s)$ for all possible samples by linear programming (LP).

In brief, his $\pi PS$ sampling procedure has the properties:

c) Prior the sample selection, the sampling design $p_m(s) = \pi_{ij}$ for all possible samples is determined by LP. It meets the $\pi PS$ requirement, that is,

$$\pi_i = \sum_{i \in s} p_m(s) = np_i .$$

d) One selection using $p_m(s)$ samples the whole sample $s$. This is a whole sample procedure.

His sampling procedure is attractive with respect to the variance reduction achieved by using the model. But his model is unusual because there is no error term. Kim *et al.* (2006) developed a theory of model-based $\pi PS$ sampling procedures using an infinite superpopulation model. They assume that a finite population of $N$ units is drawn from an infinite superpopulation with the regression model $\xi$, given by

$$y_i = \alpha + \beta x_i + \varepsilon_i, \ i = 1, \cdots, N, \tag{9}$$

where $E_\xi(\varepsilon_i | x_i) = 0$, $Var_\xi(\varepsilon_i | x_i) = \delta x_i^\gamma$ ($\delta > 0$, $\gamma \geq 0$), and $E_\xi(\varepsilon_i, \varepsilon_j | x_i, x_j) = 0$, $i \neq j$.

$E_\xi$ and $Var_\xi$ respectively denote the expected value and variance under the model $\xi$. It is also assumed that the $\varepsilon_i$ are normally distributed.

Note that many writers often prefer the model without the intercept for the purpose of the simplicity of theoretical comparison between sampling procedures, while the model in (9) has the intercept for the practical use.

The variance of the H-T estimator, given by Horvitz and Thompson (1952), is

$$Var\left(\hat{Y}_{HT}\right) = \sum_{i=1}^{N} \frac{y_i^2(1-\pi_i)}{\pi_i} + 2\sum_{i=1}^{N}\sum_{j>i}^{N} \frac{\pi_{ij}}{\pi_i \pi_j} y_i y_j - 2\sum_{i=1}^{N}\sum_{j>i}^{N} y_i y_j \tag{10}$$

A different expression on the variance of the H-T estimator, suggested by Yates and Grundy (1953), is

$$Var\left(\hat{Y}_{HT}\right) = \sum_{i=1}^{N}\sum_{j>i}^{N} \left(\pi_i \pi_j - \pi_{ij}\right) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j}\right)^2 \tag{11}$$

With respect to inference, the anticipated variance (ANV), introduced by Isaki and Fuller (1982), is used as a measure describing the variability between the total and the estimator of the total under both the sampling design and superpopulation model. If the H-T estimator is used, it simply becomes the average variance (AV), that is, the model expectation of the design variance expressed as

9

$$E_\xi E_p \left[ \left( \hat{Y}_{HT} - Y \right)^2 \right] = E_\xi \left[ Var \left( \hat{Y}_{HT} \right) \right],$$ (12)

where $E_p$ denotes the expected value under the sampling design, and both $Y$ and $\hat{Y}_{HT}$ are random variables.

Let $p_\xi(s)$ be the sampling design under model-based $\pi PS$ sampling of Kim *et al.* (2006) using the regression model $\xi$. They showed that in cases of $n = 2$, an optimal sampling design $p_\xi(s)$ in a set of possible $\pi PS$ sampling designs that minimize the AV in (12) can be obtained by using one of the following optimization problems:

*Minimize*

$$\sum_{i=1}^{N} \sum_{j>i}^{N} \frac{\alpha + \beta(x_i + x_j)}{x_i x_j} p_\xi(s),$$ (13)

or

*Minimize*

$$\sum_{i=1}^{N} \sum_{j>i}^{N} \left( \frac{1}{x_j} - \frac{1}{x_i} \right) \left( \alpha \frac{1}{x_i} + \beta \right) p_\xi(s),$$ (14)

subject to the linear equality constraints

$$\sum_{i \in s} p_\xi(s) = \pi_i , \quad i = 1, \cdots, N$$ (15)

Note that the two objective functions in (13) and (14) are induced from the expressions of $Var \left( \hat{Y}_{HT} \right)$ in (10) and (11), respectively, and hence a different form of $Var \left( \hat{Y}_{HT} \right)$ may yield a different optimization problem. We call these two optimization problems composed of (13) and (15), and (14) and (15), OP1 and OP2, respectively.

## 4. New Model-Based $\pi PS$ Sampling

We continue to focus on the design stage for the actual selection of a sample from a finite population, rather than the estimation stage. In other words, although the H-T estimator does not involve the superpopulation model $\xi$, we assume the model, and seek to find $p_\xi(s)$ to reduce $Var\left(\hat{Y}_{HT}\right)$ for the finite population as well as $E_\xi\left[Var\left(\hat{Y}_{HT}\right)\right]$ for the infinite population.

Here we first derive objective functions different from those in (13) or (14), and then construct different optimization problems by adding (15) and additional constraints, as seen later.

**THEOREM 1.** With the variance formula in (10), the AV on the H-T estimator under the superpopulation model in (9) is

$$\sum_{i=1}^{N}(X/nx_i-1)\left(\delta x_i^\gamma+\alpha^2+2\alpha\beta x_i+\beta^2 x_i^2\right)-2\sum_{i}^{N}\sum_{j>i}^{N}\left(\alpha^2+\alpha\beta(x_i+x_j)+\beta^2 x_i x_j\right)$$

$$+\frac{2\alpha^2 X^2}{n^2}\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j}\pi_{ij}+2\alpha\beta X\frac{n-1}{n}N+\beta^2 X^2\frac{n-1}{n} \tag{16}$$

***Proof.*** Consider the form of the variance of the H-T estimator in (10). Since it is $\pi PS$ sampling, $\pi_i=np_i$. Then from the first and third terms in (10) under the superpopulation model, we have

$$E_\xi\left[\sum_{i=1}^{N}\frac{y_i^2(1-\pi_i)}{\pi_i}-2\sum_{i=1}^{N}\sum_{j>i}^{N}y_i y_j\right]$$

11

$$= \sum_{i=1}^{N}(X/nx_i - 1)E_\xi(y_i^2) - 2\sum_{i=1}^{N}\sum_{j>i}^{N}E_\xi(y_i y_j)$$

$$= \sum_{i=1}^{N}(X/nx_i - 1)\left(\delta x_i^\gamma + \alpha^2 + 2\alpha\beta x_i + \beta^2 x_i^2\right) - 2\sum_{i}^{N}\sum_{j>i}^{N}\left(\alpha^2 + \alpha\beta(x_i + x_j) + \beta^2 x_i x_j\right)$$

(17)

