

CHAPTER 10

Two-Phase Sampling

10.1 INTRODUCTION

It has been explained in Chapters 5, 8, and 9 how one can use auxiliary information at the stages of estimation of parameters, selection of samples, stratification, and their combinations to increase the precision of the estimators. To use auxiliary information for estimation of the population characteristics through ratio and regression methods, one needs to know only the population total of the auxiliary variable. For using the auxiliary variable as a measure of size for the selection of sample, one needs to know the values of the auxiliary variable at unit level. For the purpose of stratification using an auxiliary variable, one needs to know the values of the auxiliary variable for a large number of units so that after formation of the strata, each unit of the population can be assigned to one of the strata. In some situations the auxiliary information is not available in advance but can be made available at little extra cost. For example, in household income and expenditure surveys, one can easily collect information regarding the household size at the time of listing the households by putting little or no additional cost.

Two-phase sampling, also known as double sampling, was introduced by Neyman (1938). In two-phase sampling, samples are selected in phases. In the first phase, a relatively large sample s' of size n' is selected by some sampling design p' , and only information on the auxiliary variable is collected. During the second phase, a sample s of size n is selected either from s' , the sample selected in the first phase, or from the entire population by using a suitable sampling design p , and information regarding the study and auxiliary variable is collected. Evidently, two-phase sampling is useful if the auxiliary information x is relatively easy and cheaper to collect than the study variable y as well as if the strength of the relationship between the variables x and y is high. If the sample is selected in more than two phases, the resulting sampling design is called a multiphase sampling. In this chapter we will consider the use of the first-phase sample at the stages of estimation of the parameter, stratification, selection of sample, and their combinations in the second phase.

10.2 TWO-PHASE SAMPLING FOR ESTIMATION

In the first phase, a sample s' of size n' be selected from a population with probability $p(s')$ using a sampling design p and the information of the auxiliary variable x is obtained only. Let

$$t(s', x) = \sum_{i \in s'} b_{s'i} x_i \quad (10.2.1)$$

be an unbiased estimator for the population total $X \left(= \sum_{i=1}^N x_i \right)$ of the auxiliary variable x based on the first-phase sample s' and $b_{s'i}$'s are known constants satisfying the unbiasedness condition

$$\sum_{s' \supset i} b_{s'i} p'(s') = 1 \quad \forall i = 1, \dots, N \quad (10.2.2)$$

In the second phase, a subsample s of size n is selected from s' and an unbiased estimator of the total $Y \left(= \sum_{i=1}^N y_i \right)$ is obtained as

$$t(s, y) = \sum_{i \in s} b_{si} y_i \quad (10.2.3)$$

Here the weights b_{si} 's are suitably chosen to satisfy

$$E\{t(s, y) | s'\} = \sum_{i \in s'} b_{s'i} y_i \quad (10.2.4)$$

The condition (Eq. 10.2.4) requires

$$\sum_{s \supset i} b_{si} p(s) = b_{s'i} \quad \forall i \in s' \quad (10.2.5)$$

10.2.1 Difference Method of Estimation

A difference estimator for the total Y is given by

$$\hat{Y}_{dt} = t(s, y) - \lambda \{t(s, x) - t(s', x)\} \quad (10.2.6)$$

where $t(s, x) = \sum_{i \in s} b_{si} x_i$ and λ is a suitably chosen constant.

Clearly, \hat{Y}_{dt} is unbiased for Y for any chosen value of λ because $E\{t(s, y)\} = E[E\{t(s, y) | s'\}] = E\{t(s', y)\} = Y$ and $E\{t(s, x)\} = E\{t(s', x)\} = X$.

The variance of \hat{Y}_{dt} is given by

$$V(\hat{Y}_{dt}) = V\{E(\hat{Y}_{dt}|s')\} + E\{V(\hat{Y}_{dt}|s')\} \quad (10.2.7)$$

$$= V\{t(s', \gamma)\} + E[V\{t(s, d_\lambda)|s'\}]$$

where $t(s', \gamma) = \sum_{i \in s'} b_{s'i} \gamma_i$, $t(s, d_\lambda) = \sum_{i \in s} b_{si} d_{\lambda i}$, and $d_{\lambda i} = \gamma_i - \lambda x_i$.

An unbiased estimator of $V(\hat{Y}_{dt})$ is given by

$$\hat{V}(\hat{Y}_{dt}) = \hat{V}\{t(s', \gamma)\} + \hat{V}\{t(s, d_\lambda)|s'\} \quad (10.2.8)$$

where $\hat{V}\{t(s', \gamma)\}$ and $\hat{V}\{t(s, d_\lambda)|s'\}$ are unbiased estimators of $V\{t(s', \gamma)\}$ and $V\{t(s, d_\lambda)|s'\}$, respectively.

10.2.1.1 Arbitrary Sampling Design

Let the initial sample s' be selected by a varying probability sampling scheme with $\pi_i (> 0)$ and $\pi_{ij} (> 0)$ as the inclusion probabilities of the i th, and i th and j th units, respectively. The subsample s is selected from s' by simple random sampling without replacement (SRSWOR) method. In this case, if

we choose, $t(s', x) = \sum_{i \in s'} \frac{x_i}{\pi_i}$ and $t(s, \gamma) = \frac{n'}{n} \sum_{i \in s} \frac{\gamma_i}{\pi_i}$, we get

$$\hat{Y}_{dt} = \frac{n'}{n} \sum_{i \in s} \frac{\gamma_i}{\pi_i} - \lambda \left(\frac{n'}{n} \sum_{i \in s} \frac{x_i}{\pi_i} - \sum_{i \in s'} \frac{x_i}{\pi_i} \right) \quad (10.2.9)$$

$$V\{E(\hat{Y}_{dt}|s')\} = V\{t(s', \gamma)\}$$

$$= V\left(\sum_{i \in s'} \frac{\gamma_i}{\pi_i}\right) \quad (10.2.10)$$

$$= \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{\gamma_i}{\pi_i} - \frac{\gamma_j}{\pi_j} \right)^2$$

and

$$\begin{aligned} E[V\{t(s, d_\lambda)|s'\}] &= n' \frac{(n' - n)}{n(n' - 1)} E \left[\sum_{i \in s'} \frac{d_{\lambda i}^2}{\pi_i^2} - n' \left(\frac{1}{n'} \sum_{i \in s'} \frac{d_{\lambda i}}{\pi_i} \right)^2 \right] \\ &= n' \frac{(n' - n)}{n(n' - 1)} \left[\left(\sum_{i \in U} \frac{d_{\lambda i}^2}{\pi_i} - \frac{D_\lambda^2}{n'} \right) - \frac{1}{n'} \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{d_{\lambda i}}{\pi_i} - \frac{d_{\lambda j}}{\pi_j} \right)^2 \right] \end{aligned} \quad (10.2.11)$$

where $d_{\lambda i} = y_i - \lambda_i x_i$ and $D_\lambda = Y - \lambda X$.

Unbiased estimators of $V\{t(s', y)\}$ and $V(t(s, d_\lambda)|s')$ are, respectively, given by

$$\widehat{V}\{t(s', y)\} = \frac{n'}{n} \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \left(\frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \right) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 \quad (10.2.12)$$

and

$$\widehat{V}\{t(s, d_\lambda)|s'\} = n' \frac{(n' - n)}{n(n - 1)} \left[\sum_{i \in s} \frac{d_{\lambda i}^2}{\pi_i^2} - \frac{1}{n} \left(\sum_{i \in s} \frac{d_{\lambda i}}{\pi_i} \right)^2 \right] \quad (10.2.13)$$

The expressions from Eq. (10.2.10) to Eq. (10.2.13) yield the variance of \widehat{Y}_{dt} and its unbiased estimators as follows:

$$\begin{aligned} V(\widehat{Y}_{dt}) &= \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 + n' \frac{(n' - n)}{n(n' - 1)} \\ &\quad \times \left[\left(\sum_{i \in U} \frac{d_{\lambda i}^2}{\pi_i} - \frac{D_\lambda^2}{n'} \right) - \frac{1}{n'} \left\{ \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{d_{\lambda i}}{\pi_i} - \frac{d_{\lambda j}}{\pi_j} \right)^2 \right\} \right] \end{aligned} \quad (10.2.14)$$

and

$$\begin{aligned} \widehat{V}(\widehat{Y}_{dt}) &= \frac{n'}{n} \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \left(\frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \right) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 \\ &\quad + n' \frac{(n' - n)}{n(n - 1)} \sum_{i \in s} \left(\frac{d_{\lambda i}}{\pi_i} - \frac{1}{n} \sum_{i \in s} \frac{d_{\lambda i}}{\pi_i} \right)^2 \end{aligned} \quad (10.2.15)$$

10.2.1.2 Simple Random Sampling Without Replacement

In case, both the samples s' and s are selected by SRSWOR method, we get $\pi_i = n'/N$ and $\pi_{ij} = n'(n' - 1)/\{N(N - 1)\}$. In this situation the expressions (Eqs. 10.2.9, 10.2.14, and 10.2.15) reduce to

$$\widehat{Y}_{dt} = N [\bar{y}(s) - \lambda(\bar{x}(s) - \bar{x}(s'))] \quad (10.2.16)$$

$$V(\widehat{Y}_{dt}) = N^2 \left[\left(\frac{N - n'}{Nn'} \right) S_y^2 + \frac{(n' - n)}{nn'} S_{d\lambda}^2 \right] \quad (10.2.17)$$

and

$$V(\widehat{Y}_{dt}) = N^2 \left[\left(\frac{N - n'}{Nn'} \right) s_y^2 + \frac{(n' - n)}{nn'} s_{d\lambda}^2 \right] \quad (10.2.18)$$

$$\begin{aligned}
\text{where } \bar{x}(s') &= \sum_{i \in s'} x_i/n', \quad \bar{x}(s) = \sum_{i \in s} x_i/n, \quad \bar{y}(s) = \sum_{i \in s} y_i/n, \\
S_y^2 &= \sum_{i \in U} (y_i - \bar{Y})^2 / (N - 1), \quad S_{d\lambda}^2 = \sum_{i \in U} (d_{\lambda i} - \bar{D}_\lambda)^2 / (N - 1), \\
s_y^2 &= \sum_{i \in s} (y_i - \bar{y}_s)^2 / (n - 1), \quad s_{d\lambda}^2 = \sum_{i \in s} (d_{\lambda i} - \bar{d}_{\lambda s})^2 / (n - 1), \\
\bar{Y} &= Y/N, \quad \bar{D}_\lambda = D_\lambda/N, \quad \text{and } \bar{d}_{\lambda s} = \sum_{i \in s} d_{\lambda i}/n.
\end{aligned}$$

10.2.1.3 Efficiency Under Simple Random Sampling Without Replacement

Let us consider the following simple cost function

$$C = c'n' + cn \quad (10.2.19)$$

where C is the total cost of the survey excluding the fixed cost, c' and c denote, respectively, the cost for collecting information on x and y from a unit.

The cost c is expected to be much higher than c' . If information regarding the study variable y is collected only through direct survey (without using two-phase sampling) keeping C as C_0 , then the sample size for the direct survey would be

$$n_0 = C_0/c = (c'/c)n' + n \quad (10.2.20)$$

The variance of the sample mean $\bar{y}_{dr} = \sum_{i \in s_0} y_i/n_0$ based on a sample s_0 of size n_0 selected by SRSWOR method is given by

$$V(\bar{y}_{dr}) = \left(\frac{1}{n_0} - \frac{1}{N} \right) S_y^2$$

Two-phase sampling will be more precise than direct sampling if

$$V(\hat{Y}_{dr}) < V(N\bar{y}_{dr})$$

$$\text{i.e., if } N^2 \left[\left(\frac{N - n'}{Nn'} \right) S_y^2 + \frac{(n' - n)}{nn'} \left(S_y^2 + \lambda^2 S_x^2 - 2\lambda\rho S_x S_y \right) \right] < N^2 \left(\frac{1}{n_0} - \frac{1}{N} \right) S_y^2$$

$$\begin{aligned}
\text{i.e., if } \rho > \frac{1}{2} \left[h + \left\{ h \left(1 - \frac{n}{n'} \right) \left(1 + \frac{nc}{n'c'} \right) \right\}^{-1} \right] \\
\quad (10.2.21)
\end{aligned}$$

where ρ is the correlation coefficient between x and y and $h = \lambda S_x / S_y$.

