

## CHAPTER 9

# Regression, Product, and Calibrated Methods of Estimation

### 9.1 INTRODUCTION

We have seen in Chapter 8 that the ratio estimator is optimal if the regression of  $y$  on  $x$  is linear and passing through the origin. The ratio estimator may not be efficient if the regression line does not pass through the origin. In this situation, one may use auxiliary information effectively through the regression estimator to improve precision of the conventional estimator. Like the ratio estimator, the regression estimator is not unbiased for the population mean or total. In this section we will study the difference, regression, product, calibrated, and related estimators along with their properties in detail.

### 9.2 DIFFERENCE ESTIMATOR

If the population total  $X$  of the auxiliary variable  $x$  is known, then the estimator

$$\hat{Y}_\lambda = \hat{Y} - \lambda(\hat{X} - X) \quad (9.2.1)$$

is unbiased for the population total  $Y$  for any known value of  $\lambda$ . The estimator (Eq. 9.2.1) is known as a difference estimator. To use the difference estimator one needs to find the optimum value of  $\lambda$ , which minimizes the variance of  $\hat{Y}_\lambda$ .

Now minimizing the variance of  $\hat{Y}_\lambda$  with respect to  $\lambda$ , the optimum value of  $\lambda$  comes out as

$$\lambda_{opt} = \frac{Cov(\hat{X}, \hat{Y})}{V(\hat{X})} = \beta \quad (\text{say}) \quad (9.2.2)$$

Putting  $\lambda = \beta$  in Eq. (9.2.1), we get the optimum estimator

$$\hat{Y}_\beta = \hat{Y} - \beta(\hat{X} - X) \quad (9.2.3)$$

The variance of the optimum estimator  $\hat{Y}_\beta$  is given by

$$V(\hat{Y}_\beta) = (1 - \rho_{\hat{X}, \hat{Y}}^2) V(\hat{Y}) \quad (9.2.4)$$

where  $\rho_{\hat{X}, \hat{Y}}$  is the correlation between  $\hat{X}$  and  $\hat{Y}$ .

Because  $|\rho_{\hat{X}, \hat{Y}}| \leq 1$ ,  $\hat{Y}_\beta$  possesses smaller variance than the convention estimator  $\hat{Y}$ . The estimator  $\hat{Y}_\beta$  is seldom used in practice because the value of  $\beta$  is unknown in most practical situations. Let  $\beta_0$  be a guess value of  $\beta$  obtained from a past experiment. Then substituting  $\beta = \beta_0$  in Eq. (9.2.3), one gets the estimator

$$\hat{Y}_{\beta_0} = \hat{Y} - \beta_0(\hat{X} - X) = \hat{Y}_\beta + (\beta - \beta_0)(\hat{X} - X) \quad (9.2.5)$$

The loss of efficiency of using  $\hat{Y}_{\beta_0}$  over the optimum estimator  $\hat{Y}_\beta$  comes out as

$$\begin{aligned} \frac{V(\hat{Y}_{\beta_0})}{V(\hat{Y}_\beta)} - 1 &= \frac{(\beta - \beta_0)^2 V(\hat{X})}{V(\hat{Y}_\beta)} \\ &= \frac{(\beta - \beta_0)^2 V(\hat{X})}{(1 - \rho_{\hat{X}, \hat{Y}}^2) V(\hat{Y})} \\ &= \left(1 - \frac{\beta_0}{\beta}\right)^2 \frac{\rho_{\hat{X}, \hat{Y}}^2}{1 - \rho_{\hat{X}, \hat{Y}}^2} \end{aligned} \quad (9.2.6)$$

To keep the loss of efficiency to a certain level  $\gamma$  (say), we must have

$$\left(1 - \frac{\beta_0}{\beta}\right)^2 < \gamma \frac{1 - \rho_{\hat{X}, \hat{Y}}^2}{\rho_{\hat{X}, \hat{Y}}^2} \quad (9.2.7)$$

For simple random sampling without replacement (SRSWOR) and simple random sampling with replacement (SRSWR), Eq. (9.2.7) reduces to

$$\left(1 - \frac{\beta_0}{\beta}\right)^2 < \gamma \frac{1 - \rho^2}{\rho^2} \quad (9.2.8)$$

where  $\rho$  is the correlation coefficient between  $x$  and  $y$ .

In particular if  $\gamma = 0.25$ ,  $\rho = 0.5$ , Eq. (9.2.8) reduces to  $0.134 < \frac{\beta_0}{\beta} < 1.866$ .

For the value of  $\lambda = 1$ , the estimator  $\hat{Y}_\lambda$  reduces to a very simple estimator viz.

$$\hat{Y}_1 = \hat{Y} - (\hat{X} - X) \quad (9.2.9)$$

This estimator  $\hat{Y}_1$  is unbiased for  $Y$  and becomes more efficient than the conventional estimator  $\hat{Y}$  if  $\rho_{\hat{X}, \hat{Y}} > \frac{1}{2} \sqrt{\frac{V(\hat{Y})}{V(\hat{X})}}$ . In particular if the variances of  $\hat{X}$  and  $\hat{Y}$  are equal, then the difference estimator given in Eq. (9.2.9) performs better than the conventional estimator  $\hat{Y}$  if  $\rho_{\hat{X}, \hat{Y}} > 1/2$ . For SRSWOR and SRSWR sampling,  $\hat{X} = N\bar{x}_s$  and  $\hat{Y} = N\bar{y}_s$ . Here the condition of superiority of the difference estimator over the conventional estimator reduces to  $\rho > 1/2$  and it is realized in most practical situations.

### 9.3 REGRESSION ESTIMATOR

The value of  $\beta$  in Eq. (9.2.2) is generally unknown and it is estimated by

$$\hat{\beta} = \frac{C\hat{o}v(\hat{X}, \hat{Y})}{\hat{V}(\hat{X})} \quad (9.3.1)$$

where  $C\hat{o}v(\hat{X}, \hat{Y})$  and  $\hat{V}(\hat{X})$  are unbiased for  $Cov(\hat{X}, \hat{Y})$  and  $V(\hat{X})$ , respectively.

Hence, replacing  $\beta$  by its estimator  $\hat{\beta}$  given in Eq. (9.3.1), we get the well known linear regression estimator (or simply regression estimator) for the total  $Y$  as

$$\hat{Y}_{reg} = \hat{Y} - \hat{\beta}(\hat{X} - X) \quad (9.3.2)$$

Like the ratio estimator, the regression estimator is not unbiased for  $Y$ , but it is more efficient than the conventional estimator  $\hat{Y}$ .

#### 9.3.1 Exact Expression of Bias

The bias of the regression estimator  $\hat{Y}_{reg}$  is given by

$$\begin{aligned} B(\hat{Y}_{reg}) &= E(\hat{Y}_{reg} - Y) \\ &= -E\{\hat{\beta}(\hat{X} - X)\} \\ &= -Cov(\hat{\beta}, \hat{X}) \end{aligned} \quad (9.3.3)$$

Because in most situations,  $Cov(\hat{\beta}, \hat{X})$  decreases as the sample size increases, the regression estimator is approximately unbiased when the sample size is large.