For the second term in (10), we have

$$E_\xi\left[2\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{\pi_{ij}}{\pi_i\pi_j}y_i y_j\right] = \frac{2X^2}{n^2}\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{E_\xi(y_i y_j)}{x_i x_j}\pi_{ij}$$

$$= \frac{2X^2}{n^2}\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{\alpha^2 + \alpha\beta(x_i + x_j) + \beta^2 x_i x_j}{x_i x_j}\pi_{ij}$$

$$= \frac{2\alpha^2 X^2}{n^2}\sum_{i}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j}\pi_{ij}$$

$$+ \frac{2\alpha\beta X^2}{n^2}\sum_{i}^{N}\sum_{j>i}^{N}\frac{x_i + x_j}{x_i x_j}\pi_{ij}$$

(18)

$$+ \beta^2 X^2\frac{n-1}{n}$$

When the second term in (18) is expanded, it gives

$$\frac{2\alpha\beta X^2}{n^2}\left[\sum_{i}^{N}\sum_{j>i}^{N}\frac{1}{x_j}\pi_{ij} + \sum_{i}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij}\right] = \frac{\alpha\beta X^2}{n^2}\left[\sum_{j}^{N}\frac{1}{x_j}\sum_{i\neq j}\pi_{ij} + \sum_{i}^{N}\frac{1}{x_i}\sum_{j\neq i}\pi_{ij}\right]$$

$$= \frac{\alpha\beta X^2}{n^2}\left[\sum_{j}^{N}\frac{1}{x_j}(n-1)\pi_j + \sum_{i}^{N}\frac{1}{x_i}(n-1)\pi_i\right]$$

$$= 2\alpha\beta X\frac{n-1}{n}N$$

(19)

12

This completes the proof.

**COROLLARY 1.** With the variance formula in (10), the AV on the H-T estimator under the superpopulation model with $\alpha = 0$ in (9) does not depend on $\pi_{ij}$, and is fixed as

$$\sum_{i=1}^{N}(X/nx_i - 1)\left(\delta x_i^{\gamma} + \beta^2 x_i^2\right) - 2\beta^2 \sum_{i}^{N}\sum_{j>i}^{N} x_i x_j + \beta^2 X^2 \frac{n-1}{n} \tag{20}$$

**COROLLARY 2.** If the superpopulation model in (9) is assumed, the minimization of the AV on the H-T estimator given in (16) is equivalent to minimizing

$$\sum_{i}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j}\sum_{i,j \in s} p_{\xi}(s) \tag{21}$$

**Proof.** In (16) only the third term depends on $\pi_{ij}$, while the other terms do not depend on $\pi_{ij}$, and are fixed. Thus, the minimization of the AV amounts to minimizing (21).

**REMARK 1.** (21) does not depend on $\alpha$, $\beta$, $\delta$, and $\gamma$, and it is a linear function of $p_{\xi}(s)$.

**COROLLARY 3.** In cases of $n = 2$, the minimization of the AV on the H-T estimator is equivalent to minimizing

$$\sum_{i}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j} p_{\xi}(s) \tag{22}$$

**REMARK 2.** As given in (22), in cases of $n = 2$, the minimization of the AV on the H-T estimator is reduced to minimizing a simple linear function of $p_{\xi}(s)$. (22) is the same

function as Des Raj (1956) induced to minimize $Var\left(\hat{Y}_{HT}\right)$ under the assumption that $y_i = \alpha + \beta x_i$ without the error term. See page 198, Des Raj (1956).

**RESULT 1.** Based on (22), in cases of $n = 2$, a simple optimization problem to find model-based $\pi PS$ sampling design $p_\xi(s)$, called OP3, can be given by:

*Minimize*

$$\sum_{i}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j}p_\xi(s) \tag{23}$$

subject to the linear equality constraints

$$\sum_{i\in s}p_\xi(s) = \pi_i, \ i = 1,\cdots,N \tag{24}$$

Now we obtain a different AV by using a variance form different from (10).

**THEOREM 2.** Using the variance expression in (11), the AV on the H-T estimator under the superpopulation model in (9) is

$$\frac{\delta X}{n}\sum_{i=1}^{N}x_i^{\gamma-1} - \delta\sum_{i=1}^{N}x_i^{\gamma} - 2\alpha\left[\sum_{i=1}^{N}\sum_{j>i}^{N}\left(x_i - x_j\right)\left(\alpha x_i^{-1} + \beta\right)\right]$$

$$+\frac{2\alpha X^2}{n^2}\left[\alpha\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j}\pi_{ij} - \alpha\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i^2}\pi_{ij} + \frac{\beta Nn(n-1)}{X} - 2\beta\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij}\right] \tag{25}$$

***Proof.*** For (11), we may write

$$Var\left(\hat{Y}_{HT}\right) = \sum_{i=1}^{N}\sum_{j>i}^{N}\left(p_i p_j - \frac{\pi_{ij}}{n^2}\right)\left(\frac{y_i}{p_i} - \frac{y_j}{p_j}\right)^2 \tag{26}$$

14

Since

$$E_{\xi}\left[\frac{y_i}{p_i}-\frac{y_j}{p_j}\right]^2 = \frac{2}{p_i^2}(\delta x_i^{\gamma}+\alpha^2+\beta^2 x_i^2+2\alpha\beta x_i) - \frac{2}{p_i p_j}(\alpha^2+\alpha\beta(x_i+x_j)+\beta^2 x_i x_j)$$

$$= 2\delta X^{\gamma}p_i^{\gamma-2}+2\alpha X^2\frac{x_j-x_i}{x_i x_j}\left(\alpha x_i^{-1}+\beta\right), \tag{27}$$

we have

$$E_{\xi}\left[\sum_{i=1}^{N}\sum_{j>i}^{N}\left(p_i p_j-\frac{\pi_{ij}}{n^2}\right)\left(\frac{y_i}{p_i}-\frac{y_j}{p_j}\right)^2\right]$$

$$= 2\delta X^{\gamma}\sum_{i=1}^{N}\sum_{j>i}^{N}p_i^{\gamma-2}\left(p_i p_j-\frac{\pi_{ij}}{n^2}\right)$$

$$+ 2\alpha X^2\left[\sum_{i=1}^{N}\sum_{j>i}^{N}\left(p_i p_j-\frac{\pi_{ij}}{n^2}\right)\frac{x_j-x_i}{x_i x_j}\left(\alpha x_i^{-1}+\beta\right)\right]$$

$$= \frac{\delta X^{\gamma}}{n}\sum_{i=1}^{N}(1-np_i)p_i^{\gamma-1} + 2\alpha X^2\left[\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{x_i x_j}{X^2}\frac{x_j-x_i}{x_i x_j}\left(\alpha x_i^{-1}+\beta\right)\right]$$

$$- \frac{2\alpha X^2}{n^2}\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{x_j-x_i}{x_i x_j}\left(\alpha x_i^{-1}+\beta\right)\pi_{ij}$$

$$= \frac{\delta X}{n}\sum_{i=1}^{N}x_i^{\gamma-1}-\delta\sum_{i=1}^{N}x_i^{\gamma} - 2\alpha\left(\sum_{i=1}^{N}\sum_{j>i}^{N}\left(x_i-x_j\right)\left(\alpha x_i^{-1}+\beta\right)\right)$$

$$+ \frac{2\alpha X^2}{n^2}\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{x_i-x_j}{x_i x_j}\left(\alpha x_i^{-1}+\beta\right)\pi_{ij} \tag{28}$$