10.2.1.4 Probability Proportional to Size With Replacement Sampling

Let the initial sample s' be selected by probability proportional to size with replacement (PPSWR) method using z_i as a measure of size for the i th unit and the second phase sample s be selected from s' by SRSWOR method, treating all the units in s' are distinct. In this case

$$\hat{Y}_{dt} = \frac{1}{n} \sum_{i \in s} \frac{y_i}{q_i} - \lambda \left(\frac{1}{n} \sum_{i \in s} \frac{x_i}{q_i} - \frac{1}{n'} \sum_{i \in s'} \frac{y_i}{q_i} \right) \quad (10.2.22)$$

where $q_i = z_i/Z$ and $Z = \sum_{i \in U} z_i$ and $\sum_{i \in s}$ denotes the sum over the units in s with repetition. Clearly, \hat{Y}_{dt} in Eq. (10.2.22) is unbiased for Y and

$$\begin{aligned} V\{E(\hat{Y}_{dt}|s')\} &= V\left(\frac{1}{n'} \sum_{i \in s'} \frac{y_i}{q_i}\right) \\ &= \frac{1}{n'} \sum_{i \in U} q_i \left(\frac{y_i}{q_i} - Y\right)^2; \end{aligned} \quad (10.2.23)$$

$$\begin{aligned} E\{V(\hat{Y}_{dt}|s')\} &= \left(\frac{1}{n} - \frac{1}{n'}\right) \frac{1}{n' - 1} E \sum_{i \in s'} \left\{ \frac{d_{\lambda i}}{q_i} - \frac{1}{n'} \left(\sum_{i \in s'} \frac{d_{\lambda i}}{q_i} \right) \right\}^2 \\ &= \left(\frac{1}{n} - \frac{1}{n'}\right) \sum_{i \in U} q_i \left(\frac{d_{\lambda i}}{q_i} - D_{\lambda}\right)^2 \end{aligned} \quad (10.2.24)$$

Eqs. (10.2.23) and (10.2.24) yield

$$V(\hat{Y}_{dt}) = \frac{1}{n'} \sum_{i \in U} q_i \left(\frac{y_i}{q_i} - Y\right)^2 + \left(\frac{1}{n} - \frac{1}{n'}\right) \sum_{i \in U} q_i \left(\frac{d_{\lambda i}}{q_i} - D_{\lambda}\right)^2 \quad (10.2.25)$$

Unbiased estimators of $V\{E(\hat{Y}_{dt}|s')\}$ and $E\{V(\hat{Y}_{dt}|s')\}$ are given, respectively, by

$$\hat{V}\{E(\hat{Y}_{dt}|s')\} = \frac{1}{n'(n-1)} \sum_{i \in s} \left(\frac{y_i}{q_i} - \frac{1}{n} \sum_{i \in s} \frac{y_i}{q_i}\right)^2$$

and

$$\hat{V}(\hat{Y}_{dt}|s') = \left(\frac{1}{n} - \frac{1}{n'}\right) \frac{1}{n-1} \sum_{i \in s} \left(\frac{d_{\lambda i}}{q_i} - \frac{1}{n} \sum_{i \in s} \frac{d_{\lambda i}}{q_i}\right)^2$$

Hence an unbiased estimator for $V(\hat{Y}_{dt})$ is

$$\begin{aligned}\hat{V}(\hat{Y}_{dt}) &= \frac{1}{n'(n-1)} \sum_{i \in s} \left(\frac{y_i}{q_i} - \frac{1}{n} \sum_{i \in s} \frac{y_i}{q_i} \right)^2 \\ &\quad + \left(\frac{1}{n} - \frac{1}{n'} \right) \frac{1}{n-1} \sum_{i \in s} \left(\frac{d_{\lambda i}}{q_i} - \frac{1}{n} \sum_{i \in s} \frac{d_{\lambda i}}{q_i} \right)^2\end{aligned}\quad (10.2.26)$$

Under the cost function, (Eq. 10.2.19) two-phase sampling will be more precise than the single phase using the same cost if

$$\begin{aligned}\frac{1}{n'} \sum_{i \in U} q_i \left(\frac{y_i}{q_i} - Y \right)^2 + \left(\frac{1}{n} - \frac{1}{n'} \right) \sum_{i \in U} q_i \left(\frac{d_{\lambda i}}{q_i} - D_{\lambda} \right)^2 &< \frac{1}{n_0} \sum_{i \in U} q_i \left(\frac{y_i}{q_i} - Y \right)^2 \\ \text{i.e., if } \delta &> \frac{1}{2} \left[h^* + \left\{ h^* \left(1 - \frac{n}{n'} \right) \left(1 + \frac{nc}{n'c'} \right) \right\}^{-1} \right]\end{aligned}\quad (10.2.27)$$

where

$$\delta = \frac{\sum_{i \in U} q_i \left(\frac{y_i}{q_i} - Y \right) \left(\frac{x_i}{q_i} - X \right)}{\sqrt{\sum_{i \in U} q_i \left(\frac{y_i}{q_i} - Y \right)^2} \sqrt{\sum_{i \in U} q_i \left(\frac{x_i}{q_i} - X \right)^2}}$$

and

$$h^* = \lambda \sqrt{\sum_{i \in U} q_i \left(\frac{x_i}{q_i} - X \right)^2} / \sqrt{\sum_{i \in U} q_i \left(\frac{y_i}{q_i} - Y \right)^2}$$

10.2.2 Ratio Method of Estimation

The ratio estimator for the population total in two-phase sampling is given by

$$\hat{Y}_{Rt} = \frac{t(s, y)}{t(s, x)} t(s', x) \quad (10.2.28)$$

Like the ordinary ratio estimator \hat{Y}_R described in Eq. (8.3.1), \hat{Y}_{Rt} is a biased estimator for Y .

Now writing,

$$\begin{aligned}
 \text{Cov}\left\{\frac{t(s, y)}{t(s, x)}, t(s, x)|s'\right\} &= E(t(s, y)|s') - \left[E\left\{\frac{t(s, y)}{t(s, x)}|s'\right\}\right] E(t(s, x)|s') \\
 &= E(t(s, y)|s') - \left[E\left\{\frac{t(s, y)}{t(s, x)}|s'\right\}\right] t(s', x)
 \end{aligned} \tag{10.2.29}$$

and taking expectation both sides of Eq. (10.2.29) we get

$$\begin{aligned}
 E\left[\text{Cov}\left\{\frac{t(s, y)}{t(s, x)}, t(s, x)|s'\right\}\right] &= Y - E\left[E\left\{\frac{t(s, y)}{t(s, x)}t(s', x)|s'\right\}\right] \\
 &= Y - E(\widehat{Y}_{Rt})
 \end{aligned} \tag{10.2.30}$$

Thus we have

Theorem 10.2.1

The bias of \widehat{Y}_{Rt} under two-phase sampling is

$$B(\widehat{Y}_{Rt}) = -E\left[\text{Cov}\left\{\frac{t(s, y)}{t(s, x)}, t(s, x)|s'\right\}\right]$$

The above expression of bias is exact. However, it is not useful as it does not provide any elegant expression for any specific sampling design.

10.2.2.1 Approximate Expression of Bias

Let us define

$$e_x = \{t(s, x) - X\}/X, e_y = \{t(s, y) - Y\}/Y, \text{ and } e'_x = \{t(s', x) - X\}/X \tag{10.2.31}$$

Then,

$$\begin{aligned}
 \widehat{Y}_{Rt} &= Y \frac{(1 + e_y)}{(1 + e_x)} (1 + e'_x) \\
 &= Y(1 + e_y)(1 - e_x + e_x^2 - \cdots)(1 + e'_x) \\
 &= Y(1 + e_y - e_x + e'_x + e_x^2 - e_x e_y + e'_x e_y - e_x e'_x + \cdots)
 \end{aligned} \tag{10.2.32}$$

The expression of bias of \hat{Y}_{Rt} is given by

$$\begin{aligned} B(\hat{Y}_{Rt}) &= E(\hat{Y}_{Rt} - Y) \\ &= YE(e_y - e_x + e'_x + e_x^2 - e_x e_y + e'_x e_y - e_x e'_x + \cdots) \end{aligned} \quad (10.2.33)$$

On neglecting the terms $E(e_x^i e_y^j e_x^{lk})$ with $i + j + k > 2$, assuming it small for large sample size n , an approximate expression of bias of \hat{Y}_{Rt} up to the first order of approximation comes out as

$$\begin{aligned} B(\hat{Y}_{Rt}) &= E(\hat{Y}_{Rt} - Y) \\ &= YE(e_y - e_x + e'_x + e_x^2 - e_x e_y + e'_x e_y - e_x e'_x) \\ &= Y \left[E(e_x^2 - e_x e_y) + E\{e'_x E(e_y|s')\} - e'_x E(e_x|s') \right] \\ &= Y \left[E(e_x^2 - e_x e_y) + E(e'_x e'_y - e_x'^2) \right] \\ &= Y \left[\frac{V\{t(s, x)\}}{X^2} - \frac{\text{Cov}\{t(s, x), t(s, y)\}}{XY} + \frac{\text{Cov}\{t(s', y), t(s', x)\}}{XY} - \frac{V\{t(s', x)\}}{X'^2} \right] \\ &= Y \left[\frac{V\{t(s, x)\} - V\{t(s', x)\}}{X^2} - \frac{\text{Cov}\{t(s, x), t(s, y)\} - \text{Cov}\{t(s', y), t(s', x)\}}{XY} \right] \end{aligned} \quad (10.2.34)$$

Now substituting

$$\begin{aligned} V\{t(s, x)\} &= E[V\{t(s, x)|s'\}] + V[E\{t(s, x)|s'\}] \\ &= E[V\{t(s, x)|s'\}] + V\{t(s', x)\} \end{aligned}$$

and

$$\begin{aligned} \text{Cov}\{t(s, x), t(s, y)\} &= E[\text{Cov}\{t(s, x), t(s, y)|s'\}] + \text{Cov}[E\{t(s, x)|s'\}, E\{t(s, y)|s'\}] \\ &= E[\text{Cov}\{t(s, x), t(s, y)|s'\}] + \text{Cov}\{t(s', x), t(s', y)\} \end{aligned}$$

in Eq. (10.2.34), the expression of $B(\hat{Y}_{Rt})$ reduces to

$$B(\hat{Y}_{Rt}) = Y \left[\frac{E[V\{t(s, x)|s'\}]}{X^2} - \frac{E[\text{Cov}\{t(s, x), t(s, y)|s'\}]}{XY} \right] \quad (10.2.35)$$

10.2.2.2 Approximate Expression of Mean Square Error

The mean square error of \hat{Y}_{Rt} is given by

$$\begin{aligned} M(\hat{Y}_{Rt}) &= E(\hat{Y}_{Rt} - Y)^2 \\ &= Y^2 E(e_y - e_x + e'_x + e_x^2 - e_x e_y + e'_x e_y - e_x e'_x + \dots)^2 \end{aligned}$$

Neglecting the terms $E(e_x^i e_y^j e_x^k)$ with $i + j + k > 2$, an approximate expression of $M(\hat{Y}_{Rt})$ up to the first order of approximation is obtained as

$$\begin{aligned} M(\hat{Y}_{Rt}) &= Y^2 E(e_y^2 + e_x^2 + e_x'^2 - 2e_y e_x - 2e_x e_x'^2 + 2e_y e_x'^2) \\ &= Y^2 E\left[(e_y^2 + e_x^2 - 2e_y e_x) + \{e_x'^2 - 2E(e_x|s')e_x'^2 + 2E(e_y|s')e_x'^2\}\right] \\ &= Y^2 E\left[(e_y^2 + e_x^2 - 2e_y e_x) - (e_x'^2 - 2e_y' e_x')\right] \\ &= \left[V\{t(s, y)\} + R^2 V\{t(s, x)\} - 2RCov\{t(s, y), t(s, x)\}\right] \\ &\quad - \left[R^2 V\{t(s', x)\} - 2RCov\{t(s', y), t(s', x)\}\right] \\ &= V\left[t(s, y) - R t(s, x)\right] - V\left[E\{t(s, y) - R t(s, x)|s'\} + V\{t(s', y)\}\right] \\ &= V\{t(s', y)\} + E\left[V\{t(s, y) - R t(s, x)|s'\}\right] \\ &= V\{t(s', y)\} + E\left[V\{t(s, d)|s'\}\right] \end{aligned} \tag{10.2.36}$$

where $t(s, d) = \sum_{i \in s} b_{si} d_i$ and $d_i = y_i - R x_i$.

An approximate expression for an unbiased estimator of $M(\hat{Y}_{Rt})$ is given by

$$\hat{M}(\hat{Y}_{Rt}) = \hat{V}\{t(s', y)\} + \hat{V}\left\{t\left(s, \hat{d}\right)|s'\right\} \tag{10.2.37}$$

where $\hat{V}\{t(s', y)\}$ is an unbiased estimator of $V\{t(s', y)\}$ and $\hat{V}\left\{t\left(s, \hat{d}\right)|s'\right\}$ is an approximate unbiased estimator of $V\{t(s, d)|s'\}$, which is obtained by replacing d_i with its estimate $\hat{d}_i = y_i - \hat{R} x_i$ in the expression of $\hat{V}\{t(s, d)|s'\}$, an unbiased estimator of $V\{t(s, d)|s'\}$.

Eqs. (10.2.35)–(10.2.37) yield

Theorem 10.2.2

Approximate expressions for

(i) Bias of \hat{Y}_{Rt} is

$$B(\hat{Y}_{Rt}) = Y \left[\frac{E[V\{t(s, x)|s'\}]}{X^2} - \frac{E[Cov\{t(s, x), t(s, y)|s'\}]}{XY} \right]$$

(ii) Mean square error of \hat{Y}_{Rt} is

$$M(\hat{Y}_{Rt}) = V\{t(s', y)\} + E[V\{t(s, d)|s'\}]$$

(iii) Estimator of $M(\hat{Y}_{Rt})$ is

$$\hat{M}(\hat{Y}_{Rt}) = \hat{V}\{t(s', y)\} + \hat{V}\{t(s, \hat{d})|s'\}$$

10.2.2.3 Simple Random Sampling Without Replacement

Consider the situation when both the samples s and s' are selected by SRSWOR method of sampling. In this case $t(s, x) = N \bar{x}(s)$, $t(s, y) = N \bar{y}(s)$, $t(s', x) = N \bar{x}(s')$, $t(s', y) = N \bar{y}(s')$, $t(s, d) = N\{\bar{y}(s) - R \bar{x}(s)\}$, $E[V\{N \bar{x}(s)|s'\}] = N^2 \left(\frac{1}{n} - \frac{1}{n'} \right) S_x^2$, $V\{N \bar{y}(s')\} = N^2 \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2$, $E[Cov\{N \bar{x}(s), N \bar{y}(s)|s'\}] = N^2 \left(\frac{1}{n} - \frac{1}{n'} \right) S_{xy}$ with $S_{xy} = \sum_{i \in U} (x_i - \bar{X})(y_i - \bar{Y}) / (N - 1)$, and $E[V(N\{\bar{y}(s) - R \bar{x}(s)\}|s')] = N^2 \left(\frac{1}{n} - \frac{1}{n'} \right) (S_y^2 + R^2 S_x^2 - 2RS_{xy})$.