### 9.3.2 Approximate Expression of Bias

Now writing  $\varepsilon_{xy} = \frac{C\hat{o}v(\hat{X}, \hat{Y})}{Cov(\hat{X}, \hat{Y})} - 1$  and  $\varepsilon_{xx} = \frac{\hat{V}(\hat{X})}{V(\hat{X})} - 1$ , we get

$$\begin{aligned}\hat{\beta} &= \frac{(1 + \varepsilon_{xy})Cov(\hat{X}, \hat{Y})}{(1 + \varepsilon_{xx})V(\hat{X})} \\ &= \beta(1 + \varepsilon_{xy})(1 + \varepsilon_{xx})^{-1} \\ &= \beta(1 + \varepsilon_{xy})(1 - \varepsilon_{xx} + \varepsilon_{xx}^2 - \dots)\end{aligned}\quad (9.3.4)$$

The expression of bias of the regression estimator comes out as

$$\begin{aligned}B(\hat{Y}_{reg}) &= E(\hat{Y}_{reg} - Y) \\ &= -\beta E[(1 + \varepsilon_{xy})(1 - \varepsilon_{xx} + \varepsilon_{xx}^2 - \dots)(\hat{X} - X)] \\ &= -\beta E[(\varepsilon_{xy} - \varepsilon_{xx} + \varepsilon_{xx}^2 - \varepsilon_{xy}\varepsilon_{xx} + \dots)(\hat{X} - X)]\end{aligned}$$

Now neglecting the terms containing  $E[\varepsilon_{xx}^i \varepsilon_{xy}^j (\hat{X} - X)]$  for  $i + j > 1$  (which is expected to be small), we get an expression of bias up to the first order of approximation as

$$\begin{aligned}B(\hat{Y}_{reg}) &\cong \beta E[(\varepsilon_{xx} - \varepsilon_{xy})(\hat{X} - X)] \\ &= \beta \left[ \frac{Cov\{\hat{V}(\hat{X}), \hat{X}\}}{V(\hat{X})} - \frac{Cov\{C\hat{o}v(\hat{X}, \hat{Y}), \hat{X}\}}{Cov(\hat{X}, \hat{Y})} \right]\end{aligned}\quad (9.3.5)$$

#### 9.3.2.1 Bias Under Simple Random Sampling Without Replacement

For an SRSWOR,  $\hat{X} = N\bar{x}_s$ ,  $V(\hat{X}) = N^2\left(\frac{1}{n} - \frac{1}{N}\right)S_x^2$ ,  $\hat{V}(\hat{X}) = N^2\left(\frac{1}{n} - \frac{1}{N}\right)s_x^2$ ,  $Cov(\hat{X}, \hat{Y}) = N^2\left(\frac{1}{n} - \frac{1}{N}\right)S_{xy}$ , and  $C\hat{o}v(\hat{X}, \hat{Y}) = N^2\left(\frac{1}{n} - \frac{1}{N}\right)s_{xy}$ . Hence the approximate bias of the regression estimator is obtained using Eq. (9.3.5) as follows:

$$\begin{aligned}B(\hat{Y}_{reg}) &\cong \beta \left[ \frac{Cov\left\{N^2\left(\frac{1}{n} - \frac{1}{N}\right)s_x^2, N\bar{x}_s\right\}}{\left\{N^2\left(\frac{1}{n} - \frac{1}{N}\right)S_x^2\right\}} - \frac{Cov\left\{N^2\left(\frac{1}{n} - \frac{1}{N}\right)s_{xy}, N\bar{x}_s\right\}}{\left\{N^2\left(\frac{1}{n} - \frac{1}{N}\right)S_{xy}\right\}} \right] \\ &= N\beta \left[ \frac{Cov(s_x^2, \bar{x}_s)}{S_x^2} - \frac{Cov(s_{xy}, \bar{x}_s)}{S_{xy}} \right]\end{aligned}$$

Now using the following results of Sukhatme (1944) (detailed derivation is given in [Appendix 9A](#)), we find

$$Cov(s_x^2, \bar{x}_s) = \frac{N(N-n)}{(N-1)(N-2)} \frac{\mu_{30}}{n} \text{ and}$$

$$Cov(s_{xy}, \bar{x}_s) = \frac{N(N-n)}{(N-1)(N-2)} \frac{\mu_{21}}{n}$$

where  $\mu_{ij} = \frac{1}{N} \sum_{i \in U} (x_i - \bar{X})^i (y_i - \bar{Y})^j$ .

The expression for the bias (Eq. 9.3.5) comes out as

$$B(\hat{Y}_{reg}) = \frac{N(N-n)}{(N-2)} \frac{\beta}{n} \left( \frac{\mu_{30}}{\mu_{20}} - \frac{\mu_{21}}{\mu_{11}} \right)$$

In case  $x$  and  $y$  follow bivariate normal distribution, we have  $\mu_{30} = \mu_{21} = 0$  and the bias of the regression estimator  $B(\hat{Y}_{reg})$  reduces to zero.

### 9.3.3 Approximate Expression of the Mean Square Error

The mean square error (MSE) of the regression estimator  $\hat{Y}_{reg}$  is given by

$$\begin{aligned} M(\hat{Y}_{reg}) &= E(\hat{Y}_{reg} - Y)^2 \\ &= E\left[(\hat{Y} - Y) - \hat{\beta}(\hat{X} - X)\right]^2 \\ &= E\left[(\hat{Y} - Y) - \beta(1 + \varepsilon_{xy})(1 - \varepsilon_{xx} + \varepsilon_{xx}^2 - \dots)(\hat{X} - X)\right]^2 \\ &= E\left[\{(\hat{Y} - Y) - \beta(\hat{X} - X)\} - \beta(\hat{X} - X) \right. \\ &\quad \left. \times \{\varepsilon_{xy} - \varepsilon_{xx} + \varepsilon_{xx}^2 - \varepsilon_{xy}\varepsilon_{xx} + \dots\}\right]^2 \end{aligned}$$

Neglecting terms  $E(\hat{X} - X)^2 \varepsilon_{xy}^i \varepsilon_{xx}^j$  for  $i + j \geq 2$ , the expression of MSE up to the first order of approximation is given by

$$\begin{aligned} M(\hat{Y}_{reg}) &\cong E[(\hat{Y} - Y) - \beta(\hat{X} - X)]^2 \\ &= V(\hat{Y}) + \beta^2 V(\hat{X}) - 2\beta Cov(\hat{X}, \hat{Y}) \end{aligned} \tag{9.3.6}$$

$$= (1 - \rho_{\hat{X}, \hat{Y}}^2) V(\hat{Y}) \tag{9.3.7}$$

Let  $\hat{X} = \sum_{i \in s} b_{si} x_i$  and  $\hat{Y} = \sum_{i \in s} b_{si} y_i$  be linear homogeneous unbiased estimators of the population totals  $X$  and  $Y$ , respectively, based on a sample

$s$  of size  $n$  selected with probability  $p(s)$ . The constants  $b_{si}$ 's satisfy the unbiasedness condition  $\sum_{s \supset i} b_{si} p(s) = 1$ . Then we can write Eq. (9.3.6) as

$$\begin{aligned}
 M(\hat{Y}_{reg}) &\cong E \left[ \sum_{i \in s} b_{si} (y_i - \beta x_i) - (Y - \beta X) \right]^2 \\
 &= V \left[ \sum_{i \in s} b_{si} (y_i - \beta x_i) \right] \\
 &= V \left( \sum_{i \in s} b_{si} E_i \right) \\
 &= \sum_i (\alpha_i - 1) E_i^2 + \sum_{i \neq j} \sum_j (\alpha_{ij} - 1) E_i E_j \\
 &= V(\hat{E})
 \end{aligned} \tag{9.3.8}$$

where  $E_i = y_i - \beta x_i$ ,  $\hat{E} = \sum_{i \in s} b_{si} E_i$ ,  $\alpha_i = \sum_{s \supset i} b_{si}^2 p(s)$ , and  $\alpha_{ij} = \sum_{s \supset i, j} b_{si} b_{sj} p(s)$ .

From Eq. (9.3.8), we set an approximate unbiased estimator for  $M(\hat{Y}_{reg})$  as

$$\begin{aligned}
 \hat{M}(\hat{Y}_{reg}) &= \sum_{i \in s} (\alpha_i - 1) \hat{E}_i^2 / \pi_i + \sum_{i \neq j} \sum_j (\alpha_{ij} - 1) \hat{E}_i \hat{E}_j / \pi_{ij} \\
 &= \hat{V}(\hat{E})
 \end{aligned} \tag{9.3.9}$$

where  $\hat{E}_i = y_i - \hat{\beta} x_i$ ,  $\pi_i$  = inclusion probability of  $i$ th unit,  $\pi_{ij}$  ( $> 0$ ) inclusion probability of the  $i$ th and  $j$  ( $\neq i$ )th unit and  $\hat{\beta}$  as in Eq. (9.3.1).