The last term in (28) can be written in the form

$$\frac{2\alpha X^2}{n^2}\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{x_i-x_j}{x_i x_j}\left(\alpha x_i^{-1}+\beta\right)\pi_{ij} = \frac{2\alpha X^2}{n^2}\left[\alpha\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{x_i-x_j}{x_i^2 x_j}\pi_{ij}+\beta\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{x_i-x_j}{x_i x_j}\pi_{ij}\right]$$

$$= \frac{\alpha X^2}{n^2}\left[2\alpha\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j}\pi_{ij}-2\alpha\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i^2}\pi_{ij}+\beta\left(2\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_j}\pi_{ij}-2\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij}\right)\right] \tag{29}$$

15

Also,

$$\beta\left(2\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_j}\pi_{ij} - 2\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij}\right) = \beta\left(2\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_j}\pi_{ij} + 2\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij} - 4\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij}\right)$$

$$= \beta\left(\sum_{i}^{N}\frac{1}{x_i}\sum_{j\neq i}^{N}\pi_{ij} + \sum_{j}^{N}\frac{1}{x_j}\sum_{i\neq j}^{N}\pi_{ij} - 4\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij}\right)$$

$$= 2\beta\left(\frac{Nn(n-1)}{X} - 2\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij}\right) \tag{30}$$

This completes the proof.

**COROLLARY 4.** Under the superpopulation model with $\alpha = 0$ in (9), (25) reduces to

$$\frac{\delta X}{n}\sum_{i=1}^{N}x_i^{\gamma-1} - \delta\sum_{i=1}^{N}x_i^{\gamma} \tag{31}$$

**REMARK 3.** (31) is different from (20), due to the different expressions for the variance of the H-T estimator.

**COROLLARY 5.** Under the superpopulation model in (9), minimizing the AV on the H-T estimator given in (25) amounts to minimizing

$$\alpha\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j}\pi_{ij} - \alpha\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i^2}\pi_{ij} - 2\beta\sum_{i=1}^{N}\sum_{j>i}^{N}\frac{1}{x_i}\pi_{ij}, \tag{32}$$

where $\pi_{ij} = \sum_{i,j\in s}p_\xi(s)$.

**REMARK 4.** (32) depends on $\alpha$ and $\beta$, and it is a linear function of $p_\xi(s)$.

**COROLLARY 6.** In cases of $n = 2$, the minimization of (32) reduces to minimizing

$$\sum_{i}^{N}\sum_{j>i}^{N}\left[\alpha\left(\frac{1}{x_i x_j}-\frac{1}{x_i^2}\right)-2\beta\frac{1}{x_i}\right]p_\xi(s) \tag{33}$$

**RESULT 2.** The different optimization problem, called OP4, to obtain a model-based $\pi PS$ sampling design $p_\xi(s)$, for the case of $n=2$, is given by:

*Minimize*

$$\sum_{i}^{N}\sum_{j>i}^{N}\left[\alpha\left(\frac{1}{x_i x_j}-\frac{1}{x_i^2}\right)-2\beta\frac{1}{x_i}\right]p_\xi(s) \tag{34}$$

subject to

$$\sum_{i\in s}p_\xi(s)=\pi_i,\ i=1,\cdots,N. \tag{35}$$

**REMARK 5.** In addition to (35), the linear inequality constraints (36) can be basically added

$$0<p_\xi(s)\le\pi_i\pi_j,\ j>i=1,\cdots,N, \tag{36}$$

since the well-known variance estimator $\widehat{Var}\left(\hat{Y}_{HT}\right)$ in (37), given by Yates and Grundy (1953) and by Sen (1953) from (11), is defined if $\pi_{ij}>\mathbf{0}$, and nonnegative if $\pi_i\pi_j\ge\pi_{ij}$.

$$\widehat{Var}\left(\hat{Y}_{HT}\right)=\sum_{i=1}^{n}\sum_{j>i}^{n}\frac{\pi_i\pi_j-\pi_{ij}}{\pi_{ij}}\left(\frac{y_i}{\pi_i}-\frac{y_j}{\pi_j}\right)^2, \tag{37}$$

Also, (38) can replace (36).

$$c\pi_i\pi_j\le p_\xi(s)\le\pi_i\pi_j,\ j>i=1,\cdots,N, \tag{38}$$

where $0<c<1$.

Note that the stability of the variance estimator in (37) may be improved if $c$ in (38) is sufficiently far from 0, as discussed by Hanurav (1967), Nigam, Kumar and Gupta (1984), and Rao and Nigam (1992). Thus, the larger value of $c$ is preferred.

Since $\pi_i = 2p_i$, (36) and (38) can be respectively expressed in forms

$$0 < p_\xi(s) \leq \frac{4}{X^2} x_i x_j, \ \ j > i = 1, \cdots, N \tag{39}$$

$$\frac{4c}{X^2} x_i x_j \leq p_\xi(s) \leq \frac{4}{X^2} x_i x_j, \ \ j > i = 1, \cdots, N \tag{40}$$

The constraints in (39) or (40) can be added to OP1, OP2, and OP3, as in OP4.

## 5. Estimation of Model Parameters

Except for OP3 in Result 1, which is the simplest problem among four optimization problems, in order to solve optimization problems such as OP4 in Result 2 as well as OP1 and OP2 on sampling design $p_\xi(s)$, the estimation of model parameters of not only $\alpha$ and $\beta$ but also $\delta$ and $\gamma$ is essential. Two approaches for the estimation can be used: (i) Maximum Likelihood (ML) Estimation, and (ii) Restricted Maximum Likelihood (REML) Estimation. For the method of ML, as discussed by Godfrey, Roshwalb and Wright (1984), and Särndal and Wright (1984), the Harvey's (1976) algorithm can be used. His algorithm uses the ordinary least squares (OLS) estimates as the starting values for the regression coefficients $\alpha$ and $\beta$, and in each iteration the values of $\alpha$ and $\beta$ depend on $\delta$ and $\gamma$, or the reverse. The REML estimation was developed by Patterson and Thompson (1971), and Harville (1977). The values of $\alpha$ and $\beta$ only depend on $\delta$ and $\gamma$. Harvey's algorithm for ML estimation can be easily implemented by direct programming. Both ML and REML estimation methods based on particular iterative algorithms are also available in statistical software programs, for example, PROC