The expressions of \hat{Y}_{Rt} , $B(\hat{Y}_{Rt})$, $\hat{B}(\hat{Y}_{Rt})$, $M(\hat{Y}_{Rt})$, and $\hat{M}(\hat{Y}_{Rt})$ appear as follows:

$$\hat{Y}_{Rt} = N \frac{\bar{y}(s)}{\bar{x}(s)} \bar{x}(s')$$

$$\begin{aligned} B(\hat{Y}_{Rt}) &= N \left(\frac{1}{n} - \frac{1}{n'} \right) \bar{Y} \left(\frac{S_x^2}{\bar{X}^2} - \frac{S_{xy}}{\bar{X} \bar{Y}} \right) \\ &= N \left(\frac{1}{n} - \frac{1}{n'} \right) \bar{Y} (C_x^2 - \rho C_x C_y) \end{aligned}$$

$$\hat{B}(\hat{Y}_{Rt}) = N \left(\frac{1}{n} - \frac{1}{n'} \right) \bar{y}(s) \left(\frac{s_x^2}{\{\bar{x}(s)\}^2} - \frac{s_{xy}}{\bar{x}(s) \bar{y}(s)} \right)$$

$$M(\hat{Y}_{Rt}) = N^2 \left[\left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n'} \right) (S_y^2 + R^2 S_x^2 - 2RS_{xy}) \right]$$

$$\hat{M}(\hat{Y}_{Rt}) = N^2 \left[\left(\frac{1}{n'} - \frac{1}{N} \right) s_y^2 + \left(\frac{1}{n} - \frac{1}{n'} \right) (s_y^2 + \hat{R}^2 s_x^2 - 2\hat{R}s_{xy}) \right]$$

$$\text{where } s_x^2 = \frac{\sum_{i \in s} \{x_i - \bar{x}(s)\}^2}{(n-1)}, \quad s_y^2 = \frac{\sum_{i \in s} \{y_i - \bar{y}(s)\}^2}{n-1},$$

$$s_{xy} = \frac{\sum_{i \in s} \{x_i - \bar{x}(s)\} \{y_i - \bar{y}(s)\}}{n-1}; \quad C_x \text{ and } C_y \text{ are the coefficient of variation of } x \text{ and } y, \text{ respectively.}$$

If the total cost of the survey is fixed, the two-phase sampling will be more precise than the single phase sampling under the cost function (Eq. 10.2.19) if

$$N^2 \left[\left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n'} \right) (S_y^2 + R^2 S_x^2 - 2RS_{xy}) \right] < N^2 \left(\frac{1}{n_0} - \frac{1}{N} \right) S_y^2$$

$$\text{i.e., if } \rho > \frac{1}{2} \left[\bar{h} + \left\{ \bar{h} \left(1 - \frac{n}{n'} \right) \left(1 + \frac{nc}{n'c'} \right) \right\}^{-1} \right] \quad (10.2.38)$$

where $\bar{h} = C_x/C_y$.

10.2.2.4 Optimal Allocation Under Simple Random Sampling Without Replacement

Here we will find the optimum values of n and n' , which minimize $M(\hat{Y}_{Rt})$ when total cost of the survey $C = c'n' + cn$ is fixed to a certain level C_0 . For this minimization problem, consider

$$\phi = M(\hat{Y}_{Rt}) - \mu(C_0 - cn - c'n')$$

where μ is an undetermined Lagrange multiplier.

For large N , we can write

$$\phi = N^2 \left[\frac{1}{n'} S_y^2 + \left(\frac{1}{n} - \frac{1}{n'} \right) S_d^2 \right] - \mu(C_0 - cn - c'n')$$

where $S_d^2 = S_y^2 + R^2 S_x^2 - 2RS_{xy}$.

Now $\frac{\partial \phi}{\partial n} = 0$ and $\frac{\partial \phi}{\partial n'} = 0$ yield

$$n = NS_d / \sqrt{c\mu} \quad \text{and} \quad n' = N \sqrt{S_y^2 - S_d^2} / \sqrt{c'\mu} \quad (10.2.39)$$

Substituting n and n' from Eq. (10.2.39) in the constraint $C_0 = c'n' + cn$, the optimum values of n and n' are derived as

$$\begin{aligned} n_{opt} &= \frac{C_0}{S_d\sqrt{c} + \sqrt{(S_y^2 - S_d^2)c'}} \frac{S_d}{\sqrt{c}} \quad \text{and} \\ n'_{opt} &= \frac{C_0}{S_d\sqrt{c} + \sqrt{(S_y^2 - S_d^2)c'}} \frac{\sqrt{S_y^2 - S_d^2}}{\sqrt{c'}} \end{aligned} \quad (10.2.40)$$

Placing the optimum values of n and n' in the expression of $M(\hat{Y}_{Rt})$ and taking $1/N \cong 0$, we arrive at the minimum value of $M(\hat{Y}_{Rt})$ for large N as

$$M_{\min}(\hat{Y}_{Rt}) = \frac{N^2 \left[S_d\sqrt{c} + \sqrt{(S_y^2 - S_d^2)c'} \right]^2}{C_0} \quad (10.2.41)$$

10.2.3 Regression Method of Estimation

The regression estimator for the population total for two-phase sampling is

$$\hat{Y}_{regt} = t(s, y) - B(s)[t(s, x) - t(s', x)] \quad (10.2.42)$$

where $B(s) = \frac{C\hat{o}v\{t(s, y), t(s, x)|s'\}}{\hat{V}\{t(s, x)|s'\}}$, $C\hat{o}v\{t(s, y), t(s, x)|s'\}$ is an unbiased estimator of $Cov\{t(s, y), t(s, x) - t(s', x)\} = E[Cov\{t(s, y), t(s, x)|s'\}]$, and $\hat{V}\{t(s, x)|s'\}$ is an unbiased estimator of $V\{t(s, x) - t(s', x)\} = E[V\{t(s, x)|s'\}]$.

10.2.3.1 Approximate Expressions of Bias and Mean Square Errors

Let $e_B = \frac{B(s) - B}{B}$ with $B = \frac{Cov\{t(s, y), t(s, x) - t(s', x)\}}{V\{t(s, x) - t(s', x)\}}$.

The bias of \hat{Y}_{regt} is

$$\begin{aligned} B(\hat{Y}_{regt}) &= E(\hat{Y}_{regt}) - Y \\ &= E[Y(1 + e_B) - B(1 + e_B)X(e_x - e'_x)] - Y \\ &= -B X E\{(e_B)(e_x - e'_x)\} \\ &= -E\{B(s), t(s, x) - t(s', x)\} \\ &= -E[Cov\{B(s), t(s, x)|s'\}] \end{aligned} \quad (10.2.43)$$

The mean square error of \hat{Y}_{regt} is given by

$$\begin{aligned} M(\hat{Y}_{regt}) &= E(\hat{Y}_{regt} - Y)^2 \\ &= E[Y(1 + e_y) - BX(1 + e_B)(e_x - e'_x) - Y]^2 \end{aligned}$$

Now neglecting terms $E\{e_B^i e_x^j e'_x{}^k e_y^l\}$ for $i + j + k + l > 2$, an approximate expression of $M(\hat{Y}_{regt})$ is obtained as

$$\begin{aligned} M(\hat{Y}_{regt}) &\cong V[t(s, y) - B\{t(s, x) - t(s', x)\}] \\ &= V\{t(s', y)\} + E[V\{t(s, \epsilon)|s'\}] \end{aligned} \quad (10.2.44)$$

where $t(s, \epsilon) = \sum_{i \in s} b_{si} E_i$ and $E_i = y_i - B x_i$.

An approximate unbiased estimator of $M(\hat{Y}_{regt})$ is

$$\hat{M}(\hat{Y}_{regt}) = \hat{V}(t(s', y)) + \hat{V}(t(s, \hat{\epsilon})|s') \quad (10.2.45)$$

where $\hat{V}(t(s', y))$ is an unbiased estimator of $V(t(s, y))$ and $\hat{V}(t(s, \hat{\epsilon})|s')$ is an approximate unbiased estimator of $V(t(s, \epsilon)|s')$. The estimator $\hat{V}(t(s, \hat{\epsilon})|s')$ is obtained by replacing E_i by its estimate $\hat{E}_i = y_i - B(s)x_i$ in the expression of $\hat{V}(t(s, \epsilon)|s')$, an unbiased estimator of $V(t(s, \epsilon)|s')$.

10.2.3.2 Arbitrary Sampling Design

Let the initial sample s' and s be selected as in [Section 10.2.1.1](#) and

$$t(s', y) = \sum_{i \in s'} \frac{y_i}{\pi_i}, t(s', x) = \sum_{i \in s'} \frac{x_i}{\pi_i} \text{ and } t(s, y) = \frac{n'}{n} \sum_{i \in s} \frac{y_i}{\pi_i}.$$

In this case

$$\hat{Y}_{regt} = \frac{n'}{n} \sum_{i \in s} \frac{y_i}{\pi_i} - B(s) \left(\frac{n'}{n} \sum_{i \in s} \frac{x_i}{\pi_i} - \sum_{i \in s'} \frac{x_i}{\pi_i} \right) \quad (10.2.46)$$

where

$$B(s) = \frac{C\hat{o}v\left(\sum_{i \in s} \frac{y_i}{\pi_i}, \sum_{i \in s} \frac{x_i}{\pi_i} | s'\right)}{\hat{V}\left(\sum_{i \in s} \frac{y_i}{\pi_i} | s'\right)} = \frac{\sum_{i \in s} \left(\frac{x_i}{\pi_i} - \frac{1}{n} \sum_{i \in s} \frac{x_i}{\pi_i}\right) \sum_{i \in s} \left(\frac{y_i}{\pi_i} - \frac{1}{n} \sum_{i \in s} \frac{y_i}{\pi_i}\right)}{\sum_{i \in s} \left(\frac{x_i}{\pi_i} - \frac{1}{n} \sum_{i \in s} \frac{x_i}{\pi_i}\right)^2}.$$

The expression of bias of \hat{Y}_{regt} is very complicated and hence it is omitted here. On the other hand, the expression of $M(\hat{Y}_{regt})$ is quite elegant and is given as follows:

$$\begin{aligned}
 M(\hat{Y}_{regt}) &\cong V\left(\sum_{i \in s'} \frac{y_i}{\pi_i}\right) + E\left\{V\left(\frac{n'}{n} \sum_{i \in s} \frac{E_i}{\pi_i} \middle| s'\right)\right\} \\
 &= \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j}\right)^2 \\
 &\quad + \frac{n'(n' - n)}{n} \frac{1}{n' - 1} E\left[\sum_{i \in s'} \left\{\frac{E_i}{\pi_i} - \frac{1}{n'} \sum_{i \in s'} \frac{E_i}{\pi_i}\right\}^2\right] \\
 &= \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j}\right)^2 + \frac{n'(n' - n)}{n} \frac{1}{n' - 1} \\
 &\quad \times \left[E\left(\sum_{i \in s'} \frac{E_i^2}{\pi_i^2}\right) - \left\{V\left(\sum_{i \in s'} \frac{E_i}{\pi_i}\right)^2 + \left(\sum_{i \in U} E_i\right)^2\right\} / n'\right] \\
 &= \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j}\right)^2 + n' \frac{(n' - n)}{n(n' - 1)} \\
 &\quad \times \left[\left\{\sum_{i \in U} \frac{E_i^2}{\pi_i} - \left(\sum_{i \in U} E_i\right)^2 / n\right\} - \frac{1}{n'} \sum_{i \neq j} \sum_{j \in U} (\pi_i \pi_j - \pi_{ij}) \left(\frac{E_i}{\pi_i} - \frac{E_j}{\pi_j}\right)^2\right]
 \end{aligned} \tag{10.2.47}$$

$$\text{where } E_i = y_i - Bx_i \text{ and } B = \frac{E\left\{Cov\left(\sum_{i \in s} \frac{y_i}{\pi_i}, \sum_{i \in s} \frac{x_i}{\pi_i}\right) \middle| s'\right\}}{E\left\{V\left(\sum_{i \in s} \frac{x_i}{\pi_i}\right) \middle| s'\right\}}.$$

An approximate unbiased estimator of $M(\hat{Y}_{regt})$ is given by

$$\begin{aligned}
 \hat{M}(\hat{Y}_{regt}) &\cong \frac{n'(n' - 1)}{n(n - 1)} \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \left(\frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}}\right) \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j}\right)^2 \\
 &\quad + n' \frac{(n' - n)}{n(n - 1)} \left\{\sum_{i \in s} \frac{\hat{E}_i^2}{\pi_i^2} - \frac{1}{n} \left(\sum_{i \in s} \frac{\hat{E}_i}{\pi_i}\right)^2\right\}
 \end{aligned} \tag{10.2.48}$$

with $\hat{E}_i = y_i - B(s) x_i$.