### 9.3.4 Mean Square Errors for Some Sampling Designs

#### 9.3.4.1 Arbitrary Fixed Effective Sample Size Design

Consider a fixed effective sample size  $n$  design with  $b_{si} = 1/\pi_i$ . In this case

$$\begin{aligned}
 \hat{X} &= \sum_{i \in s} x_i / \pi_i, \hat{Y} = \sum_{i \in s} y_i / \pi_i, V(\hat{X}) = \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2, \\
 Cov(\hat{X}, \hat{Y}) &= \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right) \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right), \\
 \beta &= \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right) \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right) / \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2
 \end{aligned}$$

$$\begin{aligned}\widehat{V}(\widehat{X}) &= \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2, \\ C \widehat{ov}(\widehat{X}, \widehat{Y}) &= \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right) \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right) \\ \widehat{\beta} &= \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right) \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right) / \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2\end{aligned}$$

where  $\Delta_{ij} = \pi_i \pi_j - \pi_{ij}$ .

The expression for the regression estimator  $\widehat{Y}_{reg}$ , its approximate MSE, and an approximate unbiased estimator of the MSE, respectively, come out as follows:

$$\widehat{Y}_{reg} = \sum_{i \in s} \frac{y_i}{\pi_i} - \widehat{\beta} \left( \sum_{i \in s} \frac{x_i}{\pi_i} - X \right) \quad (9.3.10)$$

$$M(\widehat{Y}_{reg}) = \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} \left( \frac{E_i}{\pi_i} - \frac{E_j}{\pi_j} \right)^2 \quad (9.3.11)$$

$$\widehat{M}(\widehat{Y}_{reg}) = \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( \frac{\widehat{E}_i}{\pi_i} - \frac{\widehat{E}_j}{\pi_j} \right)^2 \quad (9.3.12)$$

#### 9.3.4.2 Simple Random Sampling Without Replacement

For an SRSWOR design  $\pi_i = \frac{n}{N}$  and  $\pi_{ij} = \frac{n(n-1)}{N(N-1)}$ . In this case

$$\begin{aligned}\widehat{X} &= N \bar{x}_s, \widehat{Y} = N \bar{y}_s, V(\widehat{X}) = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) S_x^2, \widehat{V}(\widehat{X}) = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) s_x^2, \\ V(\widehat{Y}) &= N^2 \left( \frac{1}{n} - \frac{1}{N} \right) S_y^2, \widehat{V}(\widehat{Y}) = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) s_y^2, \text{Cov}(\widehat{X}, \widehat{Y}) \\ &= N^2 \left( \frac{1}{n} - \frac{1}{N} \right) S_{xy}, C \widehat{ov}(\widehat{X}, \widehat{Y}) = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) s_{xy}, \beta = \frac{S_{xy}}{S_x^2}, \widehat{\beta} = \frac{s_{xy}}{s_x^2}, \text{ and} \\ \rho_{\widehat{X}, \widehat{Y}} &= \rho = \frac{S_{xy}}{S_x S_y}, \quad \text{where} \quad \bar{x}_s = \sum_{i \in s} x_i / n, \quad \bar{y}_s = \sum_{i \in s} y_i / n, \\ S_x^2 &= \sum_{i \in U} (x_i - \bar{X})^2 / (N-1), \quad s_x^2 = \sum_{i \in s} (x_i - \bar{x}_s)^2 / (n-1), \quad \text{and} \\ S_y^2 &= \sum_{i \in U} (y_i - \bar{Y})^2 / (N-1), \quad s_y^2 = \sum_{i \in s} (y_i - \bar{y}_s)^2 / (n-1).\end{aligned}$$

So, under SRSWOR

$$\widehat{Y}_{reg} = N \left[ \bar{y}_s - \widehat{\beta} (\bar{x}_s - \bar{X}) \right] \quad (9.3.13)$$

$$M(\widehat{Y}_{reg}) = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) S_E^2 = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) (1 - \rho^2) S_y^2 \quad (9.3.14)$$

$$\widehat{M}(\widehat{Y}_{reg}) = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) s_E^2 = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) (1 - \widehat{\rho}^2) s_Y^2 \quad (9.3.15)$$

where

$$s_E^2 = \sum_{i \in U} (E_i - \bar{E})^2 / (N - 1), \quad s_E^2 = \sum_{i \in U} (\widehat{E}_i - \widehat{\bar{E}}_s)^2 / (n - 1),$$

$$\bar{E} = \sum_{i \in U} E_i / N, \quad \widehat{\bar{E}}_s = \sum_{i \in s} \widehat{E}_i / n \text{ and } \widehat{\rho} = s_{xy} / (s_x s_y).$$

### 9.3.5 Efficiency of Regression Estimator

From Eq. (9.3.7) we note that the regression estimator  $\widehat{Y}_{reg}$  is more efficient than the conventional estimator  $\widehat{Y}$  because  $V(\widehat{Y}) - M(\widehat{Y}_{reg}) = \rho_{\widehat{X}\widehat{Y}}^2 V(\widehat{Y}) \geq 0$ .

#### 9.3.5.1 Comparison With the Ratio Estimator

##### Theorem 9.3.1

The regression estimator is more efficient than the ratio estimator.

##### Proof

The difference between the approximate expressions of the MSEs of the ratio and regression estimators obtained from the Theorem 8.3.1 and Eq. (9.3.7) is

$$\begin{aligned} M(\widehat{Y}_R) - M(\widehat{Y}_{reg}) &= V(\widehat{Y}) + R^2 V(\widehat{X}) - 2\rho_{\widehat{X}\widehat{Y}} R \sqrt{V(\widehat{Y}) V(\widehat{X})} \\ &\quad - (1 - \rho_{\widehat{X}\widehat{Y}}^2) V(\widehat{Y}) \\ &= \left[ R \sqrt{V(\widehat{X})} - \rho_{\widehat{X}\widehat{Y}} \sqrt{V(\widehat{Y})} \right]^2 \\ &\geq 0 \end{aligned}$$

### 9.3.6 Optimality of the Regression Estimator

It is pointed out in Corollary 6.4.3 that the regression estimator  $\widehat{Y}_{reg} = N\{\bar{y}_s - \widehat{\beta}(\bar{x}_s - \bar{X})\}$  is the optimal in the class of linear  $\xi$ -unbiased predictors  $C_{I\xi}$  of the population  $Y$  under the superpopulation model  $E_\xi(y_i) = \beta_0 + \beta x_i$ ;  $V_\xi(y_i) = \sigma^2$  and  $C_\xi(y_i, y_j) = 0$  for  $i \neq j$ .

##### Example 9.3.1

From the data given in Example 8.5.1, estimate total production of apples ( $Y$ ) by the regression method of estimation. Estimate the standard error of



the regression estimator. Compare the relative efficiencies of the regression estimator with the ratio estimator and the expansion estimator for estimating the total  $Y$ .

From the given data, we get  $N = 125$ ,  $n = 15$ ,  $\bar{X} = 40$ ,  $\bar{x}(s) = 35$ ,  $\bar{y}(s) = 1050$ ,  $s_x^2 = 164.286$ ,  $s_y^2 = 133942.9$ , and  $s_{xy} = 4617.857$ . The sample regression coefficient  $\hat{\beta} = \frac{s_{xy}}{s_x^2} = 28.1087$  and  $\hat{\rho}$  = sample correlation coefficient  $= r = 0.984$ .