MIXED of SAS and "nlme" package of R (http://cran.r-project.org). In the next section, we will empirically compare the estimated values of model parameters using ML and REML.
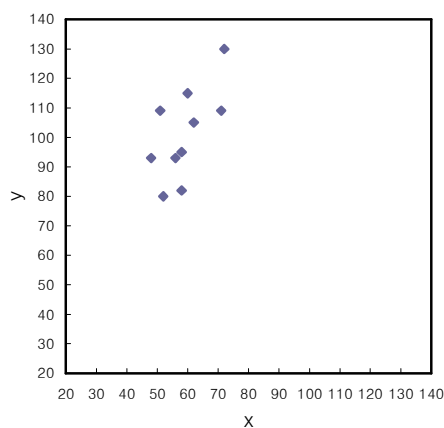
## 6. Empirical Study

The previous model-based $\pi PS$ sampling methods, OP1 and OP2, the suggested model-based $\pi PS$ sampling methods, OP3 and OP4, and the conventional design-based sampling methods of Mizuno (1952), Brewer (1963) and Murthy (1957) were compared for the case of $n = 2$. The comparison used 18 small natural populations described in the paper of Rao and Bayless (1969) and summarized in Table 1. There were originally 20 populations in their paper, but 2 populations (numbered 6 and 8 in their paper) were excluded because the linear model in (9) above was not successfully applied.

**Table 1** *Description of the natural populations.*

| No | Source | $y$ | $x$ | $N$ |
|---|---|---|---|---|
| 1 | Horvitz and Thompson (1952) | No. of Households | Eye-estimated no. of Households | 20 |
| 2 | Des Raj (1965) | No. of Households | Eye-estimated no. of Households | 20 |
| 3 | Rao (1963) | Corn acreage in 1960 | Corn acreage in 1958 | 14 |
| 4 | Kish (1965) | No. of rented dwelling units | Total no. of dwelling units | 10 |
| 5 | Kish (1965) | No. of rented dwelling units | Total no. of dwelling units | 10 |
| 6 | Hanurav (1967) | Population in 1967 | Population in 1957 | 20 |
| 7 | Hanurav (1967) | Population in 1967 | Population in 1957 | 16 |
| 8 | Hanurav (1967) | Population in 1967 | Population in 1957 | 17 |
| 9 | Cochran (1963) | No. of persons per block | No. of rooms per block | 10 |
| 10 | Cochran (1963) | No. of people in 1930 | No. of people in 1920 | 16 |
| 11 | Cochran (1963) | No. of people in 1930 | No. of people in 1920 | 16 |
| 12 | Cochran (1963) | No. of people in 1930 | No. of people in 1920 | 17 |
| 13 | Sukhatme (1954) | No. of wheat acres in 1937 | No. of wheat acres in 1936 | 10 |
| 14 | Sukhatme (1954) | No. of wheat acres in 1937 | No. of wheat acres in 1936 | 10 |
| 15 | Sampford (1962) | Oats acreage in 1957 | Oats acreage in 1947 | 35 |

19

| 16 | Sukhatme (1954) | Wheat acreage | No. of villages | 20 |
| 17 | Sukhatme (1954) | Wheat acreage | No. of villages | 20 |
| 18 | Sukhatme (1954) | Wheat acreage | No. of villages | 9 |

As an illustration, Figure 1 shows a $x - y$ scatter plot of elements for population 9 in Table 1. The superpopulation model in (9) may be applied and the parameters of the model can be estimated by ML estimation or REML estimation.



**Figure 1.** *Scatter plot of population 9,*
*Cochran (1963)*

For the 18 populations in Table 1, Table 2 shows the estimated values of the four parameters $\alpha$, $\beta$, $\delta$, and $\gamma$ in the superpopulation model (9) when using ML estimation and REML estimation. If one judges "the estimates from ML estimation and REML estimation are the same" when they coincide by the second digit of each value, of 18 populations, 8 (44%) for $\alpha$, 15(83%) for $\beta$, 4(22%) for $\delta$, and 3(17%) for $\gamma$ are the same. Therefore, ML estimation and REML estimation in these populations tend to give differing estimates with the exception of $\beta$. Also, although Rao and Bayless (1969) assumed the superpopulation model (9) with $\alpha = 0$ for the empirical comparison between

20

**Table 2** *Comparison of estimates of parameters by ML estimation and REML estimation.*

| Population | $\alpha$ ML | $\alpha$ REML | $\beta$ ML | $\beta$ REML | $\sigma^2$ ML | $\sigma^2$ REML | $\gamma$ ML | $\gamma$ REML |
|---|---|---|---|---|---|---|---|---|
| 1 | 1.1425# | 1.1128# | 1.0381# | 1.0402# | 0.0038# | 0.0052# | 2.7462 | 2.6700 |
| 2 | 0.6482 | 0.7158 | 1.1086# | 1.1032# | 0.0016# | 0.0031# | 3.4297 | 3.2270 |
| 3 | 25.9265# | 25.9040# | 1.0272# | 1.0275# | 0.1988 | 0.2746 | 1.6100 | 1.5724 |
| 4 | -0.5254# | -0.5887# | 0.5056# | 0.5478# | 0.0202# | 0.0515# | 2.2976 | 1.9877 |
| 5 | -0.8129# | -0.8188# | 0.5951# | 0.5948# | 0.1321 | 0.2054 | 1.4111 | 1.3227 |
| 6 | 185718.5700 | 176849.3000 | 1015.3631# | 1022.8110# | 0.5370 | 21.6865 | 3.2817 | 2.7914 |
| 7 | 71558.7950 | 73130.1200 | 1284.6306# | 1283.0830# | 1701.0314 | 4176.2820 | 2.2624 | 2.1570 |
| 8 | 12177.1530# | 12146.5600# | 1264.0638# | 1264.6180# | 0.0463 | 128.7146 | 3.5951 | 2.4377 |
| 9 | 21.0523 | 23.3574 | 1.3594# | 1.3213# | 70351.356 | 5664.530 | -1.5640 | -0.8888 |
| 10 | -0.7908 | 28.9407 | 1.1797 | 0.7786 | 400.9461 | 0.0011 | -0.1311 | 2.7342 |
| 11 | 5.7650 | 7.5682 | 1.2257# | 1.2164# | 1251.5417 | 937.8721 | -0.1378 | -0.0338 |
| 12 | 18.3981# | 18.6041# | 1.0601# | 1.0589# | 1112.3704 | 865.6584 | -0.1974# | -0.1084# |
| 13 | -16.4124# | -16.1364# | 1.0562# | 1.0547# | 0.0801 | 0.1600 | 1.6751 | 1.5885 |
| 14 | 0.8136 | 0.9784 | 0.9665# | 0.9651# | 5.9355 | 15.2708 | 0.8118 | 0.6624 |
| 15 | 2.3077# | 2.3128# | 0.2313# | 0.2313# | 0.0092# | 0.0100# | 1.9390# | 1.9347# |
| 16 | 220.7632 | 225.4950 | 221.2283 | 219.6092 | 52872.4720 | 64098.1500 | 1.1941# | 1.1171# |
| 17 | 131.3070 | 128.5208 | 258.6386# | 259.7728# | 27433.806 | 34193.550 | 1.3894 | 1.2814 |
| 18 | 627.8519 | 655.1650 | 283.4577 | 275.7315 | 278243.9400 | 317033.9000 | -0.4755 | -0.3158 |

# Coincided by the second digit between ML estimate and REML estimate

21

five sampling methods including the methods of Brewer and Murthy, most estimates of $\alpha$ are not close to zero, as shown in Table 2. In addition, though $\gamma \geq 0$ in the model (9), the estimates of $\gamma$ lie in the interval (-1.6, 3.6), which includes the negative values, and more than half of the estimates of $\gamma$ do not lie in the interval (0, 2) or (1, 2) often assumed in many references.