10.2.3.3 Simple Random Sampling Without Replacement

In case samples s' and s are selected by SRSWOR method we have

$$\begin{aligned}
 t(s', y) &= N \frac{1}{n'} \sum_{i \in s'} y_i = N \bar{y}(s'), t(s', x) = N \frac{1}{n'} \sum_{i \in s'} x_i = N \bar{x}(s'), \\
 t(s, y) &= n' \frac{1}{n} \sum_{i \in s} y_i = n' \bar{y}(s), \\
 B(s) &= \frac{\widehat{Cov}\{\bar{y}(s), \bar{x}(s)|s'\}}{\widehat{V}\{\bar{x}(s)|s'\}} = \frac{s_{xy}}{s_x^2} \text{ and} \\
 \hat{Y}_{regt} &= N \left[\bar{y}(s) - \frac{s_{xy}}{s_x^2} (\bar{x}(s) - \bar{x}(s')) \right] \quad (10.2.49)
 \end{aligned}$$

An approximate bias of \hat{Y}_{regt} is obtained from Eq. (10.2.43) as

$$\begin{aligned}
 B(\hat{Y}_{regt}) &= -NE \left[Cov \left\{ \frac{s_{xy}}{s_x^2}, (\bar{x}(s) - \bar{x}(s')) | s \right\} \right] \\
 &= -NE \left[\frac{s_{xy}}{s_x^2} (\bar{x}(s) - \bar{x}(s')) \right]
 \end{aligned}$$

Now writing $\varepsilon_{xy} = \frac{s_{xy} - S_{xy}}{S_{xy}}$, $\varepsilon_{xx} = \frac{s_x^2 - S_x^2}{S_x^2}$, and $B = \beta = \frac{S_{xy}}{S_x^2}$, we get

$$\begin{aligned}
 B(\hat{Y}_{regt}) &= -N\beta E \left[(1 + \varepsilon_{xy})(1 + \varepsilon_{xx})^{-1} (\bar{x}(s) - \bar{x}(s')) \right] \\
 &\cong -N\beta E \left[(\varepsilon_{xy} - \varepsilon_{xx}) (\bar{x}(s) - \bar{x}(s')) \right] \quad (10.2.50) \\
 &= -N\beta \left[\frac{E[Cov\{s_{xy}, \bar{x}(s)|s'\}]}{S_{xy}} - \frac{E[\{s_x^2, \bar{x}(s)|s'\}]}{S_x^2} \right]
 \end{aligned}$$

Now using Sukhatme et al. (1984), we find

$$\begin{aligned}
 E[Cov\{s_{xy}, \bar{x}(s)|s'\}] &= \frac{n'(n' - n)}{(n' - 1)(n' - 2)n} E \left[\frac{1}{n'} \sum_{i \in s'} \{x_i - \bar{x}(s')\}^2 \{y_i - \bar{y}(s')\} \right] \\
 &= \frac{n'(n' - n)}{(n' - 1)(n' - 2)n} \mu_{21} \quad (10.2.51)
 \end{aligned}$$

and

$$\begin{aligned} E[\{s_x^2, \bar{x}(s)\} | s'] &= \frac{n'(n' - n)}{(n' - 1)(n' - 2)n} E \left[\frac{1}{n'} \sum_{i \in s'} \{x_i - \bar{x}(s')\}^3 \right] \\ &= \frac{n'(n' - n)}{(n' - 1)(n' - 2)n} \mu_{30} \end{aligned} \quad (10.2.52)$$

where $\mu_{ij} = \sum_{i \in U} (x_i - \bar{X})^i (y_i - \bar{Y})^j / N$

Now substituting Eqs. (10.2.51) and (10.2.52) in Eq. (10.2.50), we find

$$B(\hat{Y}_{regt}) \cong -\beta \frac{Nn'(n' - n)}{(n' - 1)(n' - 2)n} \left(\frac{\mu_{21}}{S_{xy}} - \frac{\mu_{30}}{S_x^2} \right)$$

If n' is large compared to n , the bias $B(\hat{Y}_{regt})$ can be approximated as

$$B(\hat{Y}_{regt}) \cong -N\beta \left(\frac{1}{n} - \frac{1}{n'} \right) \left(\frac{\mu_{21}}{S_{xy}} - \frac{\mu_{30}}{S_x^2} \right) \quad (10.2.53)$$

Following Eqs. (10.2.44) and (10.2.45), approximate expressions for the mean square error of \hat{Y}_{regt} and its unbiased estimator are obtained as

$$\begin{aligned} M(\hat{Y}_{regt}) &\cong N^2 \left[\left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 + \left(\frac{1}{n} - \frac{1}{n'} \right) S_E^2 \right] \\ &\cong N^2 \left[\left(\frac{1}{n'} - \frac{1}{N} \right) + \left(\frac{1}{n} - \frac{1}{n'} \right) (1 - \rho^2) \right] S_y^2 \end{aligned} \quad (10.2.54)$$

and

$$\begin{aligned} \hat{M}(\hat{Y}_{regt}) &\cong N^2 \left[\left(\frac{1}{n'} - \frac{1}{N} \right) s_y^2 + \left(\frac{1}{n} - \frac{1}{n'} \right) s_{\hat{E}}^2 \right] \\ &= N^2 \left[\left(\frac{1}{n'} - \frac{1}{N} \right) + \left(\frac{1}{n} - \frac{1}{n'} \right) (1 - r^2) \right] s_y^2 \end{aligned} \quad (10.2.55)$$

where $S_E^2 = \sum_{i \in U} (E_i - \bar{E})^2 / (N - 1)$, $s_{\hat{E}}^2 = \sum_{i \in s} (\hat{E}_i - \bar{\hat{E}})^2 / (n - 1)$, $\bar{E} = \sum_{i \in U} E_i / N$, $\bar{\hat{E}} = \sum_{i \in s} \hat{E}_i / n$, $\hat{E}_i = y_i - \hat{\beta}x_i$, and $r = s_{xy} / (s_x s_y) =$ sample correlation coefficient.

10.2.3.4 Optimum Allocation

For large N , the expression (Eq. 10.2.54) reduces to

$$M(\hat{Y}_{regt}) = N^2 \left[\frac{(1 - \rho^2)}{n} + \frac{\rho^2}{n'} \right] S_y^2 \quad (10.2.56)$$

The optimum values of n and n' that minimize (Eq. 10.2.56) for a given cost $C = C_0 = cn + c'n'$ are obtained by minimizing $N^2 \left[\frac{(1 - \rho^2)}{n} + \frac{\rho^2}{n'} \right] S_y^2$, subject to $C_0 = cn + c'n'$. For minimization consider the function

$$\phi = N^2 \left[\frac{(1 - \rho^2)}{n} + \frac{\rho^2}{n'} \right] S_y^2 - \mu(C_0 - cn - c'n')$$

where μ is an undetermined Lagrange multiplier.

$\frac{\partial \phi}{\partial n} = 0$ and $\frac{\partial \phi}{\partial n'} = 0$ yield the optimum values of n and n' are as follows:

$$n_{opt} = \frac{C_0}{(\sqrt{(1 - \rho^2)} \sqrt{c} + \rho \sqrt{c'})} \frac{\sqrt{1 - \rho^2}}{\sqrt{c}}$$

and

$$n'_{opt} = \frac{C_0}{(\sqrt{(1 - \rho^2)} \sqrt{c} + \rho \sqrt{c'})} \frac{\rho}{\sqrt{c'}}$$

The value of $M(\hat{Y}_{regt})$ with the optimum values of n and n' is given by

$$M_{min} = N^2 \frac{(\sqrt{(1 - \rho^2)} \sqrt{c} + \rho \sqrt{c'})^2}{C_0} S_y^2 \quad (10.2.57)$$

For the fixed cost C_0 , the value of $n_0 = C_0/c$. Hence, two-phase sampling with optimum allocation will be more precise than the single-phase sampling with the same cost for a large N

$$\begin{aligned} & \text{if } M_{min} < N^2 \frac{1}{n_0} S_y^2 \\ & \text{i.e. if } \left(\sqrt{(1 - \rho^2)} \sqrt{c} + \rho \sqrt{c'} \right)^2 < c \\ & \text{i.e. if } 2\sqrt{1 - \rho^2} \sqrt{cc'} < \rho^2(c - c') \\ & \text{i.e. if } \rho^2 > \frac{4cc'}{(c + c')^2} \end{aligned} \quad (10.2.58)$$

Example 10.2.1

A sample of 50 plants was selected from a garden of 500 plants by SRSWOR method, and eye estimates of the heights of the selected plants were noted. From the selected sample of 50 plants, a subsample of 10 plants was selected by SRSWOR method and the heights of the selected plants were measured accurately. The data are presented as follows. Estimate the average height of the plants by using (i) ratio and (ii) regression method of estimation using the information that the sample mean height of eye estimation of the 50 plants was 10.5 m. Compare standard errors of the estimators used. (iii) Given that the cost of eye estimation and actual estimation of measuring height are \$10 and \$50, respectively, estimate the optimum proportions of actual measurement that should be made for the ratio and regression estimators.

Eye estimates and actual heights of plants.

Plants	Eye estimate of height (in m) (x)	Actual height (in m) (y)
1	15	13.85
2	10	9.75
3	8	10
4	8.5	10.5
5	12	14
6	6	7
7	6.5	8
8	15	14
9	10	12
10	8	7.5

Here $N = 500$, $n' = 50$, $n = 10$, $\bar{x}(s') = 10.5$, $\bar{x}(s) = 9.9$, $\bar{y}(s) = 10.66$, $s_x^2 = 10.267$, $s_y^2 = 7.337$, $s_{xy} = 7.907$, $\hat{R} = 10.66/9.9 = 1.077$, $s_d^2 = 2.213$, and sample correlation coefficient $= r = 0.911$.

(i) Ratio estimation:

Estimated mean height of the plants under two-phase sampling scheme is

$$\hat{\bar{Y}}_{Rt} = \frac{\bar{y}(s)}{\bar{x}(s)} \bar{x}(s') = (10.66/9.9) \times 10.5 = 11.306 \text{ m}$$

$$\begin{aligned}\text{Estimated MSE of } \widehat{\bar{Y}}_R &= \widehat{M}(\widehat{\bar{Y}}_{Rt}) = \left(\frac{1}{n'} - \frac{1}{N}\right)s_y^2 + \left(\frac{1}{n} - \frac{1}{n'}\right)s_d^2 \\ &= \left[\left(\frac{1}{50} - \frac{1}{500}\right) \times 7.336 + \left(\frac{1}{10} - \frac{1}{50}\right) \times 2.213\right] \\ &= 0.3091\end{aligned}$$

$$\text{Estimated standard error of } \widehat{\bar{Y}}_{Rt} = SE(\widehat{\bar{Y}}_{Rt}) = \sqrt{\widehat{M}(\widehat{\bar{Y}}_{Rt})} = 0.556 \text{ m.}$$

(ii) Regression Estimation:

Estimated mean height of plants under two-phase sampling is

$$\begin{aligned}\widehat{\bar{Y}}_{regt} &= \bar{y}(s) - \frac{s_{xy}}{s_x^2} \{\bar{x}(s) - \bar{x}(s')\} \\ &= 10.66 - (7.116/10.266) \times (9.9 - 10.5) = 11.122 \text{ m}\end{aligned}$$

$$\begin{aligned}\text{Estimated MSE of } \widehat{\bar{Y}}_{regt} &= \widehat{M}(\widehat{\bar{Y}}_{regt}) = \left[\left(\frac{1}{n'} - \frac{1}{N}\right) + \left(\frac{1}{n} - \frac{1}{n'}\right)(1 - r^2)\right]s_y^2 \\ &= \left[\left(\frac{1}{50} - \frac{1}{500}\right) + \left(\frac{1}{10} - \frac{1}{50}\right)(1 - 0.911^2)\right] \times 7.336 = 0.232\end{aligned}$$

$$\begin{aligned}\text{Estimated standard error of } \widehat{\bar{Y}}_{regt} &= SE(\widehat{\bar{Y}}_{regt}) = \sqrt{\widehat{M}(\widehat{\bar{Y}}_{regt})} \\ &= 0.482 \text{ m.}\end{aligned}$$

Efficiency of the regression estimator compared to the ratio estimator is

$$\frac{M(\widehat{\bar{Y}}_{Rt})}{M(\widehat{\bar{Y}}_{regt})} \times 100 = 133.307\%.$$

(iii) Optimum proportion of second-phase sample:

The optimum proportion for the ratio estimator is

$$n_{opt}/n'_{opt} = \sqrt{\frac{c'}{c} \frac{s_d^2}{s_y^2 - s_d^2}} = \sqrt{\frac{c'}{c} \frac{s_d^2}{s_y^2 - s_d^2}} = \sqrt{\frac{10}{50} \times \frac{2.213}{5.214}} = 0.291$$

The optimum proportion for the regression estimator is

$$n_{opt}/n'_{opt} = \sqrt{\frac{c'}{c} \frac{(1 - r^2)}{r^2}} = \sqrt{\frac{10}{50} \times \frac{(1 - 0.911^2)}{0.911^2}} = 0.202$$

10.3 TWO-PHASE SAMPLING FOR STRATIFICATION

It is well known that stratification is an important tool to increase precision of estimators. In stratification, the entire population is divided into a number of disjoint and exhaustive strata so that the each stratum becomes homogeneous with respect to the character under study. If no criterion for formation of strata is known in advance, we may use two-phase sampling for forming the strata. In the first phase, a relatively large sample s' of size n' is selected from the entire population by SRSWOR method, and information on the auxiliary variable (which can be used for stratification) is obtained. Here it is supposed that by noting the value of the auxiliary variable, one can classify the sample s' into a predetermined number of " H " strata. Let s'_h be the set of units of size n'_h falling in the h th stratum. Here we make the assumption that the sample size n' is so large that n'_h 's always exceeds 1 for every $h = 1, \dots, H$; in other words, $\text{Prob}(n'_h \geq 1) = 1$. From each of the selected samples s'_h 's, subsamples s_h 's of sizes $n_h = \gamma_h n'_h$ are selected independently by SRSWOR method. Here γ_h 's are predetermined fractions and n'_h 's are assumed to be integers. Hence, for proportional allocation with fixed n , we get $\gamma_h = \gamma = n/n'$ and $n_h = nn'_h/n'$. For example, in an agricultural survey, one can select a large sample s' to acquire information about the size of the farm x (auxiliary variable) only, and on the basis of the observed x values, the farms can be classified into small, medium, large, or very large strata.