The regression estimator for the population total  $Y$  is

$$\begin{aligned}\hat{Y}_{reg} &= N \left[ \bar{y}_s - \hat{\beta}(\bar{x}_s - \bar{X}) \right] = 125 \times \{1050 - 28.1087 \times (35 - 40)\} \\ &= 148,817.9\end{aligned}$$

An estimate of the MSE of the regression estimator is

$$\begin{aligned}\hat{M}(\hat{Y}_{reg}) &= N^2(1-f)(1-\hat{\rho}^2)s_y^2/n \\ &= (125)^2(1-15/125)(1-0.9844^2) \times 13394.9/15 \\ &= 3975,839.8\end{aligned}$$

Estimated standard error of  $\hat{Y}_{reg}$  is  $\widehat{SE}(\hat{Y}_{reg}) = \sqrt{\hat{M}(\hat{Y}_{reg})} = 1948.29$ .

From the Example 8.5.1, we get  $\hat{Y}_R = N\hat{R}\bar{X} = 125 \times 30 \times 40 = 150,000$  and the estimated MSE is  $\hat{M}(\hat{Y}_R) = N^2\hat{M}(\hat{Y}_R) = 125^2 \times 277.389 = 4334203.13$ .

An estimate of  $Y$  based on the expansion estimator is  $\hat{Y} = N\bar{y}_s = 125 \times 1050 = 131,250$ .

Estimated variance of  $\hat{Y}$  is  $\hat{V}(\hat{Y}) = N^2(1-f)s_y^2/n = 125^2(1-15/125) \times 133,942.9/15 = 122,780,952.4$ .

Estimated efficiency of the regression estimator compared to the ratio estimator is

$$\left\{ \hat{M}(\hat{Y}_R) / \hat{M}(\hat{Y}_{reg}) \right\} \times 100 = 114.18\%$$

Estimated efficiency of the regression estimator compared to the expansion estimator is

$$\left\{ \hat{V}(\hat{Y}) / \hat{M}(\hat{Y}_{reg}) \right\} \times 100 = 3234.61\%$$

### 9.3.7 Unbiased Regression Estimator

Singh and Srivastava (1980) proposed a sampling scheme for which the regression estimator (Eq. 9.3.13) is unbiased for the total  $Y$ . The sampling scheme is described as follows.

#### 9.3.7.1 Singh and Srivastava Sampling Scheme

Under the Singh and Srivastava (SS) sampling scheme, we first select two units  $i$  and  $j$  ( $i \neq j$ ) with probability proportional to  $(x_i - x_j)^2$ , then select  $(n - 2)$  units from the remaining  $(N - 2)$  units by SRSWOR method. So, for the SS sampling scheme, probability of selection a sample  $s$  of  $n$  units is

$$p(s) \propto \sum_{i \neq j} \sum_{j \in s} (x_i - x_j)^2 \quad (9.3.16)$$

$$\text{i.e. } p(s) = \sum_{i \neq j} \sum_{j \in s} (x_i - x_j)^2 / \sum_s \sum_{i \neq j} \sum_{j \in s} (x_i - x_j)^2$$

Now noting  $\sum_{i \neq j} \sum_{j \in s} (x_i - x_j)^2 = 2n(n - 1)s_x^2$  and  $\sum_s \sum_{i \neq j} \sum_{j \in s} (x_i - x_j)^2 = 2N(N - 1)S_x^2 \binom{N - 2}{n - 2}$ , we find

$$p(s) = s_x^2 / \left\{ \binom{N}{n} S_x^2 \right\} \quad (9.3.17)$$

Now noting  $\hat{\beta} = (s_{xy}/s_x^2)$ , we find for the SS sampling scheme

$$\begin{aligned} E(\hat{Y}_{reg}) &= N \sum_s \{ \bar{y}_s - (s_{xy}/s_x^2)(\bar{x}_s - \bar{X}) \} s_x^2 / \left\{ s_x^2 \binom{N}{n} \right\} \\ &= \frac{N}{S_x^2} \sum_s \{ \bar{Y} s_x^2 + (\bar{y}_s - \bar{Y}) s_x^2 - (s_{xy})(\bar{x}_s - \bar{X}) \} / \binom{N}{n} \end{aligned}$$

$$\begin{aligned} &= N [ \bar{Y} E(s_x^2 | \text{SRSWOR}) + \text{Cov}(\bar{y}_s, s_x^2 | \text{SRSWOR}) \\ &\quad - \text{Cov}(s_{xy}, \bar{x}_s | \text{SRSWOR}) ] / S_x^2 \end{aligned}$$

where

$$E(s_x^2 | \text{SRSWOR}) = S_x^2, \quad (9.3.18)$$

$$\text{Cov}(\bar{y}_s, s_x^2 | \text{SRSSWOR}) = \frac{1}{n} \frac{N(N-n)}{(N-1)(N-2)} \mu_{21} \quad (9.3.19)$$

$$\text{Cov}(s_{xy}, \bar{x}_s | \text{SRSSWOR}) = \frac{1}{n} \frac{N(N-n)}{(N-1)(N-2)} \mu_{21} \quad (9.3.20)$$

$$\text{and } \mu_{21} = \frac{1}{N} \sum_{i \in U} (x_i - \bar{X})^2 (y_i - \bar{Y}).$$

The results (Eqs. 9.3.19 and 9.3.20) were obtained by Sukhatme (1944) (details given in Appendix 9A). From Eqs. (9.3.18–9.3.20), we note that  $\hat{Y}_{reg}$  is an unbiased estimator of the total  $Y$ .

The expression of the exact variance of  $\hat{Y}_{reg}$  for this sampling scheme is very complex. However, an approximate expression of variance and its approximate unbiased estimator are stated in the following theorem without derivation.

### Theorem 9.3.2

Under the SS sampling scheme

(i)  $\hat{Y}_{reg} = N\{\bar{y}_s - (s_{xy}/s_x^2)(\bar{x}_s - \bar{X})\}$  is unbiased for the  $Y$

$$(ii) \quad V(\hat{Y}_{reg}) \cong N^2 \left[ (1 - \rho^2) \left( \frac{1}{n} + \frac{1}{n^2} \right) \mu_{02} + \frac{1}{n^2 \mu_{20}} \left( \frac{2\mu_{11}\mu_{31}}{\mu_{20}} - \frac{\mu_{11}^2 \mu_{40}}{\mu_{20}^2} - \mu_{22} \right) \right]$$

$$(iii) \quad \hat{V}(\hat{Y}_{reg}) \cong Y_{reg}^2 - \frac{NS_x^2}{ns_x^2} \left( \sum_{i \in s} y_i^2 + \frac{N-1}{n-1} \sum_{i \neq j} \sum_{i \in s} y_i y_j \right)$$

$$\text{where } \mu_{jk} = \frac{1}{N} \sum_{i \in U} (x_i - \bar{X})^j (y_i - \bar{Y})^k$$

## 9.3.8 Stratified Regression Estimator

Here we suppose that a population is stratified into  $K$  strata. Let  $y_{ij}$  and  $x_{ij}$  be the value of the study ( $y$ ) and auxiliary ( $x$ ) variables for the  $j$ th unit of the  $i$ th stratum  $U_i$ ,  $i = 1, \dots, K$ . From each of the strata, samples are selected independently. Let the sample  $s_i$  of size  $n_i$  be selected from the  $i$ th stratum of size  $N_i$  by some sampling procedure. Let  $\hat{Y}_i$  and  $\hat{X}_i$  be the unbiased estimators of the totals  $Y_i$  and  $X_i$  of the study and auxiliary variables of the  $i$ th stratum, respectively. Then we may construct the following two types of regression estimators, viz. separate and combined regression estimators for the total  $Y$  analogous to separate and combined ratio estimators.