Figure 2 may be helpful to understand an appreciable difference between model-based sampling and design-based sampling, although only results for the model-based sampling using OP3 are presented for a population. This figure, which is for population 9, shows a comparison of the distribution of sampling design $p(s)$ (i.e., $p_\xi(s)$ or $p_d(s)$) by $x_i$ and $x_j$, the values of the auxiliary variable, and the corresponding variances for model-based sampling using OP3 with $c = 0$, 0.1, 0.2, 0.3, 0.4, 0.5 and the conventional design-based sampling methods of Mizuno, Brewer, and Murthy. The model-based sampling design $p_\xi(s)$ was obtained from OP3 consisting of (41), (42) and (43) or (44), and "LP procedure" (or "OPTMODEL procedure") in SAS/OR was used to find the solution to $p_\xi(s)$. OP3 was infeasible for the cases with $c = 0.6$, 0.7, 0.8, and 0.9.

*Minimize*

$$\sum_{i}^{N}\sum_{j>i}^{N}\frac{1}{x_i x_j}p_\xi(s) \tag{41}$$

subject to

$$\sum_{i \in s} p_\xi(s) = 2p_i, \ i = 1, \cdots, N \tag{42}$$

and for $0 < c < 1$,

$$\frac{4c}{X^2}x_i x_j \leq p_\xi(s) \leq \frac{4}{X^2}x_i x_j, \ j > i = 1, \cdots, N \tag{43}$$

or for $c = 0$,

$$0 < p_\xi(s) \le \frac{4}{X^2} x_i x_j, \quad j > i = 1, \cdots, N \tag{44}$$

The design-based sampling design $p_d(s)$ for the methods of Mizuno and Brewer was calculated by (3) and (5), respectively. The $p_d(s)$ for Murthy's method were computed as:
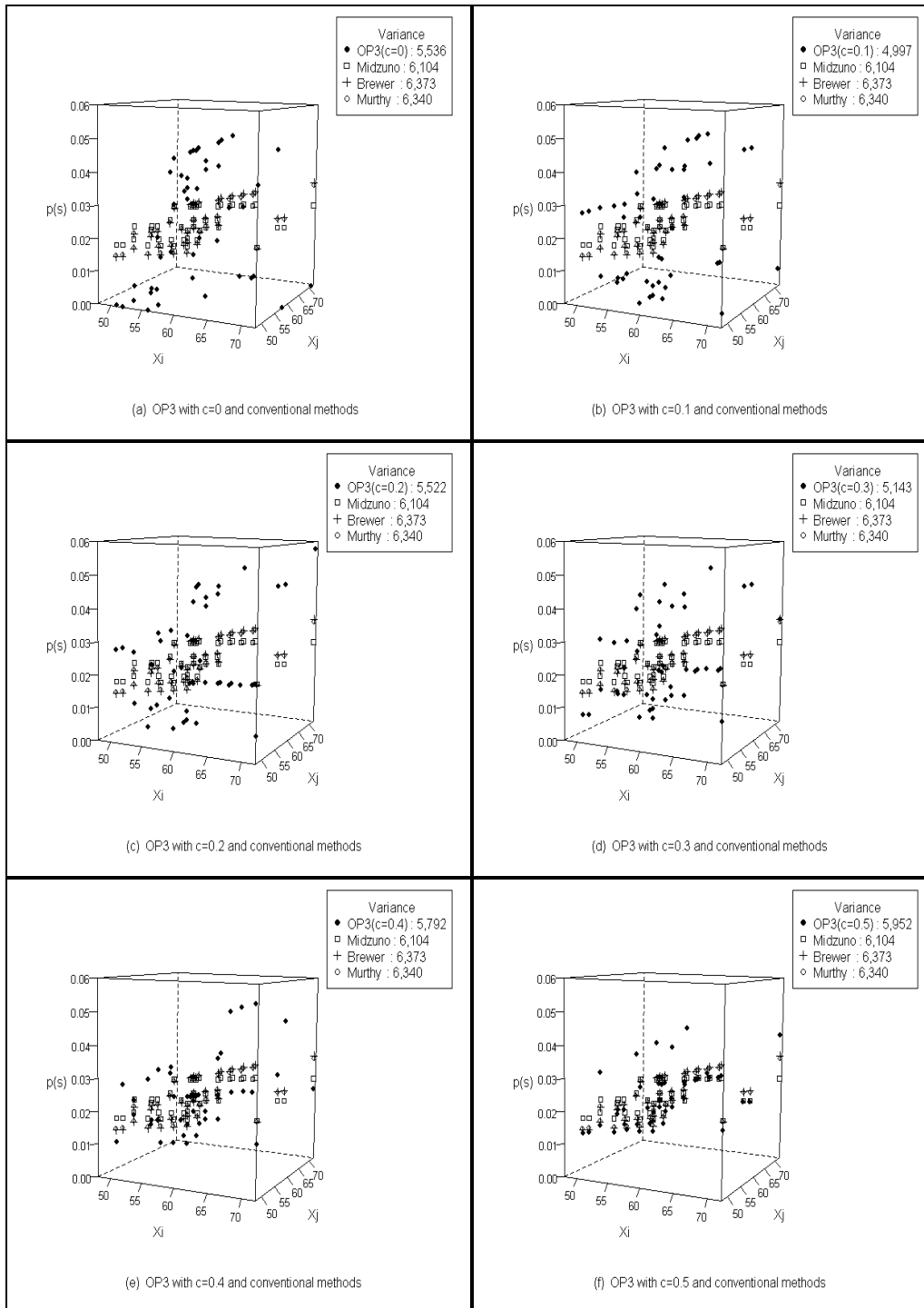
$$p_d(s) = \frac{p_i p_j (2 - p_i - p_j)}{(1 - p_i)(1 - p_j)} \tag{45}$$