10.3.1 Estimation of Mean and Variance

Let y_{hi} and x_{hi} be the value of the study (y) and auxiliary variable (x) for the i th unit of the h th stratum $i = 1, \dots, N_h$, $h = 1, \dots, H$. Here the strata sizes N_h 's and weights $W_h = N_h/N$'s are not known. Let $\bar{Y}_h = \sum_{i=1}^{N_h} y_{hi}/N_h$ be the population mean for the h th stratum and $\bar{Y} = \sum_{h=1}^H W_h \bar{Y}_h$ be the mean for the entire population.

Theorem 10.3.1

Let $w_h = n'_h/n'$, then

- (i) $E(w_h) = W_h$, (ii) $V(w_h) = g W_h(1 - W_h)/n'$, and
 - (iii) $\text{Cov}(w_h, w_{h'}) = -g W_h W_{h'}/n'$ for $h \neq h'$
- where $g = (N - n')/(N - 1)$.

Proof

For a fixed $\sum_{k=1}^H n'_k = n'$, $\mathbf{n}' = (n'_1, \dots, n'_k, \dots, n'_H)$ follows the generalized hypergeometric distribution

$$\begin{aligned} \text{Prob} \left(n'_1 = t_1, \dots, n'_k = t_k, \dots, n'_H = t_H \middle| \sum_{k=1}^H n'_k = n' \right) \\ = \frac{\binom{N_1}{t_1} \dots \binom{N_k}{t_k} \dots \binom{N_H}{t_H}}{\binom{N}{n'}}; t_k = 0, 1, \dots, \min(n', N_k) \end{aligned}$$

For this generalized hypergeometric distribution,

$$E(n'_k) = n' W_k, V(n'_k) = n' g W_k (1 - W_k) \text{ and } \text{Cov}(n'_k, n'_{k'}) = -n' g W_k W_{k'}$$

Theorem follows from the aforementioned results.

Finally using Rao (1973), we have the following theorem.

Theorem 10.3.2

(i) The estimator $\widehat{\bar{Y}}_{stt} = \sum_{h=1}^H w_h \bar{y}(s_h)$ is unbiased \bar{Y}

$$\text{where } \bar{y}(s_h) = \sum_{j \in s_h} \gamma_{hj} / n_h$$

(ii) The variance of $\widehat{\bar{Y}}_{stt}$ is

$$\begin{aligned} V(\widehat{\bar{Y}}_{stt}) &= \frac{1}{n'} \sum_{h=1}^H W_h \left\{ \left(\frac{1}{\gamma_h} - 1 \right) + g \left(1 - \frac{1}{N_h} \right) \right\} S_{hy}^2 \\ &\quad + \frac{g}{n'} \sum_{h=1}^H W_h (\bar{Y}_h - \bar{Y})^2 \\ &= \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 + \sum_{h=1}^H \frac{W_h}{n'} \left(\frac{1}{\gamma_h} - 1 \right) S_{hy}^2 \end{aligned}$$

where $(N - 1)S_y^2 = \sum_{h=1}^H \sum_{i=1}^{N_h} (\gamma_{hi} - \bar{Y})^2$ and $(N_h - 1)S_{hy}^2 =$

$$\sum_{i=1}^{N_h} (\gamma_{hi} - \bar{Y}_h)^2$$

(iii) A nonnegative unbiased estimator of $V(\widehat{\bar{Y}}_{stt})$ is

$$\begin{aligned}\widehat{V}(\widehat{\bar{Y}}_{stt}) &= \frac{1}{Nn'} \left[\frac{N-1}{n'-1} \sum_{h=1}^H n'_h \left(\frac{1}{\gamma_h} - 1 \right) s_{hy}^2 \right. \\ &\quad \left. + \frac{N-n'}{n'-1} \left(\sum_{h=1}^H \frac{1}{\gamma_h} \sum_{i \in s_h} y_{hi}^2 - n' \widehat{\bar{Y}}_{stt}^2 \right) \right] \\ &= \frac{N-1}{N} \sum_{h=1}^H \left(\frac{n'_h-1}{n'-1} - \frac{n_h-1}{N-1} \right) \frac{w_h s_{hy}^2}{n_h} + \frac{N-n'}{N(n'-1)} \sum_{h=1}^H w_h \left(\bar{y}(s_h) - \widehat{\bar{Y}}_{stt} \right)^2 \\ \text{where } g &= \frac{N-n'}{N-1} \text{ and } (n_h-1)s_{hy}^2 = \sum_{i \in s_h} \{y_{hi} - \bar{y}(s_h)\}^2.\end{aligned}$$

Proof

Let $\mathbf{n}' = (n'_1, \dots, n'_h, \dots, n'_H)$. Then,

$$\begin{aligned}\text{(i)} \quad E(\widehat{\bar{Y}}_{stt}) &= E \left[E \left(\sum_{h=1}^H w_h \bar{y}(s_h) | \mathbf{n}' \right) \right] \\ &= E \left[\sum_{h=1}^H w_h E(\bar{y}(s_h) | \mathbf{n}') \right] \\ &= E \left[\sum_{h=1}^H w_h \bar{Y}_h \right] \\ &= \sum_{h=1}^H \{E(w_h)\} \bar{Y}_h \\ &= \sum_{h=1}^H W_h \bar{Y}_h \quad (\text{using Theorem 10.3.1}) \\ &= \bar{Y}\end{aligned}$$

$$\begin{aligned}\text{(ii)} \quad V(\widehat{\bar{Y}}_{stt}) &= E \left[V \left(\sum_{h=1}^H w_h \bar{y}(s_h) | \mathbf{n}' \right) \right] + V \left[E \left(\sum_{h=1}^H w_h \bar{y}(s_h) | \mathbf{n}' \right) \right] \\ &\quad (10.3.1)\end{aligned}$$

Now,

$$\begin{aligned}
 E \left[V \left(\sum_{h=1}^H w_h \bar{y}(s_h) | \mathbf{n}' \right) \right] &= E \left[\sum_{h=1}^H w_h^2 V \{ \bar{y}(s_h) | \mathbf{n}' \} \right. \\
 &\quad \left. + \sum_{h \neq h'}^H \sum_{h'=1}^H w_h w_{h'} \text{Cov} \{ \bar{y}(s_h), \bar{y}(s_{h'}) | \mathbf{n}' \} \right] \\
 &= E \left[\sum_{h=1}^H w_h^2 \left(\frac{1}{n_h} - \frac{1}{N_h} \right) S_{hy}^2 \right] \quad (10.3.2)
 \end{aligned}$$

because $V \{ \bar{y}(s_h) | \mathbf{n}' \} = \left(\frac{1}{n_h} - \frac{1}{N_h} \right) S_{hy}^2$ and $\text{Cov} \{ \bar{y}(s_h), \bar{y}(s_{h'}) | \mathbf{n}' \} = 0$ as the samples s_h and $s_{h'}$ are selected independently for $h \neq h'$.

Writing $n_h = \gamma_h n'_h$ in Eq. (10.3.2) we get

$$\begin{aligned}
 E \left[V \left(\sum_{h=1}^H w_h \bar{y}(s_h) | \mathbf{n}' \right) \right] &= \sum_{h=1}^H \left(\frac{E(w_h)}{n'_h \gamma_h} - \frac{E(w_h^2)}{N_h} \right) S_{hy}^2 \\
 &= \sum_{h=1}^H W_h \left\{ \left(\frac{1}{n'_h \gamma_h} - \frac{1}{N} \right) - \frac{g}{n'_h} \left(\frac{1 - W_h}{N_h} \right) \right\} S_{hy}^2 \quad (10.3.3) \\
 &\quad (\text{using Theorem 10.3.1.}) \\
 &= \frac{1}{n'_h} \sum_{h=1}^H W_h \left\{ \left(\frac{1}{\gamma_h} - 1 \right) + g \left(1 - \frac{1}{N_h} \right) \right\} S_{hy}^2
 \end{aligned}$$

Furthermore,

$$\begin{aligned}
 V \left[E \left(\sum_{h=1}^H w_h \bar{y}(s_h) | \mathbf{n}' \right) \right] &= V \left(\sum_{h=1}^H w_h \bar{Y}_h \right) \\
 &= \frac{g}{n'_h} \left(\sum_{h=1}^H W_h (1 - W_h) \bar{Y}_h^2 - \sum_{h \neq h'}^H \sum_{h'=1}^H W_h W_{h'} \bar{Y}_h \bar{Y}_{h'} \right) \\
 &= \frac{g}{n'_h} \sum_{h=1}^H W_h (\bar{Y}_h - \bar{Y})^2 \quad (10.3.4)
 \end{aligned}$$

Substituting Eqs. (10.3.3) and (10.3.4) in Eq. (10.3.1) and noting

$(N - 1)S_y^2 = \sum_{h=1}^H (N_h - 1)S_{hy}^2 + \sum_{h=1}^H N_h (\bar{Y}_h - \bar{Y})^2$ we can verify (ii) of the theorem.

(iii) Now, noting $\frac{1}{N-1} \left[N \sum_{h=1}^H \frac{w_h}{n_h} \sum_{i \in s_h} y_{hi}^2 - N \widehat{\bar{Y}}_{stt}^2 \left\{ -\widehat{V}(\widehat{\bar{Y}}_{stt}) \right\} \right]$ and $\sum_{h=1}^H \frac{w_h}{n'} \left(\frac{1}{\gamma_h} - 1 \right) s_{hy}^2$ are unbiased estimators of S_y^2 and $\sum_{h=1}^H \frac{W_h}{n'} \left(\frac{1}{\gamma_h} - 1 \right) S_h^2$, respectively, we find an unbiased estimator of

$V(\widehat{\bar{Y}}_{stt})$ from the equation

$$\begin{aligned} \widehat{V}(\widehat{\bar{Y}}_{stt}) &= \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{1}{N-1} \left[N \sum_{h=1}^H \frac{w_h}{n_h} \sum_{i \in s_h} y_{hi}^2 - N \left\{ \widehat{\bar{Y}}_{stt}^2 - \widehat{V}(\widehat{\bar{Y}}_{stt}) \right\} \right] \\ &\quad + \sum_{h=1}^H \frac{w_h}{n'} \left(\frac{1}{\gamma_h} - 1 \right) s_{hy}^2 \end{aligned} \quad (10.3.5)$$

Eq. (10.3.5) yields

$$\begin{aligned} \text{i.e., } \widehat{V}(\widehat{\bar{Y}}_{stt}) &= \frac{1}{Nn'} \left[\frac{N-1}{n'-1} \sum_{h=1}^H n'_h \left(\frac{1}{\gamma_h} - 1 \right) s_{hy}^2 \right. \\ &\quad \left. + \frac{N-n'}{n'-1} \left(\sum_{h=1}^H \frac{1}{\gamma_h} \sum_{i \in s_h} y_{hi}^2 - n' \widehat{\bar{Y}}_{stt}^2 \right) \right] \end{aligned} \quad (10.3.6)$$

Furthermore, writing $\sum_{i \in s_h} y_{hi}^2 = (n_h - 1)s_{hy}^2 + n_h \{\bar{y}(s_h)\}^2$ in Eq. (10.3.6) we find

$$\begin{aligned} \widehat{V}(\widehat{\bar{Y}}_{stt}) &= \frac{N-1}{N} \sum_{h=1}^H \left(\frac{n'_h - 1}{n' - 1} - \frac{n_h - 1}{N - 1} \right) \frac{w_h s_{hy}^2}{n_h} \\ &\quad + \frac{N-n'}{N(n'-1)} \sum_{h=1}^H w_h \left(\bar{y}(s_h) - \widehat{\bar{Y}}_{stt} \right)^2 \end{aligned}$$

10.3.2 Proportional Allocation

In case the samples are selected from each of the strata by proportional allocation, keeping the total sample size for the second phase fixed as n , then γ_h becomes equal to $\gamma = n/n'$. Furthermore, if the stratum sizes are large for each of the strata, we derive the following result straight from the Theorem 10.3.2.

Theorem 10.3.3

For proportional allocation with large strata sizes

$$(i) \quad V(\bar{Y}_{stt}) \cong \sum_{h=1}^H W_h \left\{ \frac{(1-f_2)}{n} + \frac{(1-f_1)}{n'} \right\} S_{hy}^2 \\ + \frac{(1-f_1)}{n'} \sum_{h=1}^H W_h (\bar{Y}_h - \bar{Y})^2$$

and

(ii) An unbiased estimator of $V(\hat{\bar{Y}}_{stt})$ is

$$\hat{V}(\hat{Y}) = \frac{(1-f_1)}{n'} \sum_{h=1}^H w_h (\bar{y}(s_h) - \hat{\bar{Y}}_{stt})^2 + \frac{1}{n} \sum_{h=1}^H (1-f) w_h s_{hy}^2$$

where $f = n/N$, $f_1 = n'/N$ and $f_2 = n/n'$.