### 9.3.8.1 Separate Regression Estimator

Here we assume that the linear regressions of  $y$  and  $x$  for the different strata are different and that the population totals of the auxiliary information  $X_i$ 's

of all the strata are known. Under this situation, we construct a regression estimator of the total  $Y_i$  as

$$\hat{Y}_{i(\text{reg})} = \hat{Y}_i - \hat{\beta}_i(\hat{X}_i - X_i) \quad (9.3.21)$$

where  $\hat{\beta}_i = C \hat{\sigma}_v(\hat{Y}_i, \hat{X}_i) / \hat{V}(\hat{X}_i)$  is an estimate of the population regression coefficient  $\beta_i = \text{Cov}(\hat{Y}_i, \hat{X}_i) / V(\hat{X}_i)$  of the  $i$ th stratum.

The separate regression estimators for the population total  $Y$  is given by

$$\hat{Y}_{\text{reg}}(s) = \sum_{i=1}^K \hat{Y}_{i(\text{reg})} \quad (9.3.22)$$

The bias of  $\hat{Y}_{\text{reg}}(s)$  is

$$B(s) = \sum_{i=1}^K B(\hat{Y}_{i(\text{reg})})$$

where  $B(\hat{Y}_{i(\text{reg})})$  is the bias of  $\hat{Y}_{i(\text{reg})}$ .

Eqs. (9.3.3) and (9.3.5) yield

$$\begin{aligned} B(s) &= - \sum_{i=1}^K \text{Cov}(\hat{\beta}_i, \hat{X}_i) \\ &\equiv \sum_{i=1}^K \beta_i \left[ \frac{\text{Cov}\{\hat{V}(\hat{X}_i), \hat{X}_i\}}{V(\hat{X}_i)} - \frac{\text{Cov}\{\hat{Cov}(\hat{X}_i, \hat{Y}_i), \hat{X}_i\}}{\text{Cov}(\hat{X}_i, \hat{Y}_i)} \right] \end{aligned} \quad (9.3.23)$$

An approximate expression for the MSE of  $\hat{Y}_{\text{reg}}(s)$  is obtained from Eqs. (9.3.6) and (9.3.7) as

$$M[\hat{Y}_{\text{reg}}(s)] \cong \sum_{i=1}^K V\{(\hat{Y}_i) - 2\beta_i \text{Cov}(\hat{X}_i, \hat{Y}_i) + \beta_i^2 V(X_i)\} \quad (9.3.24)$$

$$= \sum_{i=1}^K (1 - \rho_{\hat{X}_i, \hat{Y}_i}^2) V(\hat{Y}_i) \quad (9.3.25)$$

where  $\rho_{\hat{X}_i, \hat{Y}_i}$  is the correlation between  $\hat{X}_i$  and  $\hat{Y}_i$ .

In case samples are selected from each of the strata by SRSWOR, we have

$$\hat{Y}_{\text{reg}}(s) = \sum_{i=1}^K N_i \left\{ \bar{y}_i - \hat{\beta}_i(\bar{x}_i - \bar{X}_i) \right\} \quad (9.3.26)$$

$$B(s) = \sum_{i=1}^K \frac{N_i(N_i - n_i)}{(N_i - 2)} \frac{\beta_i}{n_i} \left( \frac{\mu_{30}^{(i)}}{\mu_{20}^{(i)}} - \frac{\mu_{21}^{(i)}}{\mu_{11}^{(i)}} \right) \quad (9.3.27)$$

$$\begin{aligned}
M[\hat{Y}_{reg}(s)] &\cong \sum_{i=1}^K N_i^2(1-f_i) \left( S_{yi}^2 - 2\beta_i S_{xyi} + \beta_i^2 S_{xi}^2 \right) / n_i \\
&= \sum_{i=1}^K N_i^2(1-f_i) S_{qi}^2 / n_i \\
&= \sum_{i=1}^K N_i^2(1-f_i) (1-\rho_i^2) S_{yi}^2 / n_i
\end{aligned} \tag{9.3.28}$$

$$\begin{aligned}
\hat{M}[\hat{Y}_{reg}(s)] &\cong \sum_{i=1}^K N_i^2(1-f_i) s_{\hat{q}i}^2 / n_i \\
&= \sum_{i=1}^K N_i^2(1-f_i) (1-\hat{\rho}_i^2) s_{yi}^2 / n_i
\end{aligned} \tag{9.3.29}$$

where  $f_i = n_i/N_i$ ,  $\bar{y}_i = \sum_{j \in s_i} y_{ij}/n_i$ ,  $\bar{x}_i = \sum_{j \in s_i} x_{ij}/n_i$ ,

$$\hat{\beta}_i = \sum_{j \in s_i} (x_{ij} - \bar{x}_i)(y_{ij} - \bar{y}_i) / \sum_{j \in s_i} (x_{ij} - \bar{x}_i)^2, S_{yi}^2 = \sum_{j \in U_i} (y_{ij} - \bar{Y}_i)^2 / (N_i - 1),$$

$$S_{xi}^2 = \sum_{j \in U_i} (x_{ij} - \bar{X}_i)^2 / (N_i - 1), S_{xyi} = \sum_{j \in U_i} (y_{ij} - \bar{Y}_i)(x_{ij} - \bar{X}_i) / (N_i - 1),$$

$$\bar{X}_i = \sum_{j \in U_i} x_{ij} / N_i, \bar{Y}_i = \sum_{j \in U_i} y_{ij} / N_i, \beta_i = S_{xyi} / S_{xi}, \rho_i = S_{xyi} / (S_{xi} S_{yi}),$$

$$s_{xi}^2 = \sum_{j \in s_i} (x_{ij} - \bar{x}_i)^2 / (n_i - 1), s_{yi}^2 = \sum_{j \in s_i} (y_{ij} - \bar{y}_i)^2 / (n_i - 1),$$

$$s_{xyi} = \sum_{j \in s_i} (x_{ij} - \bar{x}_i)(y_{ij} - \bar{y}_i) / (n_i - 1), \bar{x}_i = \sum_{j \in s_i} x_{ij} / n_i, \bar{y}_i = \sum_{j \in s_i} y_{ij} / n_i,$$

$$\hat{\rho}_i = s_{xyi} / (s_{xi} s_{yi}), S_{qi}^2 = \sum_{j \in U_i} (q_{ij} - \bar{Q}_i)^2 / (N_i - 1), q_{ij} = y_{ij} - \beta_i x_{ij},$$

$$\bar{Q}_i = \sum_{j \in U_i} q_{ij} / N_i, s_{\hat{q}i}^2 = \sum_{j \in s_i} (\hat{q}_{ij} - \bar{q}_i)^2 / (n_i - 1), \hat{q}_{ij} = y_{ij} - \hat{\beta}_i x_{ij},$$

$$\bar{q}_i = \sum_{j \in s_i} \hat{q}_{ij} / n_i, \text{ and } \mu_{kl}^{(i)} = \sum_{j \in U_i} (x_{ij} - \bar{X}_i)^k (y_{ij} - \bar{Y}_i)^l / N_i.$$

### 9.3.8.2 Combined Regression Estimator

Here we assume that a single regression line may be fitted over the entire population. The combined regression estimator for the total  $Y$  is given by

$$\hat{Y}_{reg}(c) = \hat{Y}_{st} - \hat{\beta}_c (\hat{X}_{st} - X) \tag{9.3.30}$$

where  $\hat{X}_{st} = \sum_{i=1}^K \hat{X}_i$ ,  $\hat{Y}_{st} = \sum_{i=1}^K \hat{Y}_i$  and  $\hat{\beta}_c = C\hat{o}v(\hat{X}_{st}, \hat{Y}_{st}) / \hat{V}(\hat{X}_{st}) =$  an estimate of the combined regression coefficient  $\beta_c = Cov(\hat{X}_{st}, \hat{Y}_{st}) / V(\hat{X}_{st})$ .