The variance reported in the box located in the upper right hand corner of each panel in the figure is interpreted as follows. "OP3 ($c = 0$): 5,536," indicates that the value of the variance $Var\left(\hat{Y}_{HT}\right)$ based on sampling design $p_\xi(s)$ obtained from the model-based sampling using OP3 with $c = 0$ is 5,536. "Mizuno: 6,104" and "Brewer: 6,373" respectively denote the value of $Var\left(\hat{Y}_{HT}\right)$ calculated using sampling design $p_d(s)$ from the methods of Mizuno and Brewer. Also, "Murthy: 6,340" denotes the value of the variance $Var_M\left(\hat{Y}\right)$ in (46) calculated using $p_d(s)$ from the method of Murthy.

$$Var_M\left(\hat{Y}\right) = \sum_{i=1}^{N} \sum_{j>i}^{N} \frac{p_i p_j (1 - p_i - p_j)}{2 - p_i - p_j} \left( \frac{y_i}{p_i} - \frac{y_j}{p_j} \right)^2 \tag{46}$$

As shown in Figure 2, there is a clear difference between the model-based and design-based sampling methods. First, the distribution of $p(s)$ (i.e., $p_\xi(s)$) for the model-based sampling method varies by the value of $c$. Also, as seen in panels (a), (b), (c), and (d), $p(s)$ (i.e., $p_\xi(s)$) from model-based sampling with $c = 0, 0.1, 0.2, 0.3$ are scattered over a wide range according to the values of $x_i$ and $x_j$, while $p(s)$ (i.e., $p_d(s)$) from the methods of Mizuno, Brewer, and Murthy tend to concentrate in a small range, in spite of

**Figure 2.** *Comparison of sampling designs by the values of the auxiliary variable and the corresponding variances between model-based sampling method using OP3 with a different value of c and three conventional design-based sampling methods*; *those obtained from conventional methods are repeatedly shown in each panel for the convenience of comparison.*

the changes of those of $x_i$ and $x_j$. This causes a smaller variance for model-based sampling and a larger variance for design-based sampling, as seen in the values of the variance in the boxes of panels. In contrast, as in panels (e) and (f), the spread of $p(s)$ (i.e., $p_\xi(s)$) from model-based sampling with $c = 0.4$ and 0.5 and that from the design-based sampling methods are more similar, yielding more equal variances between them.

When comparing the six panels in Figure 2, it seems that as regards the value of $c$, there is a trade-off between the reduction of the variance and the stability of the variance estimator. The larger value of $c$ indicates the larger stability of the variance estimator, as noted in (38). However, when the value of $c$ is relatively low, as in panels (a), (b), (c), and (d), $p(s)$ (i.e., $p_\xi(s)$) obtained from model-based sampling method using OP3 tend to be dispersed, resulting in a large reduction in variance, compared to the cases where $c = 0.4$ or $c = 0.5$.

Note that with respect to any value of $c$, model-based sampling using OP3 gives a smaller variance than the three design-based sampling methods. Also, it is flexible in terms of $c$. If one pursues the larger variance reduction rather than the stability of the variance estimator, using a lower value of $c$ may be appropriate. But if we prefer the stability of the variance estimator, a higher value of $c$ can be used, but for the price is the larger variance. Anyway, it would offer an optimal sampling design under the chosen constraints on the value of $c$.

Next, we turn to Table 3, which shows the summary on results of empirical comparison on the relative efficiency (RE) for 18 populations for model-based sampling methods using OP1, OP2, OP3, and OP4 with $c = 0$, 0.1, 0.2, 0.3, 0.4, 0.5 and the three design-based sampling methods. The model-based sampling using OP1, OP2, and OP4 denote that only (41) is replaced by (13), (14), and (34), respectively, in OP3 consisting

of (41), (42) and (43) or (44). Note that those optimization problems were consistently infeasible for the cases with $c = 0.6$, 0.7, 0.8, and 0.9. The details on Table 3 are illustrated as follows:

For example, "OP1 M" in the table denotes OP1 consisting of the estimates of the model from ML estimation, while "OP1 R" indicates OP1 by the estimates from REML estimation. Here, the RE for model-based $\pi PS$ sampling is denoted by

$$RE_{\xi,\pi PS} = \left[ Var_{PPS}\left(\hat{Y}\right) \Big/ Var\left(\hat{Y}_{HT}\right) \right] \times \mathbf{100}, \tag{47}$$

where $Var_{PPS}\left(\hat{Y}\right) = \dfrac{\mathbf{1}}{n} \sum\limits_{i=\mathbf{1}}^{N} \sum\limits_{j>i}^{N} p_i p_j \left( \dfrac{y_i}{p_i} - \dfrac{y_j}{p_j} \right)^{\mathbf{2}}$ , which is the variance of the estimate of the population total under PPS sampling with replacement.

The REs for the design-based $\pi PS$ sampling methods of Mizuno or Brewer are also computed by (47), and with a distinction, the REs are denoted by $RE_{d,\pi PS}$ instead of $RE_{\xi,\pi PS}$. The RE for Murthy's method, which is a non-$\pi PS$ sampling method, is calculated by:

$$RE_M = \left[ Var_{PPS}\left(\hat{Y}\right) \Big/ Var_M\left(\hat{Y}\right) \right] \times \mathbf{100}. \tag{48}$$

According to the empirical study of Rao and Bayless (1969), for 18 populations, PPS sampling with replacement always had a larger variance than Brewer's method. Also, it is theoretically clear that $Var_M\left(\hat{Y}\right) < Var_{PPS}\left(\hat{Y}\right)$.

The frequencies in column "f" in the table denote the number of populations where

$$RE_{\xi,\pi PS} > RE_{d,\pi PS} \tag{49}$$

or

$$RE_{\xi,\pi PS} > RE_M \tag{50}$$

**Table 3** *Comparison of frequency of populations that model-based sampling shows a better efficiency than design-based sampling.*