10.3.3 Estimation of Proportion

To estimate the proportion of persons who possess some attribute A , such as being a smoker or wealthy, we define $y_{hi} = 1$ if the i th unit of the h th stratum belongs to A and $y_{hi} = 0$ otherwise. Then \bar{Y} equals to

$\pi = \sum_{h=1}^H W_h \pi_h$, the proportion of persons belonging to the group A , and π_h

becomes the proportion in the h th stratum. On substituting $\bar{y}(s_h) = \sum_{j \in s_h} y_{hj} / n_h = \hat{\pi}_h$ sample proportion for the h th stratum,

$S_{hy}^2 = \frac{N_h \pi_h (1 - \pi_h)}{N_h - 1}$ and $s_{hy}^2 = \frac{n_h \hat{\pi}_h (1 - \hat{\pi}_h)}{n_h - 1}$ in [Theorem 10.3.2](#), we get the following.

Theorem 10.3.4

(i) An unbiased estimator of π is $\hat{\pi} = \sum_{h=1}^H w_h \hat{\pi}_h$

(ii) The variance of $\hat{\pi}$ is

$$V(\hat{\pi}) = \frac{1}{n'} \sum_{h=1}^H W_h \left\{ \left(\frac{1}{\gamma_h} - 1 \right) + g \left(1 - \frac{1}{N_h} \right) \right\} \frac{N_h \pi_h (1 - \pi_h)}{N_h - 1} \\ + \frac{g}{n'} \sum_{h=1}^H W_h (\pi_h - \pi)^2$$

(iii) An unbiased estimator of $V(\hat{\pi})$ is

$$\hat{V}(\hat{\pi}) = \frac{N-1}{N} \sum_{h=1}^H \left(\frac{n'_h - 1}{n' - 1} - \frac{n_h - 1}{N - 1} \right) \frac{w_h \hat{\pi}_h (1 - \hat{\pi}_h)}{n_h - 1} \\ + \frac{N - n'}{N(n' - 1)} \sum_{h=1}^H w_h (\hat{\pi}_h - \hat{\pi})^2$$

10.3.4 Optimum Allocation of Sample Sizes

Following Rao (1973), we consider the simple cost function

$$C = n'c' + \sum_{h=1}^H n_h c_h$$

where c' = cost of surveying a unit in the first phase and c_h = cost of surveying a unit of the h th stratum in the second phase. Clearly, c_h 's are expected to be much larger than c' 's. Because C is a random variable, we determine the optimum values of γ_h 's, which minimize the variance $V(\hat{\bar{Y}}_{stt})$, keeping the expected cost $E(C) = n'c' + n' \sum_{h=1}^H \gamma_h W_h c_h$ to a certain level C^* .

$$\text{i.e., } C^* = n' \left(c' + \sum_{h=1}^H \gamma_h W_h c_h \right) \quad (10.3.7)$$

Here we perform minimization in two stages. In the first stage, we minimize

$$V(\hat{\bar{Y}}_{stt}) = \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 + \sum_{h=1}^H \frac{W_h}{n'} \left(\frac{1}{\gamma_h} - 1 \right) S_{hy}^2 \quad (10.3.8)$$

subject to Eq. (10.3.8) keeping n' as fixed. For minimization consider

$$\phi = \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 + \sum_{h=1}^H \frac{W_h}{n'} \left(\frac{1}{\gamma_h} - 1 \right) S_{hy}^2 - \lambda \left(C^* - n'c' - n' \sum_{h=1}^H \gamma_h W_h c_h \right)$$

where λ is an undetermined Lagrange multiplier.

$$\frac{\partial \phi}{\partial \gamma_h} = 0 \text{ implies}$$

$$\gamma_h = \frac{1}{n' \sqrt{\lambda}} \frac{S_{hy}}{\sqrt{c_h}} \quad (10.3.9)$$

Eqs. (10.3.7) and (10.3.9) yield

$$\gamma_h = \frac{C^* - n'c'}{n' \left(\sum_{h=1}^H W_h S_{hy} \sqrt{c_h} \right)} \frac{S_{hy}}{\sqrt{c_h}} \quad (10.3.10)$$

Finally, putting Eq. (10.3.10) in Eq. (10.3.8) and minimizing it with respect to n' , the optimum value of n' comes out as

$$n' = \frac{C^*}{\sqrt{c'}} \left[\sqrt{c'} + \frac{\sum_{h=1}^H W_h S_{hy} \sqrt{c_h}}{\sqrt{S_y^2 - \sum_{h=1}^H W_h S_{hy}^2}} \right]^{-1} \quad (10.3.11)$$

Substituting Eq. (10.3.11) in Eq. (10.3.10) the optimum value of γ_h is obtained as

$$\gamma_h = \gamma_{h0} = S_{hy} \sqrt{\frac{c'}{c_h}} / \sqrt{S_y^2 - \sum_{h=1}^H W_h S_{hy}^2} \quad (10.3.12)$$

provided $\gamma_h \leq 1$ for all $h = 1, \dots, H$. It should be noted that $\gamma_h > 0$ because $\left(S_y^2 - \sum_{h=1}^H W_h S_{hy}^2 \right) > 0$ except in pathological situations. If some γ_h 's exceed 1, we should set them the value at 1 and determine the rest γ_h 's so that all of them are less than or equal to 1.

Now, substituting the optimum values of n' and γ_h from Eqs. (10.3.11) and (10.3.12) in the expression of $V(\hat{\bar{Y}}_{stt})$, we get the minimum value of $V(\hat{\bar{Y}}_{stt})$ as

$$V_{\min}(\hat{\bar{Y}}_{stt}) = \frac{1}{C^*} \left[\sum_{h=1}^H W_h S_{hy} \sqrt{c_h} + \sqrt{\left(S_y^2 - \sum_{h=1}^H W_h S_{hy}^2 \right) c'} \right]^2 - \frac{S_y^2}{N} \quad (10.3.13)$$

Example 10.3.1

A pilot sample of 1000 households was selected from 10,000 households in a certain region by SRSWOR method. On the basis of the pilot survey report, the selected households were classified into three categories: The high, middle, and low income groups. From each of the income groups subsamples were selected by SRSWOR method. The sample mean, sample

variance of the monthly income, and sample proportion of unemployed persons for the final survey are presented in the following table.

Income group	Pilot survey	Final survey			
		No. of households	No. of households	Sample mean (monthly income in \$)	Sample variance (monthly income)
h	n'_h	n_h	$\bar{y}(s_h)$	s_{hy}^2	$\hat{\pi}_h$
Low	325	100	1575	4000	0.75
Middle	475	75	5750	3500	0.40
High	200	75	8000	2500	0.20

Estimate the average income and proportion of unemployed persons in the region along with their standard error.

Estimation of average income and its standard error

h	$w_h = n'_h/n'$	$w_h \bar{y}(s_h)$	$w_h (\bar{y}(s_h) - \hat{\bar{Y}}_{stt})^2$	$\left(\frac{n'_h-1}{n'-1} - \frac{n_h-1}{N-1}\right) \frac{w_h s_{hy}^2}{n_h}$
1	0.325	511.875	3,471,208.330	4.087
2	0.475	2731.250	390,650.576	10.353
3	0.200	1600.000	1,993,171.953	1.278
Total	1.000	4843.125	5,855,030.859	15.719

$$\text{Estimated average monthly income} = \hat{\bar{Y}}_{stt} = \sum_{h=1}^H w_h \bar{y}(s_h) = \$4843.125.$$

$$\text{Estimated variance of } \hat{\bar{Y}}_{stt} = \hat{V}(\hat{\bar{Y}}_{stt}) \cong \sum_{h=1}^H \left(\frac{n'_h-1}{n'-1} - \frac{n_h-1}{N-1} \right) \frac{w_h s_{hy}^2}{n_h}$$

$$+ \frac{N-n'}{N(n'-1)} \sum_{h=1}^H w_h (\bar{y}(s_h) - \hat{\bar{Y}}_{stt})^2$$

$$= 15.719 + \{(10,000 - 1000)/(10,000 \times 999)\} \times 5855030.859$$

$$= \$5290.521$$

$$\text{Estimated standard error of } \hat{\bar{Y}}_{stt} = \sqrt{\hat{V}(\hat{\bar{Y}}_{stt})} = \$72.74$$

Estimation of proportion of unemployed persons and its standard error

h	$w_h = n'_h/n'$	$w_h \hat{\pi}_h$	$w_h(\hat{\pi}_h - \hat{\pi})^2$	$\left(\frac{n'_h - 1}{n' - 1} - \frac{n_h - 1}{N - 1}\right) \frac{w_h \hat{\pi}_h(1 - \hat{\pi}_h)}{n_h - 1}$
1	0.325	0.244	0.0248	0.00019
2	0.475	0.190	0.0025	0.00072
3	0.200	0.040	0.0149	0.00008
Total	1.000	0.474	0.0423	0.00099

$$\text{Estimated proportion of unemployed persons} = \hat{\pi} = \sum_{h=1}^H w_h \hat{\pi}_h = 0.4737$$

$$\begin{aligned} \text{Estimated variance of } \hat{\pi} = \hat{V}(\hat{\pi}) &\cong \sum_{h=1}^H \left(\frac{n'_h - 1}{n' - 1} - \frac{n_h - 1}{N - 1} \right) \frac{w_h \hat{\pi}_h(1 - \hat{\pi}_h)}{n_h - 1} \\ &\quad + \frac{N - n'}{N(n' - 1)} \sum_{h=1}^H w_h(\hat{\pi}_h - \hat{\pi})^2 \\ &= 0.00099 + \{(10000 - 1000)/(10000 \times 999)\} \times 0.0423 = 0.0010 \end{aligned}$$

$$\text{Estimated standard error of } \hat{\pi} = \sqrt{\hat{V}(\hat{\pi})} = 0.032$$

10.4 TWO-PHASE SAMPLING FOR SELECTION OF SAMPLE

In case the auxiliary information (x) selected in the first phase is positive, it can be used as a measure of size for the selection of sample in the second phase. In this section we describe the use of PPSWR and Rao—Hartley—Cochran (1962, RHC) methods of sampling for the selection of samples in the second phase when the first-phase sample is selected by SRSWOR method.

10.4.1 Probability Proportional to Size With Replacement Sampling

Let us suppose that in the first phase, a sample s' of size n' be selected by SRSWOR method and the values of the auxiliary variable x are obtained. In the second phase, a subsample s of size n is selected from s' by using PPSWR method using x 's as measure of size, i.e., the normed size measure for the i th unit is $p'_i = x_i/X'$ for $i \in s'$, where $X' = \sum_{i \in s'} x_i$. The estimator for the population total Y is given by

$$\hat{Y}_{ppst} = \frac{N}{n'} \left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} \right) \quad (10.4.1)$$

Theorem 10.4.1

- (i) \hat{Y}_{ppst} is an unbiased estimator for Y .
(ii) Variance of \hat{Y}_{ppst} is

$$V(\hat{Y}_{ppst}) = N^2 \left[\frac{(n' - 1)}{N(N - 1)n'} \sum_{i \in U} p_i \left(\frac{y_i}{p_i} - Y \right)^2 + \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 \right]$$

where $p_i = x_i/X$ and $X = \sum_{i \in U} x_i$.

- (iii) An unbiased estimator of $V(\hat{Y}_{ppst})$ is

$$\begin{aligned} \hat{V}(\hat{Y}_{ppst}) = N^2 & \left[\frac{N - 1}{N(n' - 1)n'} \hat{V} \left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} \middle| s' \right) \right. \\ & \left. + \frac{N - n'}{N(n' - 1)} \left(\frac{1}{n'n} \sum_{i \in s} \frac{y_i^2}{p'_i} - \frac{\hat{Y}_{ppst}^2}{N^2} \right) \right] \end{aligned}$$

$$\text{where } \hat{V} \left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} \middle| s' \right) = \frac{1}{n(n - 1)} \sum_{i \in s} \left(\frac{y_i}{p'_i} - \frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} \right)^2.$$

Proof

$$\begin{aligned} \text{(i)} \quad E(\hat{Y}_{ppst}) &= E \left[\frac{N}{n'} E \left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} \middle| s' \right) \right] \\ &= NE \left(\frac{1}{n'} \sum_{i \in s'} y_i \right) \\ &= NE(\bar{y}(s')) = Y \end{aligned}$$

$$\begin{aligned} \text{(ii)} \quad V(\hat{Y}_{ppst}) &= E \left[\frac{N^2}{n'^2} V \left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} \middle| s' \right) \right] + V \left[\frac{N}{n'} E \left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} \middle| s' \right) \right] \\ &= E \left[\frac{N^2}{n'^2} \left\{ \frac{1}{2n} \sum_{i \neq j} \sum_{j \in s'} p'_i p'_j \left(\frac{y_i}{p'_i} - \frac{y_j}{p'_j} \right)^2 \right\} \right] + V(N \bar{y}(s')) \\ &= \frac{N^2}{n'^2} \frac{1}{2n} \sum_{i \neq j} \sum_{j \in U} p_i p_j \left(\frac{y_i}{p_i} - \frac{y_j}{p_j} \right)^2 \pi_{ij}^0 + N^2 \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 \end{aligned}$$