The bias of  $\hat{Y}_{reg}(c)$  is obtained by using Eqs. (9.3.3) and (9.3.5) as

$$B(c) = -Cov\left(\hat{\beta}_c, \hat{X}_{st}\right) \\ \cong \beta_c \left[ \frac{Cov\left\{\hat{V}(\hat{X}_{st}), \hat{X}_{st}\right\}}{V(\hat{X}_{st})} - \frac{Cov\left\{\hat{Cov}(\hat{X}_{st}, \hat{Y}_{st}), \hat{X}_{st}\right\}}{Cov(\hat{X}_{st}, \hat{Y}_{st})} \right] \quad (9.3.31)$$

An approximate MSE of  $\hat{Y}_{reg}(c)$  upto first order of approximation is obtained from Eq. (9.3.6) as

$$M[\hat{Y}_{reg}(c)] \cong V(\hat{Y}_{st} - \beta_c \hat{X}_{st}) \\ = \sum_{i=1}^K \{V(\hat{Y}_i) - 2\beta_c \text{cov}(\hat{X}_i, \hat{Y}_i) + \beta_c^2 V(X_i)\} \\ = \sum_{i=1}^K V(\hat{Q}_{ci}) \quad (9.3.32)$$

where  $\hat{Q}_{ci} = \hat{Y}_i - \beta_c \hat{X}_i$ .

An approximate unbiased estimator of  $M[\hat{Y}_{reg}(c)]$  is

$$\hat{M}[\hat{Y}_{reg}(c)] = \sum_{i=1}^K \left\{ \hat{V}(\hat{Y}_i) - 2\hat{\beta}_c \hat{Cov}(\hat{X}_i, \hat{Y}_i) + \hat{\beta}_c^2 \hat{V}(X_i) \right\} \quad (9.3.33)$$

For SRSWOR sampling,  $\hat{Y}_{reg}(c)$  reduces to

$$\hat{Y}_{reg}(c) = \left( \sum_{i=1}^K N_i \bar{y}_i \right) - \hat{\beta}_c \left( \sum_{i=1}^K N_i (\bar{x}_i - \bar{X}) \right) \quad (9.3.34)$$

where  $\hat{\beta}_c = \frac{\sum_{i=1}^K N_i^2 (1 - f_i) s_{xyi} / n_i}{\sum_{i=1}^K N_i^2 (1 - f_i) s_{xi}^2 / n_i}$ .

The expressions for MSE and its approximate unbiased estimator of  $\hat{Y}_{reg}(c)$  under SRSWOR are, respectively, given by

$$M[\hat{Y}_{reg}(c)] \cong \sum_{i=1}^K N_i^2 (1 - f_i) \left\{ S_{yi}^2 - 2\beta_c S_{xyi} + \beta_c^2 S_{xi}^2 \right\} / n_i \\ = \sum_{i=1}^K N_i^2 (1 - f_i) S_{gi}^2 / n_i \quad (9.3.35)$$

and

$$\hat{M}[\hat{Y}_{reg}(c)] \cong \sum_{i=1}^K N_i^2 (1 - f_i) \left\{ s_{yi}^2 - 2\hat{\beta}_c s_{xyi} + \hat{\beta}_c^2 s_{xi}^2 \right\} / n_i \\ = \sum_{i=1}^K N_i^2 (1 - f_i) s_{gi}^2 / n_i \quad (9.3.36)$$

where  $S_{gi}^2 = \sum_{j \in U_i} (g_{ij} - \bar{G}_i)^2 / (N_i - 1)$ ,  $g_{ij} = y_{ij} - \beta_c x_{ij}$ ,  $\bar{G}_i = \sum_{j \in U_i} g_{ij} / N_i$ ,  $s_{\hat{g}_i}^2 = \sum_{j \in s_i} (\hat{g}_{ij} - \bar{g}_i)^2 / (n_i - 1)$ ,  $\hat{g}_{ij} = y_{ij} - \hat{\beta}_c x_{ij}$ , and  $\bar{g}_i = \sum_{j \in s_i} \hat{g}_{ij} / n_i$ .

### 9.3.8.3 Comparison Between Combined and Separate Regression Estimators

Approximate expressions of biases and MSEs for separate and combined regression estimators are obtained assuming sample sizes are large for each of the strata. The bias of a regression estimator decreases with the increase in sample size. The total sample size  $n$  may be large, but the sample sizes  $n_i$ 's for each of the strata may not be large. Hence the bias of  $\hat{Y}_{i(reg)}$  may not be negligible compared to its standard error if the sample size  $n_i$  is small. Hence the amount of bias of the separate regression estimator is expected to be much larger than the combined regression estimator.

It is obvious that the separate regression estimator  $\hat{Y}_{reg}(s)$  is more efficient than the combined regression estimator  $\hat{Y}_{reg}(c)$  because

$$M[\hat{Y}_{reg}(c)] - M[\hat{Y}_{reg}(s)] \cong \sum_{i=1}^K \{V(\hat{Y}_i - \beta_c \hat{X}_i) - V(\hat{Y}_i - \beta_i \hat{X}_i)\} \geq 0$$

as  $V(\hat{Y}_i - \beta_c \hat{X}_i)$  attains a minimum when  $\beta_c = \beta_i$ .

In practice, the values of  $\beta_i$ 's are unknown and they are estimated from the samples. In case the sample sizes are not large enough for each of the strata, the estimate  $\hat{\beta}_i$  becomes unreliable and hence the separate regression estimator becomes inefficient. But, if the regression coefficients are approximately equal, then  $\hat{\beta}_c$ , the estimator of  $\beta_c$  becomes more efficient than each of the separate estimators  $\hat{\beta}_i$  because the former is computed using a relatively large sample size. In this situation combined regression estimator certainly becomes more efficient than the separate regression estimators. As a rule of thumb, a combined regression estimator should be used in case regression lines for each of the strata are the same. Otherwise, separate regression estimator should be used.

#### Example 9.3.2

From the data given in Example 8.7.1, estimate the average wage of the factory workers by using (i) separate and (ii) combined regression estimators when average years of services of the three categories are known to be 15, 8, and 3 years, respectively. Estimate the relative efficiency of the separate regression estimator with respect to the combined regression estimator.

Here we prepare the following tables:

Strata	$N_i$	$\bar{X}_i$	$N_i \bar{X}_i$	$n_i$	$\bar{x}(s_i)$	$\bar{y}(s_i)$	$N_i \bar{x}(s_i)$	$N_i \bar{y}(s_i)$	$s_{x_i}^2$	$s_{y_i}^2$	$s_{xy_i}$	$\hat{\beta}_i$
Skilled	50	15	750	10	16.3	610	815	30,500	59.122	13,222.220	824.444	13.944
Semiskilled	75	8	600	12	9	425	675	31,875	6.727	6,590.909	150.000	22.297
Unskilled	100	3	300	15	3.8	260	380	26,000	2.171	12,571.430	134.285	61.842
Total	225		1650				1,870	88,375				
$\hat{Y}_{i(reg)}$	$\hat{\rho}_i = r_i$		$\frac{N_i^2(1-f_i)(1-r_i^2)s_{yi}^2}{n_i}$				$\frac{N_i^2(1-f_i)s_{xyi}}{n_i}$		$\frac{N_i^2(1-f_i)s_{xi}^2}{n_i}$		$\frac{N_i^2(1-f_i)\{s_{yi}^2 - 2\hat{\beta}_c s_{xyi} + \hat{\beta}_c^2 s_{xi}^2\}}{n_i}$	
29,593.591	0.932		345,110.568				164,888.889		11,824.440		660,175.087	
30,202.702	0.712		1,278,236.333				59,062.500		2,648.864		1,305,202.464	
21,052.631	0.813		2,417,919.799				76,095.238		1,230.476		4,665,162.568	
80,848.925			4,041,266.701				300,046.627		15,703.780		6,630,540.119	