| Design-based | Model-based | c=0 f | f1 | f2 | f3 | c=0.1 f | f1 | f2 | f3 | c=0.2 f | f1 | f2 | f3 | c=0.3 f | f1 | f2 | f3 | c=0.4 f | f1 | f2 | f3 | c=0.5 f | f1 | f2 | f3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mizuno | OP1M | 16** | 2 | 3 | 11 | 16** | 2 | 1 | 13 | 16** | 5 | 0 | 11 | 14** | 3 | 1 | 10 | 13** | 3 | 1 | 9 | 13** | 3 | 1 | 9 |
|  | OP1R | 16** | 0 | 4 | 12 | 16** | 3 | 1 | 12 | 16** | 4 | 1 | 11 | 15** | 4 | 1 | 10 | 13** | 3 | 1 | 9 | 13** | 2 | 2 | 9 |
|  | OP2M | 14** | 3 | 1 | 10 | 14** | 2 | 1 | 11 | 14** | 2 | 2 | 10 | 14** | 2 | 2 | 10 | 13** | 2 | 2 | 9 | 13** | 2 | 2 | 9 |
|  | OP2R | 14** | 3 | 1 | 10 | 14** | 2 | 1 | 11 | 14** | 2 | 2 | 10 | 14** | 2 | 2 | 10 | 13** | 2 | 2 | 9 | 13** | 2 | 2 | 9 |
|  | OP3 | 15** | 0 | 4 | 11 | 16** | 3 | 1 | 12 | 16** | 3 | 2 | 11 | 14** | 3 | 1 | 10 | 14** | 3 | 1 | 10 | 14** | 4 | 1 | 9 |
|  | OP4M | 14** | 3 | 1 | 10 | 13** | 2 | 1 | 10 | 14** | 3 | 1 | 10 | 14** | 3 | 1 | 10 | 13** | 2 | 2 | 9 | 13** | 2 | 2 | 9 |
|  | OP4R | 15** | 4 | 1 | 10 | 14** | 3 | 1 | 10 | 14** | 2 | 2 | 10 | 14** | 2 | 2 | 10 | 13** | 2 | 2 | 9 | 12** | 1 | 2 | 9 |
| Brewer | OP1M | 12** | 5 | 6 | 1 | 12** | 5 | 3 | 4 | 8 | 4 | 4 | 0 | 12** | 9 | 0 | 3 | 6 | 5 | 1 | 0 | 5 | 5 | 0 | 0 |
|  | OP1R | 12** | 5 | 5 | 2 | 10** | 3 | 3 | 4 | 8 | 3 | 5 | 0 | 8 | 4 | 1 | 3 | 7 | 6 | 1 | 0 | 6 | 6 | 0 | 0 |
|  | OP2M | 8 | 6 | 2 | 0 | 8 | 5 | 3 | 0 | 9* | 7 | 2 | 0 | 7 | 7 | 0 | 0 | 7 | 7 | 0 | 0 | 8 | 8 | 0 | 0 |
|  | OP2R | 8 | 6 | 1 | 1 | 6 | 4 | 2 | 0 | 10** | 9 | 1 | 0 | 8 | 8 | 0 | 0 | 9* | 9 | 0 | 0 | 7 | 7 | 0 | 0 |
|  | OP3 | 10** | 6 | 2 | 2 | 11** | 7 | 1 | 3 | 13** | 11 | 2 | 0 | 9* | 7 | 1 | 1 | 9* | 9 | 0 | 0 | 9* | 9 | 0 | 0 |
|  | OP4M | 10** | 8 | 2 | 0 | 11** | 9 | 2 | 0 | 8 | 6 | 2 | 0 | 8 | 8 | 0 | 0 | 7 | 7 | 0 | 0 | 7 | 7 | 0 | 0 |
|  | OP4R | 9* | 7 | 2 | 0 | 10** | 8 | 2 | 0 | 9* | 7 | 2 | 0 | 8 | 8 | 0 | 0 | 7 | 7 | 0 | 0 | 9* | 9 | 0 | 0 |
| Murthy | OP1M | 8 | 3 | 4 | 1 | 8 | 2 | 3 | 3 | 6 | 4 | 2 | 0 | 7 | 4 | 0 | 3 | 5 | 3 | 2 | 0 | 3 | 3 | 0 | 0 |
|  | OP1R | 9* | 3 | 4 | 2 | 8 | 2 | 3 | 3 | 7 | 3 | 4 | 0 | 7 | 3 | 1 | 3 | 5 | 3 | 2 | 0 | 4 | 4 | 0 | 0 |
|  | OP2M | 6 | 3 | 3 | 0 | 6 | 4 | 2 | 0 | 8 | 7 | 1 | 0 | 3 | 3 | 0 | 0 | 5 | 5 | 0 | 0 | 6 | 6 | 0 | 0 |
|  | OP2R | 6 | 3 | 3 | 0 | 5 | 4 | 1 | 0 | 8 | 7 | 1 | 0 | 3 | 3 | 0 | 0 | 5 | 5 | 0 | 0 | 6 | 6 | 0 | 0 |
|  | OP3 | 9* | 5 | 2 | 2 | 10** | 7 | 1 | 2 | 8 | 6 | 2 | 0 | 7 | 6 | 0 | 1 | 9* | 7 | 2 | 0 | 6 | 6 | 0 | 0 |
|  | OP4M | 7 | 5 | 2 | 0 | 8 | 7 | 1 | 0 | 5 | 3 | 2 | 0 | 4 | 4 | 0 | 0 | 4 | 4 | 0 | 0 | 7 | 7 | 0 | 0 |
|  | OP4R | 7 | 5 | 2 | 0 | 8 | 7 | 1 | 0 | 6 | 4 | 2 | 0 | 3 | 3 | 0 | 0 | 4 | 4 | 0 | 0 | 7 | 7 | 0 | 0 |

*exactly half of 18 populations, **over half of 18 populations

For example, the first "16" in terms of "OP1 M" and "Mizuno" in the column of "f" in the table indicates that of 18 populations, 16 populations satisfy (49). More specifically, for 16 populations, the REs for model-based sampling using "OP1 M" are larger than in design-based sampling of Mizuno, whereas for 2 populations they are smaller.

The frequencies in "f1," "f2," and "f3" in the table respectively denote the number of populations that are

$$0 < RE_{\xi,\pi PS} - RE_{d,\pi PS} \leq 10 \tag{51}$$

or

$$0 < RE_{\xi,\pi PS} - RE_M \leq 10, \tag{52}$$

$$11 \leq RE_{\xi,\pi PS} - RE_{d,\pi PS} \leq 20 \tag{53}$$

or

$$11 \leq RE_{\xi,\pi PS} - RE_M \leq 20, \tag{54}$$

and

$$RE_{\xi,\pi PS} - RE_{d,\pi PS} \geq 21 \tag{55}$$

or

$$RE_{\xi,\pi PS} - RE_M \geq 21. \tag{56}$$

Here, (51) or (52), (53) or (54), and (55) or (56) denote that the REs on model-based sampling are respectively "slightly better," "much better," and "very much better," than those on design-based sampling. Note that f = f1 + f2 + f3. For example, the first "2" in the column of "f1," the first "3" in "f2," and the first "11" in "f3" in the table indicates that of 16 populations in "f," 2 populations satisfy (51), 3 populations do (53), and 11 populations do (55).

The findings from Table 3 are summarized as follows:

(1) Model-based sampling methods (using OP1 M, OP1 R, OP2 M, OP2 R, OP3, OP4 M, and OP4 R) are consistently more efficient than Mizuno's method, regardless of the value of $c$. For at least half of 18 populations, they show "very much better" efficiency.

(2) When the value of $c$ is low, model-based sampling methods are overall more efficient relative to Brewer's method. For some populations, when the value of $c$ is low, the methods using OP1 or OP3 show "very much better" efficiency. Taken overall, model-based sampling using OP3 shows a better efficiency than the other model-based methods.