[where π_{ij}^0 = inclusion probability of i th and j th unit ($i \neq j$) in $s' = n'(n' - 1)/\{N(N - 1)\}$]

$$= N^2 \left[\frac{(n' - 1)}{n'nN(N - 1)} \sum_{i \in U} p_j \left(\frac{y_i}{p_i} - Y \right)^2 + \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 \right]$$

(iii) Unbiased estimators of $V\left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} | s'\right)$ and S_y^2 are given, respectively, by

$$\widehat{V}\left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} | s'\right) = \frac{1}{n(n-1)} \sum_{i \in s} \left(\frac{y_i}{p'_i} - \frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} \right)^2$$

and

$$\widehat{S}_y^2 = \frac{1}{(N-1)} \left\{ \frac{N}{n'n} \sum_{i \in s} \frac{y_i^2}{p'_i} - \frac{\widehat{Y}_{ppst}^2 - \widehat{V}(\widehat{Y}_{ppst})}{N} \right\}$$

Then an unbiased estimator of $V(\widehat{Y}_{ppst})$ is obtained from the following equation

$$\begin{aligned} \widehat{V}(\widehat{Y}_{ppst}) &= N^2 \left[\frac{1}{n'^2} \widehat{V}\left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} | s'\right) + \left(\frac{1}{n'} - \frac{1}{N}\right) \widehat{S}_y^2 \right] \\ &= N^2 \left[\frac{1}{n'^2} \widehat{V}\left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} | s'\right) \right. \\ &\quad \left. + \left(\frac{1}{n'} - \frac{1}{N}\right) \frac{1}{(N-1)} \left\{ \frac{N}{n'} \frac{1}{n} \sum_{i \in s} \frac{y_i^2}{p'_i} - \frac{\widehat{Y}_{ppst}^2 - \widehat{V}(\widehat{Y}_{ppst})}{N} \right\} \right] \end{aligned} \quad (10.4.2)$$

Finally, Eq. (10.4.2) yields

$$\begin{aligned} \widehat{V}(\widehat{Y}_{ppst}) &= N^2 \left[\frac{N-1}{N(n'-1)n'} \widehat{V}\left(\frac{1}{n} \sum_{i \in s} \frac{y_i}{p'_i} | s'\right) \right. \\ &\quad \left. + \frac{N-n'}{N(n'-1)} \left(\frac{1}{n'n} \sum_{i \in s} \frac{y_i^2}{p'_i} - \frac{\widehat{Y}_{ppst}^2}{N^2} \right) \right] \end{aligned}$$

10.4.2 Rao—Hartley—Cochran Sampling

Avadhani and Sukhatme (1970) considered two-phase sampling, where in the first phase, a sample s' of size n' is selected by SRSWOR method and in the second phase, a subsample s of size n is selected from s' by using the RHC method (see Section 5.6) using normed size measure $p'_i = x_i/X'$ for the i th unit, $i \in s'$. Let P'_i be the sum of the p' values of the group

containing the i th unit in selecting the subsample s by RHC scheme from s' . The proposed estimator for Y is

$$\hat{Y}_{rhc} = \frac{N}{n'} \sum_{j \in s} \frac{y_j}{p_j} P'_j \quad (10.4.3)$$

Theorem 10.4.2

$$\begin{aligned} \text{(i)} \quad & E(\hat{Y}_{rhc}) = Y \\ \text{(ii)} \quad & V(\hat{Y}_{rhc}) = N^2 \left[\frac{(n' - n)}{n' n N (N - 1)} \sum_{i \in U} p_i \left(\frac{y_i}{p_i} - Y \right)^2 + \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 \right] \\ \text{(iii)} \quad & \hat{V}(\hat{Y}_{rhc}) = N^2 \left[\frac{N - 1}{n' N (n' - 1)} \hat{V} \left(\sum_{j \in s} \frac{y_j}{p_j} P'_j \right) \right. \\ & \quad \left. + \frac{N - n'}{N (n' - 1)} \left(\frac{1}{n'} \sum_{i \in s} \frac{y_i^2}{p_i} P'_i - \frac{\hat{Y}_{rhc}^2}{N^2} \right) \right] \end{aligned}$$

$$\text{where } \hat{V} \left(\sum_{j \in s} \frac{y_j}{p_j} P'_j \right) = \frac{n' - n}{n' (n - 1)} \sum_{j \in s} P'_j \left(\frac{y_j}{p'_j} - \sum_{j \in s} \frac{y_j}{p'_j} P'_j \right)^2.$$

Proof

$$\begin{aligned} \text{(i)} \quad & E(\hat{Y}_{rhc}) = E \left[\frac{N}{n'} E \left(\sum_{j \in s} \frac{y_j}{p_j} P'_j \middle| s' \right) \right] \\ & = E \left(\frac{N}{n'} \sum_{j \in s'} y_j \right) \\ & = Y \\ \text{(ii)} \quad & V(\hat{Y}_{rhc}) = E \left[\frac{N^2}{n'^2} V \left(\sum_{j \in s} \frac{y_j}{p_j} P'_j \middle| s' \right) \right] + V \left[\frac{N}{n'} E \left(\sum_{j \in s} \frac{y_j}{p_j} P'_j \middle| s' \right) \right] \\ & = \frac{N^2}{n'^2} \frac{n' - n}{n (n' - 1)} \frac{1}{2} E \sum_{i \neq j} \sum_{j \in s'} p_i p_j \left(\frac{y_i}{p_i} - \frac{y_j}{p_j} \right)^2 + N^2 V \left(\frac{1}{n'} \sum_{j \in s'} y_j \right) \\ & = N^2 \left[\frac{(n' - n)}{n' n N (N - 1)} \sum_{i \in U} p_i \left(\frac{y_i}{p_i} - Y \right)^2 + \left(\frac{1}{n'} - \frac{1}{N} \right) S_y^2 \right] \end{aligned}$$

(iii) Now, noting $\widehat{V} \left(\sum_{j \in s} \frac{y_j}{p'_j} P'_j | s' \right) = \frac{n' - n}{n'(n - 1)} \sum_{j \in s} P'_j \left(\frac{y_j}{p'_j} - \sum_{j \in s} \frac{y_j}{p'_j} P'_j \right)^2$

and

$\widehat{S}_{rht}^2 = \frac{1}{N - 1} \left\{ \frac{N}{n'} \sum_{i \in s} \frac{y_i^2}{p'_i} P'_i - \frac{\widehat{Y}_{rht}^2 - \widehat{V}(\widehat{Y}_{rht})}{N} \right\}$ are unbiased estimators of $V \left(\sum_{j \in s} \frac{y_j}{p'_j} P'_j | s' \right)$ and S_y^2 , respectively, we find $\widehat{V}(\widehat{Y}_{rht})$ as an unbiased estimator of $V(\widehat{Y}_{rht})$ from the following equation:

$$\widehat{V}(\widehat{Y}_{rht}) = \frac{N^2}{n'^2} \widehat{V} \left(\sum_{j \in s} \frac{y_j}{p'_j} P'_j | s' \right) + N^2 \left(\frac{1}{n'} - \frac{1}{N} \right) \widehat{S}_{rht}^2$$

$$\text{i.e., } \widehat{V}(\widehat{Y}_{rht}) = N^2 \left[\frac{1}{n'^2} \widehat{V} \left(\sum_{j \in s} \frac{y_j}{p'_j} P'_j | s' \right) + \left(\frac{1}{n'} - \frac{1}{N} \right) \frac{1}{N - 1} \left\{ \frac{N}{n'} \sum_{i \in s} \frac{y_i^2}{p'_i} P'_i - \frac{\widehat{Y}_{rht}^2 - \widehat{V}(\widehat{Y}_{rht})}{N} \right\} \right]$$

$$\text{i.e., } \widehat{V}(\widehat{Y}_{rht}) = N^2 \left[\frac{N - 1}{n'N(n' - 1)} \widehat{V} \left(\sum_{j \in s} \frac{y_j}{p'_j} P'_j | s' \right) + \frac{N - n'}{N(n' - 1)} \left(\frac{1}{n'} \sum_{i \in s} \frac{y_i^2}{p'_i} P'_i - \frac{\widehat{Y}_{rht}^2}{N^2} \right) \right]$$

10.5 TWO-PHASE SAMPLING FOR STRATIFICATION AND SELECTION OF SAMPLE

Arnab (1991) considered a two-phase sampling procedure, where in the first phase, a PPSWR sample s' of size n' is selected from the entire population U with known size measures z_i 's (>0 for $i = 1, \dots, N$), and information of the auxiliary variable x_i 's ($i \in s'$) is obtained. By using the information of the x_i 's gathered in the initial sample s' , the n' sampled units are assigned to a predetermined number of L strata. Let \widetilde{s}_h be the sample of size n'_h corresponding to the h th stratum U_h ($h = 1, \dots, L$). In the second phase, subsamples s_h 's of sizes $n_h = \gamma_h n'_h$ (assuming integer and $0 \leq \gamma_h \leq 1$) are selected independently from each of the h th stratum by PPSWR method with normed size measure $q_{hi} = (x_{hi}/z_{hi}) / \sum_{i \in s'_h} x_{hi}/z_{hi}$ for the i th

unit ($i \in \tilde{s}_h$), where z_{hi} and x_{hi} are the values of the variables z and x for the i th unit of the h th stratum, respectively. Here it is assumed that n' is so large that $\text{Prob}(n'_h \geq 1) = 1$. The proposed estimator for the population total Y is

$$\hat{Y}_A = \sum_h w_h \hat{Y}_{Ah} \quad (10.5.1)$$

where
$$\hat{Y}_{Ah} = \frac{1}{n_h} \sum_{j \in s_h} \frac{\gamma_{hj} - \lambda_h x_{hj}}{(n'_h p_{hj}) q_{hj}} + \lambda_h \frac{1}{n'_h} \sum_{j \in \tilde{s}_h} \frac{x_{hj}}{p_{hj}}, \quad p_{hj} = z_{hj}/Z,$$

$Z = \sum_{i \in U} z_i$, λ_h is a known constant, and $\sum_{j \in s}$ denotes sum over the units in s with repetition.

Theorem 10.5.1

(i) $E(\hat{Y}_A) = Y$

(ii)
$$V(\hat{Y}_A) = \frac{1}{n'} \left[\sum_h \left(1 - \frac{1}{n' P_h} \right) \frac{\sigma_h^2(d_{\lambda h}|x)}{P_h \gamma_h} + \sigma^2(y|z) \right]$$

where $w_h = n'_h/n'$, $P_h = Z_h/Z$, $Z_h = \sum_{i \in U_h} z_{hi}$, $\sigma_h^2(d_{\lambda h}|x) = \sum_{j \in U_h} \frac{d_{\lambda hj}^2}{x_{hi}} X_h - D_{\lambda h}^2$,
 $X_h = \sum_{i \in U_h} x_{hi}$, $d_{\lambda hj} = \gamma_{hj} - \lambda_h x_{hj}$, $D_{\lambda h} = \sum_{j \in U_h} d_{\lambda hj} = Y_h - \lambda_h X_h$, and

$$\sigma^2(y|z) = \sum_{h=1}^H \sum_{j \in U_h} \frac{\gamma_{hj}^2}{p_{hj}} - Y^2.$$

Proof

(i)
$$\begin{aligned} E(\hat{Y}_A) &= E \sum_h w_h E(\hat{Y}_{Ah} | s'_h) \\ &= E \sum_h w_h \sum_{j \in \tilde{s}_h} \frac{\gamma_{hj}}{(n'_h p_{hj})} \\ &= E \sum_h \frac{w_h}{P_h} E \left(\sum_{j \in \tilde{s}_h} \frac{\gamma_{hj} P_h}{(n'_h p_{hj})} \middle| \mathbf{w}_h \right) \\ &\quad (\text{where } \mathbf{w}_h = (w_1, \dots, w_L)) \\ &= E \sum_h \frac{w_h}{P_h} Y_h \\ &= Y \\ &\quad (\text{noting } E(w_h) = P_h) \end{aligned}$$

$$(ii) \quad V(\hat{Y}_A) = E \left[\sum_h w_h^2 V(\hat{Y}_{Ah} | \tilde{s}_h) \right] + V \left[\sum_h w_h E(\hat{Y}_{Ah} | \tilde{s}_h) \right]$$