Separate regression estimate for the average wage  $\bar{Y}$  is

$$\hat{\bar{Y}}_{reg}(s) = \hat{Y}_{reg}(s)/N = \sum_{i=1}^K \hat{Y}_{i(reg)}/N = 80848.925/225 = 359.329$$

$$\begin{aligned} \text{Estimated MSE of } \bar{Y}_{reg}(s) &= \hat{M}\{\hat{\bar{Y}}_{reg}(s)\}/N^2 \\ &= \left\{ \sum_{i=1}^K N_i^2(1-f_i)(1-\hat{\rho}_i^2)s_{yi}^2/n_i \right\} / N^2 \\ &= 4041266.701/(225)^2 = 79.827 \end{aligned}$$

An estimate of the combined regression  $\beta_c$  is

$$\hat{\beta}_c = \frac{\sum_{i=1}^3 N_i^2(1-f_i)s_{xyi}/n_i}{\sum_{i=1}^3 N_i^2(1-f_i)s_{xi}^2/n_i} = 300046.627/15703.780 = 19.107$$

Combined regression estimate of the average wage  $\bar{Y}$  is

$$\begin{aligned} \hat{\bar{Y}}_{reg}(c) &= \hat{Y}_{reg}(c)/N = \left( \sum_{i=1}^K N_i \bar{y}_i / N \right) - \hat{\beta}_c \left( \sum_{i=1}^K N_i \bar{x}_i / N - \bar{X} \right) \\ &= 88375/225 - 19.107 \times (1870/225 - 1650/225) = 374.096 \end{aligned}$$

Estimated MSE of  $\hat{\bar{M}}[\hat{\bar{Y}}_{reg}(c)]$  is

$$\begin{aligned} \hat{M}[\hat{\bar{Y}}_{reg}(c)] &= \left[ \sum_{i=1}^K N_i^2(1-f_i)(s_{yi}^2 - 2\hat{\beta}_c s_{xyi} + \hat{\beta}_c^2 s_{xi}^2)/n_i \right] / N^2 \\ &= 6630540.119/(225)^2 = 130.974 \end{aligned}$$

Estimated relative efficiency of separate regression estimator with respect to the combined regression estimator for estimating  $\bar{Y} = \hat{\bar{M}}[\hat{\bar{Y}}_{reg}(c)] \times 100 / \hat{M}[\hat{\bar{Y}}_{reg}(s)] = (130.974/79.827) \times 100 = 164.071$  percent.

### 9.3.9 Regression Estimator for Several Auxiliary Variables

Suppose that we have  $p$  auxiliary variables  $x_1, \dots, x_p$ , each of which is well related to the study variable  $y$ . Let  $x_{ij}$  be the value of  $i$ th unit of the  $j$ th auxiliary variable  $x_j$ . In this case, we may set the following two types of regression estimators for the total  $Y$  when the totals of each of the auxiliary variables are known.

### 9.3.9.1 Multivariate Regression Estimator

$$\hat{Y}_{reg}(1) = \hat{Y} - \lambda_1(\hat{X}_1 - X_1) - \cdots - \lambda_j(\hat{X}_j - X_j) - \cdots - \lambda_p(\hat{X}_p - X_p) \quad (9.3.37)$$

where  $\hat{X}_j = \sum_{i \in S} b_{si} x_{ji}$  is an unbiased estimator of the total  $X_j$  and  $\lambda_1, \dots, \lambda_p$  are known constants. The optimum values of  $\lambda_1, \dots, \lambda_p$  are obtained by minimizing the variance of  $\hat{Y}_{reg}(1)$ . The variance of  $\hat{Y}_{reg}(1)$  is

$$\begin{aligned} V\{\hat{Y}_{reg}(1)\} = & V(\hat{Y}) + \sum_{i=1}^p \lambda_i^2 V(\hat{X}_i) + \sum_{i \neq j}^p \sum_{j=p}^p \lambda_i \lambda_j Cov(\hat{X}_i, \hat{X}_j) \\ & - 2 \sum_{i=1}^p \lambda_i Cov(\hat{Y}, \hat{X}_i) \end{aligned} \quad (9.3.38)$$

Differentiating  $V\{\hat{Y}_{reg}(1)\}$  with respect to  $\lambda_j$  and equating to zero for  $j = 1, \dots, p$ , the optimum values of  $\lambda_j = \beta_j$  (say) are obtained as

$$\begin{aligned} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_j \\ \vdots \\ \beta_p \end{pmatrix} = & \begin{pmatrix} V(\hat{X}_1) & \dots & Cov(\hat{X}_1, \hat{X}_j) & \dots & Cov(\hat{X}_1, \hat{X}_p) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ Cov(\hat{X}_j, \hat{X}_1) & \dots & V(\hat{X}_j) & \dots & Cov(\hat{X}_j, \hat{X}_p) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ Cov(\hat{X}_p, \hat{X}_1) & \dots & Cov(\hat{X}_p, \hat{X}_j) & \dots & V(\hat{X}_p) \end{pmatrix}^{-1} \\ & \times \begin{pmatrix} Cov(\hat{Y}, \hat{X}_1) \\ \vdots \\ Cov(\hat{Y}, \hat{X}_j) \\ \vdots \\ Cov(\hat{Y}, \hat{X}_p) \end{pmatrix} \end{aligned} \quad (9.3.39)$$

The solution (Eq. 9.3.39) indicates that the optimum value  $\beta_j$  is the partial regression coefficient of  $\hat{Y}$  on  $\hat{X}_j$ .

In particular for SRSWOR sampling design  $V(\hat{X}_j) = N^2(1-f)S_{x_j}^2/n$ ,  $Cov(\hat{X}_j, \hat{X}_k) = N^2(1-f)S_{x_j, x_k}^2/n$  and  $Cov(\hat{Y}, \hat{X}_k) = N^2(1-f)S_{y, x_k}^2/n$  where  $S_{x_j}^2 = \sum_{i \in U} (x_{ji} - \bar{x}_j)^2 / (N-1)$ ,

$$S_{x_j, x_k}^2 = \sum_{i \in U} (x_{ji} - \bar{x}_j)(x_{ki} - \bar{x}_k) / (N - 1),$$

$$S_{y, x_k} = \sum_{i \in U} (y_i - \bar{y})(x_{ki} - \bar{x}_k) / (N - 1), \text{ and } \bar{x}_j = \sum_{i \in U} x_{ji} / N.$$

The value of  $\hat{Y}_{reg}(1)$  with the optimum value of  $\lambda_i = \beta_i$  comes out as

$$V_{\min}(\hat{Y}_{reg}(1)) = N^2(1 - f)(1 - R_{y|x_1, \dots, x_k}^2) / n \quad (9.3.40)$$

where  $R_{y|x_1, \dots, x_k}^2$  = multiple correlation between  $y$  and  $x_1, \dots, x_p$ .

The optimum values  $\beta_j$  are generally unknown. Hence regression estimator of  $Y$  is obtained by replacing  $\beta_j$  by its suitable estimate  $\hat{\beta}_j$  and it is given as

$$\hat{Y}_{reg}(1) = \hat{Y} - \hat{\beta}_1(\hat{X}_1 - X_1) - \dots - \hat{\beta}_j(\hat{X}_j - X_j) - \dots - \hat{\beta}_p(\hat{X}_p - X_p)$$

### 9.3.9.2 Two Auxiliary Variables

For  $p = 2$ , Eq. (9.3.39) reduces to

$$\begin{aligned} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} &= \begin{pmatrix} S_{x_1}^2 & \rho_{x_1, x_2} S_{x_1} S_{x_2} \\ \rho_{x_1, x_2} S_{x_1} S_{x_2} & S_{x_2}^2 \end{pmatrix}^{-1} \begin{pmatrix} \rho_{yx_1} S_{x_1} S_y \\ \rho_{yx_2} S_{x_2} S_y \end{pmatrix} \\ &= \begin{pmatrix} \frac{\rho_{yx_1} - \rho_{x_1, x_2} \rho_{yx_2}}{1 - \rho_{x_1, x_2}^2} \frac{S_y}{S_{x_1}} \\ \frac{\rho_{yx_2} - \rho_{x_1, x_2} \rho_{yx_1}}{1 - \rho_{x_1, x_2}^2} \frac{S_y}{S_{x_2}} \end{pmatrix} \end{aligned}$$

where  $\rho_{uv}$  is the correlation coefficient between  $u$  and  $v$ .