(3) Model-based sampling method using OP3 compares favorably with the method of Murthy, when the value of $c$ is low, and for some populations, it has "very much better" efficiency as well as "much better" efficiency. Other model-based sampling methods are less efficient than the one of Murthy.

(4) As presented in Table 2, ML estimation and REML estimation give different estimates of the model in (9), and it seems that model-based methods using optimization problems involving these different estimates of the model may yield different efficiencies.

(5) For model-based sampling methods, there is a trade-off between the reduction of variance and the stability of the variance estimator because the REs tend to be reduced as the value of $c$ is increased.


## 7. Conclusion Remarks

We have suggested two model-based $\pi PS$ sampling strategies using the optimization problems of OP3 and OP4. The method using OP3 is empirically preferable to the method using OP4, as well as the previous methods using OP1 and OP2. Compared to

others, OP3 is the simpler optimization problem, and it does not depend on the parameters in the superpopulation model.

Those four model-based $\pi PS$ sampling methods are flexible in terms of the choice of sampling design because one may choose the value of $c$, which is related to the stability of variance estimator. But one should be careful in choosing the value, since there is a trade-off between the variance reduction and the stability of the variance estimator. With regard to the efficiency, regardless of the value of $c$, the model-based methods are shown empirically to be superior to design-based $\pi PS$ sampling of Mizuno, and when the value of $c$ is low, they are preferable to the one of Brewer. Also, in such a case, the method using OP3 is comparable to the method of Murthy.

There are several issues for a future study. First, in this paper, we assumed only one superpopulation model, which may be appropriate for some populations, but may be not so for the others. For example, as seen in Figure 2, it seems that the model was suitable for the population, because model-based sampling was working well for the variance reduction as well as the stability of the variance estimator, compared to the conventional design-based sampling methods including Murthy's method. But there might be certain populations where a different superpopulation model is required. For example, a polynomial model might be adopted to improve the efficiency of model-based sampling. For such a model, we need to develop different optimization problems. Second, we should note that it might not be feasible to solve a chosen optimization problem. In such cases, a different model assumption should be pursued, likewise. Third, a study on the efficiency of model-based sampling methods in for larger sample sizes should be conducted. In addition, the comparison of the efficiency of the H-T estimator under the model-based sampling and the GREG estimator in the conventional sampling method might be another interesting issue.

# References

Brewer, K. R. W. (1963). A model of systematic sampling with unequal probabilities. *Aust. J. Stat.*, **5**, 5-13.

Brewer, K. R. W. and Hanif, M. (1983). *Sampling with Unequal Probabilities*. New York: Springer-Verlag.

Cochran, W.G. (1963). *Sampling Techniques*. 2nd ed. New York: Wiley.

Cochran, W.G. (1977). *Sampling Techniques*. 3rd ed. New York: Wiley.

Des Raj (1956). A note on the determination of optimum probabilities in sampling without replacement. *Sankhyā*, **17**, 197-200.

Des Raj (1965). Variance estimation in randomized systematic sampling with probability proportional to size. *J. Amer. Statist. Assoc.*, **60**, 278-284.

Godambe, V. P. (1955). A unified theory of sampling from finite populations. *J. R. Stat. Soc. B Stat. Methodol.*, **17**, 269-278.

Godfrey, J., Roshwalb, A., and Wright, R. L. (1984). Model-based stratification in inventory cost estimation. *J. Bus. Econom. Statist.*, **2**, 1-9.

Hansen, M. H., and Hurwitz, W. N. (1943). On the theory of sampling from finite populations. *Annals of Mathematical Statistics*, **14**, 333-362.

Hanurav, T. V. (1967). Optimum utilization of auxiliary information: $\pi PS$ sampling of two units from a stratum. *J. R. Stat. Soc. B Stat. Methodol.*, **29**, 374-391.

Harvey, A. C. (1976). Estimating regression models with multiplicative heteroscedasticity. *Econometrica*, **44**, 461-465.

Harville, D. A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. *J. Amer. Statist. Assoc.*, **72**, 320-340.

Horvitz, D. G. and Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe. *J. Amer. Statist. Assoc.*, **47**, 663-685.

Isaki, C. T. and Fuller, W. A. (1982). Survey design under the regression superpopulation model. *J. Amer. Statist. Assoc.*, **77**, 89-96.

Kim, S. W., Heeringa, S. G., and Solenberger, P. W. (2006). Model-based sampling designs for optimum estimation. In *JSM Proceedings*, *Survey Research Methods Section, American Statistical Association*, 3245-3252.

Kish, L. (1965). *Survey Sampling*. New York: Wiley.

Mizuno, H. (1952). On the sampling system with probability proportional to sum of sizes. *Ann. Inst. Statist. Math.*, **3**, 99-107.

Murthy, M. N. (1957). Ordered and unordered estimators in sampling without replacement. *Sankhyā*, **18**, 379-390.

Nigam, A. K., Kumar, P., and Gupta, V. K. (1984). Some methods of inclusion probability proportional to size sampling. *J. R. Stat. Soc. B Stat. Methodol.*, **46**, 564-571.

Patterson, H. D., and Thompson, R. (1971). Recovery of interblock information when block sizes are unequal. *Biometrika*, **58**, 545-554.

Rao, J. N. K. (1963). On three procedures of unequal probability sampling without replacement. *J. Amer. Statist. Assoc.*, **58**, 202-215.

Rao, J. N. K. and Bayless, D. L. (1969). An empirical study of the stabilities of estimators and variance estimators in unequal probability sampling of two units per stratum. *J. Amer. Statist. Assoc.*, **64**, 540-559.

Rao, J. N. K. and Nigam, A. K. (1992). Optimal controlled sampling: a unified approach. *International Statistical Review*, **60**, 89-98.

Sampford, M. R. (1962). *An Introduction to Sampling Theory*. Edinburgh and London: Oliver and Boyd Ltd.

Sampford, M. R. (1967). On sampling without replacement with unequal probabilities of selection. *Biometrika*, 54, 499-513.

Särndal, C. E. (1996). Efficient estimators with simple variance in unequal probability sampling. *J. Amer. Statist. Assoc.*, 91, 1289-1300.

Särndal, C. E. and Wright, R. L. (1984). Cosmetic form of estimators in survey sampling. *Scand. J. Statist.*, 11, 146-156.

SAS/OR (2023). *Documentation*. Retrieved from https://support.sas.com/en/software/sas-or-support.html

SAS/STAT (2023). *Sample Survey Design and Analysis*. Retrieved from https://support.sas.com/rnd/app/stat/topics/survey-analysis.html

Sukhatme, P. V. (1954). *Sampling Theory of Surveys with Applications*. Ames, Iowa State College Press.

Wu, C., and Thompson, M. E. (2020). Sampling Theory and Practice. Cham, Springer International Publishing.