Now,

$$\begin{aligned} V \left\{ \sum_h w_h E(\hat{Y}_{Ah} | \tilde{s}_h) \right\} &= V \left\{ \sum_h w_h \sum_{j \in s_h} \frac{y_{hj}}{\left(\frac{n'_h}{p_{hj}} \right)} \right\} \\ &= E \left\{ \sum_h \frac{w_h^2}{P_h^2} V \left(\frac{1}{n'_h} \sum_{j \in \tilde{s}_h} \frac{y_{hj}}{p_{hj}} | \mathbf{w}_h \right) \right\} \\ &\quad + V \left\{ \sum_h \frac{w_h}{P_h} E \left(\frac{1}{n'_h} \sum_{j \in \tilde{s}_h} \frac{y_{hj}}{p_{hj}} | \mathbf{w}_h \right) \right\} \\ &\quad (\text{where } p'_{hi} = p_{hi}/P_h) \\ &= E \left[\sum_h \frac{w_h^2}{P_h^2} \frac{1}{n'_h} \left(\sum_{j \in U_h} \frac{y_{hj}^2}{p'_{hj}} - Y_h^2 \right) + V \left(\sum_h \frac{w_h}{P_h} Y_h \right) \right] \\ &= \frac{1}{n'} \left[\sum_h \frac{1}{P_h} \left(\sum_{j \in U_h} \frac{y_{hj}^2}{p'_{hj}} - Y_h^2 \right) + \left(\sum_h \frac{Y_h^2}{P_h} - Y^2 \right) \right] \\ &= \frac{1}{n'} \left(\sum_{h=1}^H \sum_{j \in U_h} \frac{y_{hj}^2}{p_{hj}} - Y^2 \right) \\ &= \frac{\sigma^2(y|z)}{n'} \end{aligned} \tag{10.5.2}$$

and

$$\begin{aligned}
E \left\{ \sum_h w_h^2 V(\hat{Y}_{Ah} | \tilde{s}_h) \right\} &= E \left\{ \sum_h w_h^2 V \left(\sum_{j \in s_h} \frac{\gamma_{hj} - \lambda_h x_{hj}}{\left(\frac{n'_h p_{hj}}{n_h q_{hj}} \right)} \middle| \tilde{s}_h \right) \right\} \\
&= E \sum_h \frac{w_h^2}{n_h} \left\{ \sum_{j \in \tilde{s}_h} \frac{d_{\lambda hj}^2}{\left(\frac{n'_h p_{hj}}{n_h q_{hj}} \right)^2} - \left(\sum_{j \in \tilde{s}_h} \frac{d_{\lambda hj}}{n'_h p_{hj}} \right)^2 \right\} \\
&= E \sum_h \frac{w_h^2}{n_h} E \left\{ \sum_{j \in \tilde{s}_h} \frac{d_{\lambda hj}^2}{\left(\frac{n'_h p_{hj}}{n_h q_{hj}} \right)^2} \frac{p_{hj}}{q'_{hj}} \sum_{j \in s_h} \frac{q'_{hj}}{p_{hj}} - \left(\sum_{j \in \tilde{s}_h} \frac{d_{\lambda hj}}{n'_h p_{hj}} \right)^2 \middle| \mathbf{w}_h \right\} \\
&\quad \left(\text{where } q'_{hj} = x_{hj} / X_h \right) \\
&= E \sum_h \frac{w_h^2}{n_h} E \left\{ \sum_{j \in \tilde{s}_h} \frac{d_{\lambda hj}^2}{\left(\frac{n'_h p_{hj}}{n_h q_{hj}} \right)^2} + \sum_{j \in \tilde{s}_h} \sum_{k(\neq j) \in \tilde{s}_h} \frac{d_{\lambda hj}^2}{\left(\frac{n'_h p_{hj}}{n_h q_{hj}} \right)^2} \frac{p_{hj}}{q'_{hj}} \frac{q'_{hk}}{p_{hk}} \right. \\
&\quad \left. - \left(\sum_{j \in \tilde{s}_h} \frac{d_{\lambda hj}}{n'_h p_{hj}} \right)^2 \middle| \mathbf{w}_h \right\} \\
&= E \sum_h \frac{w_h^2}{n_h} \left[\frac{1}{(n'_h)^2} \left\{ \frac{n'_h}{P_h} \sum_{j \in U_h} \frac{d_{\lambda hj}^2}{p_{hj}} + \frac{n'_h(n'_h - 1)}{P_h^2} \left(\sum_{j \in U_h} \frac{d_{\lambda hj}^2}{q'_{hj}} \right) \right\} \right. \\
&\quad \left. - \frac{1}{P_h^2} \left\{ \frac{1}{n'_h} \left(\sum_{j \in U_h} \frac{d_{\lambda hj}^2}{p'_{hj}} - D_{\lambda h}^2 \right) + D_{\lambda h}^2 \right\} \right] \\
&\quad \left(\text{noting } \sum_{j \in U_h} q'_{hj} = 1 \right) \\
&= E \sum_h \frac{w_h^2}{n_h} \frac{n'_h - 1}{n'_h P_h^2} \left(\sum_{j \in U_h} \frac{d_{\lambda hj}^2}{x_{hj}} X_h - D_{\lambda h}^2 \right) \\
&\quad \left(\text{noting } p'_{hj} = \frac{z_{hj}}{Z_h} = \frac{p_{hj}}{P_h} \text{ and } q'_{hj} = \frac{x_{hj}}{X_h} \right)
\end{aligned}$$

Finally, writing $n_h = \gamma_h n'_h$, we get

$$E \sum_h w_h^2 V(\hat{Y}_{Ah} | \tilde{s}_h) = \sum_h \left(1 - \frac{1}{n' P_h}\right) \frac{\sigma_h^2(d_{Ah} | x)}{n' \gamma_h P_h} \quad (10.5.3)$$

The proof (ii) follows from Eqs. (10.5.2) and (10.5.3).

10.6 EXERCISES

10.6.1 Let an initial sample s' of size n' be selected from a finite population of size N by SRSWOR method and information on the auxiliary variable x is collected. Then an independent sample s of size n is selected by SRSWOR method from the entire population, and information on the auxiliary (x) and study variable (y) is obtained. Let $t = \bar{y}(s) - k\{\bar{x}(s) - \bar{x}(s')\}$.

Show that

(i) $E(t) = \bar{Y}$,

(ii) $V(t) = \left(\frac{1}{n} - \frac{1}{N}\right) (S_y^2 + k^2 S_x^2 - 2k S_{xy}) + k^2 \left(\frac{1}{n'} - \frac{1}{N}\right) S_x^2$

(iii) $\hat{V}(t) = \left(\frac{1}{n} - \frac{1}{N}\right) (s_y^2 + k^2 s_x^2 - 2k s_{xy}) + k^2 \left(\frac{1}{n'} - \frac{1}{N}\right) s_x^2$ is an unbiased estimator for $V(t)$

(iv) Find the value of k for which $V(t)$ attains a minimum (Raj, 1968).

10.6.2 Initially, a sample s' of size n' is selected from a finite population of size N by PPSWR method using p_i as a normed size measure for the i th unit, and information on the auxiliary variable x is collected. Then an independent sample s of size n is selected by PPSWR method using the same normed size measure for the i th unit, and information on the auxiliary (x) and study variable (y) is obtained.

Let $t = \frac{1}{n} \sum_{i \in s} \frac{y_i}{p_i} - k \left(\frac{1}{n} \sum_{i \in s} \frac{x_i}{p_i} - \frac{1}{n'} \sum_{i \in s'} \frac{x_i}{p_i} \right)$ and $\sum_{i \in s}$ is

the sum over units in s with repetition. Show that

(i) $E(t) = \bar{Y}$

(ii) $V(t) = \frac{\sigma^2(d_k | p)}{n} + k^2 \frac{\sigma^2(x | p)}{n'}$, where $\sigma^2(d_k | p) = \sum_{i \in U} \frac{d_{ki}^2}{p_i} - D_k^2$,

$$\sigma^2(x | p) = \sum_{i \in U} \frac{x_i^2}{p_i} - X^2, \quad d_{ki} = y_i - kx_i, \quad \text{and} \quad D_k = \sum_{i \in U} d_{ki}$$

(iii) $\hat{V}(t) = \frac{\hat{\sigma}^2(d_k | p)}{n} + k^2 \frac{\hat{\sigma}^2(x | p)}{n'}$ is an unbiased estimator for $V(t)$, where

$$\hat{\sigma}^2(d_k|p) = \frac{\sum_{i \in s} \left(\frac{d_{ki}}{p_i} - \frac{1}{n} \sum_{i \in s} \frac{d_{ki}}{p_i} \right)^2}{(n-1)} \text{ and } \hat{\sigma}^2(x|p) = \frac{\sum_{i \in s'} \left(\frac{x_i}{p_i} - \frac{1}{n'} \sum_{i \in s'} \frac{x_i}{p_i} \right)^2}{(n'-1)}$$

(iv) Find the value of k for which $V(t)$ attains the minimum (Raj, 1968).

- 10.6.3** In the first phase an initial sample s' of size n' is selected by SRSWOR method, and information on the auxiliary variable x is collected. At the second phase, a subsample s of size n is selected from s' by Lahiri—Midzuno—Sen (1951, 1952, 1953) sampling procedure using x as a measure of size.

$$\text{Show that } T = N \frac{\sum_{i \in s} y_i}{\sum_{i \in s} x_i} \left(\sum_{i \in s'} x_i / n' \right) \text{ is an unbiased estimator}$$

for the population total Y . Derive the variance of T and find an unbiased estimator of $V(T)$ (Raj, 1968).

- 10.6.4** In the first phase a sample s' of size n' is selected by RHC sampling scheme using a normed size measure p_i attached to the i th unit. At the second phase a subsample s of size n is selected from s' by SRSWOR method. Let $T = \frac{n'}{n} \sum_{i \in s} \frac{y_i}{p_i} P_i - k \left(\frac{n'}{n} \sum_{i \in s} \frac{x_i}{p_i} P_i - \sum_{i \in s'} \frac{x_i}{p_i} P_i \right)$, where P_i is the sum of p_j values for the group (which was formed for the selection of the sample) containing the i th unit.

(a) Verify the following results:

(i) $E(T) = Y$

(ii)
$$V(T) = \left(\frac{1}{n} - \frac{1}{n'} \right) \frac{1}{(N-1)} [n' \overline{V}(d_k) + (N - n') \sigma^2(d_k|p)]$$

$$+ \frac{N - n'}{n'(N-1)} \sigma^2(y|p)$$

where $d_{ki} = y_i - kx_i$, $\sigma^2(y|p) = \sum_{i \in U} \frac{y_i^2}{p_i} - Y^2$, $\sigma^2(d_k|p) =$

$\sum_{i \in U} \frac{d_{ki}^2}{p_i} - D_k^2$, $D_k = \sum_{i \in U} d_{ki}$, and $\overline{V}(d_k) = N \sum_{i \in U} \frac{d_{ki}^2}{p_i} - D_k^2$.

(b) Find an unbiased estimator of $\widehat{V}(T)$ and the optimum value of value of k that minimizes $V(T)$ (Arnab, 1979).

- 10.6.5** Let π_{1i} and π_{1ij} be the first-order and second-order inclusion probabilities for the first-phase sample s_1 , and $\pi_{2ij|s_1}$ and $\pi_{2ij|s_1}$ be the

conditional first-order and second-order inclusion probabilities for the second-phase sample s_2 given s_1 . Show that

(i) $\hat{Y}_{DE} = \sum_{i \in s_2} \frac{y_i}{\pi_i^*}$ with $\pi_i^* = \pi_{1i}\pi_{2i|s_1}$ is an unbiased estimator of the

total Y .

(ii) The variance of \hat{Y}_{DE} can be estimated unbiasedly by using any of the following estimators:

$$(a) \quad \hat{V}_1 = \sum_{i \in s_2} \frac{\Delta_{1ij}}{\pi_{ij}^*} \frac{y_i}{\pi_{1i}} \frac{y_j}{\pi_{1j}} + \sum_{i \in s_2} \sum_{j \in s_2} \frac{\Delta_{2ij|s_1}}{\pi_{ij|s_1}} \frac{y_i}{\pi_i^*} \frac{y_j}{\pi_j^*}$$

and

$$(b) \quad \hat{V}_2 = \sum_{i \in s_2} \sum_{\substack{j \in s_2 \\ i < j}} \frac{\Delta_{1ij}}{\pi_{ij}^*} \left(\frac{y_i}{\pi_{1i}} - \frac{y_j}{\pi_{1j}} \right)^2 + \sum_{i \in s_2} \sum_{\substack{j \in s_2 \\ i < j}} \frac{\Delta_{2ij|s_1}}{\pi_{ij|s_1}} \left(\frac{y_i}{\pi_i^*} - \frac{y_j}{\pi_j^*} \right)^2$$

where $\pi_{ij}^* = \pi_{1ij}\pi_{2ij|s_1}$, $\Delta_{1ij} = \pi_{1i}\pi_{1j} - \pi_{1ij}$, and $\Delta_{2ij|s_1} = \pi_{2i|s_1}\pi_{2j|s_1} - \pi_{2ij|s_1}$ (Haziza et al., 2011).

10.6.6 To estimate the average price of a two-bedroom flat in a locality, a sample of 75 two-bedroom flats was selected by SRSWOR from the locality of 790 flats and eye estimates of the price of the selected flats were obtained. In phase 2, a random sample of 15 flats was selected from the selected 75 flats by SRSWOR method and the prices of the flats were determined by an appropriate evaluation agency. The findings are given in the following table:

Serial number of flats	Price of the flat by eye estimation (in 000 \$) (x)	Price of the flat by valuation agency (in 000 \$) (y)
1	450	550
2	565	600
3	700	650
4	850	900
5	650	750
6	700	820
7	850	970
8	550	700
9	600	650
10	750	700
11	900	850
12	850	950
13	840	900
14	700	850
15	760	800

Estimate the average price of the flat along with its standard error. If the cost evaluation of a flat by an agency is five times that of the eye estimation, determine the proportion of the second-phase sample that should be selected from the first-phase sample.

10.6.7 To estimate the average daily wages of a factory of 800 workers, an initial sample of 125 workers was selected at random by SRSWOR method. The workers were then stratified according to their skill. Finally, samples of workers were selected from each of the strata by SRSWOR method, and information regarding the daily wages was obtained. The data are given below.

- (i) Give an estimate of the average daily wages of the factory workers.
- (ii) Estimate the standard error of your estimator.

Workers type	No. of workers in the initial sample	No. of workers in the final sample	Daily wages (in \$)
Highly skilled	10	5	625, 750, 1025, 1500, 1000
Skilled	60	10	320, 350, 375, 200, 300, 325, 325, 475, 350, 300
Semiskilled	40	8	200, 175, 150, 240, 250, 180, 120, 170
Unskilled	20	8	150, 175, 100, 140, 200, 160, 100, 100