Hence the variance of  $\hat{Y}_{reg}(1) = \hat{Y} - \lambda_1(\hat{X}_1 - X_1) - \lambda_2(\hat{X}_2 - X_2)$  attains a minimum when  $\lambda_1 = \beta_1$  and  $\lambda_2 = \beta_2$ . The corresponding minimum variance is

$$V_{\min}\{\hat{Y}_{reg}(1)\} = N^2(1 - f) \left(1 - R_{y|x_1, x_2}^2\right) S_y^2 / n$$

$$\text{where } R_{y|x_1, x_2}^2 = \frac{\rho_{yx_1}^2 + \rho_{yx_2}^2 - 2\rho_{yx_1}\rho_{yx_2}\rho_{x_1, x_2}}{1 - \rho_{x_1, x_2}^2}$$

### 9.3.9.3 Raj's Regression Estimator

Raj (1965a,b) computed  $p$  difference estimators based on each of the auxiliary variables and then obtained a combined regression estimator by taking the weighted average of these separate estimators. Raj's estimator is given by

$$\hat{Y}_{reg}(R) = \sum_{j=1}^P w_j \hat{Y}_{reg}(j) \quad (9.3.41)$$

where  $\hat{Y}_{reg}(j) = \hat{Y} - \lambda_j^* (\hat{X}_j - X_j)$ ,  $\lambda_j^*$  is a suitably chosen constant, and  $\sum_{j=1}^K w_j = 1$ .

The optimum value of  $\lambda_j^*$  that minimizes the variance of  $\hat{Y}_{reg}(j)$  is  $\beta_j^* = \rho_{yx_j} S_y / S_{x_j}$  and the corresponding minimum variance of  $\hat{Y}_{reg}(j)$  is  $N^2(1-f) \left(1 - \rho_{yx_j}^2\right) S_y^2 / n$ . The optimum values of  $w_j$ 's that minimize the variance of  $\hat{Y}_{reg}(R)$  are obtained following Section 8.8. It is obvious that the minimum variance of Eq. (9.3.41) is higher than that of Eq. (9.3.40). Hence Raj's estimator is less efficient than the estimator given in Eq. (9.3.40). For  $p = 2$ , details are given by Sukhatme et al. (1984).

## 9.4 PRODUCT METHOD OF ESTIMATION

The ratio method of estimation is more efficient than the conventional estimator  $\hat{Y}$  when the study variable  $y$  has a high positive correlation with the auxiliary variable  $x$  and  $R = Y/X$  is positive. Goodman (1960) and Murthy (1964) proposed the product estimator that is used when the study and the auxiliary variables are negatively related. The product estimator is given by

$$\hat{Y}_p = \frac{\hat{Y} \hat{X}}{X} \quad (9.4.1)$$

The product estimator  $\hat{Y}_p$  is not unbiased. The bias is negligible for a large sample size. The expressions of bias and MSE of the product estimator are given as follows.

### 9.4.1 Bias of the Product Estimator

Writing  $\varepsilon_x = \frac{\hat{X}}{X} - 1$  and  $\varepsilon_y = \frac{\hat{Y}}{Y} - 1$ , we get the bias of the product estimator as

$$\begin{aligned} B(\hat{Y}_p) &= E \left[ \frac{XY(1 + \varepsilon_x)(1 + \varepsilon_y)}{X} - Y \right] \\ &= YE(\varepsilon_x \varepsilon_y) \\ &= \frac{Cov(\hat{X}, \hat{Y})}{X} \end{aligned} \quad (9.4.2)$$

An estimator of the bias  $B(\hat{Y}_p)$  is given by

$$\hat{B}(\hat{Y}_p) = \frac{C \hat{v}(\hat{X}, \hat{Y})}{\hat{X}} \quad (9.4.3)$$

### 9.4.2 Mean Square Error of the Product Estimator

The MSE of the product estimator is given by

$$\begin{aligned} M(\hat{Y}_p) &= E \left[ \frac{XY(1 + \varepsilon_x)(1 + \varepsilon_y)}{X} - Y \right]^2 \\ &= Y^2 E(\varepsilon_x + \varepsilon_y + \varepsilon_x \varepsilon_y)^2 \end{aligned} \quad (9.4.4)$$

Neglecting terms  $E(\varepsilon_x^i \varepsilon_y^j)$  with  $i + j > 2$  in Eq. (9.4.4), we get the expression of the MSE up to the first order of approximation as

$$M(\hat{Y}_p) = V(\hat{Y}) + R^2 V(\hat{X}) + 2RCov(\hat{X}, \hat{Y}) \quad (9.4.5)$$

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

$$\hat{M}(\hat{Y}_p) = \hat{V}(\hat{Y}) + \hat{R}^2 \hat{V}(\hat{X}) + 2\hat{R}C\hat{ov}(\hat{X}, \hat{Y}) \quad (9.4.6)$$

### 9.4.3 Comparison With the Conventional Estimator

The product estimator  $\hat{Y}_p$  becomes superior to the conventional estimator  $\hat{Y}$  if

$$\begin{aligned} V(\hat{Y}) - M(\hat{Y}_p) &\geq 0 \\ \text{i.e. if } \rho_{\hat{X}, \hat{Y}} &\leq -\frac{1}{2} \frac{C_{\hat{X}}}{C_{\hat{Y}}} \quad \text{if } R > 0 \end{aligned}$$

and

$$\rho_{\hat{X}, \hat{Y}} \geq \frac{1}{2} \left| \frac{C_{\hat{X}}}{C_{\hat{Y}}} \right| \quad \text{if } R < 0$$

where  $C_{\hat{X}} = \sqrt{V(\hat{X})}/X$  and  $C_{\hat{Y}} = \sqrt{V(\hat{Y})}/Y$  are the coefficient of variation of  $\hat{X}$  and  $\hat{Y}$ , respectively.

### 9.4.4 Product Estimator for a Few Sampling Designs

#### 9.4.4.1 Fixed Effective Sample Size Design

Consider a fixed effective sample size  $n$  design and  $b_{si} = 1/\pi_i$ . In this case

$$\begin{aligned} \hat{X} &= \sum_{i \in s} x_i / \pi_i, \quad \hat{Y} = \sum_{i \in s} y_i / \pi_i, \quad \hat{Y}_p = \frac{\left( \sum_{i \in s} x_i / \pi_i \right) \left( \sum_{i \in s} y_i / \pi_i \right)}{X} \\ B(\hat{Y}_p) &= \frac{1}{X} \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right) \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right), \\ \hat{B}(\hat{Y}_p) &= \frac{1}{\hat{X}} \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( \frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right) \left( \frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right) \end{aligned}$$

$$M(\hat{Y}_p) = \frac{1}{2} \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} \left( \frac{F_i}{\pi_i} - \frac{F_j}{\pi_j} \right)^2, \quad \hat{M}(\hat{Y}_p) = \frac{1}{2} \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( \frac{\hat{F}_i}{\pi_i} - \frac{\hat{F}_j}{\pi_j} \right)^2$$

where  $F_i = \gamma_i + R x_i$  and  $\hat{F}_i = \gamma_i + \hat{R} x_i$ .

#### 9.4.4.2 Simple Random Sampling Without Replacement

For SRSWOR,  $\hat{Y}_p = \frac{N \bar{y}_s \bar{x}_s}{\bar{X}}$ ,  $B(\hat{Y}_p) = \frac{N(1-f)}{n} \frac{S_{xy}}{\bar{X}}$ ,

$\hat{B}(\hat{Y}_p) = \frac{N(1-f)}{n} \frac{s_{xy}}{\bar{x}_s}$ ,  $M(\hat{Y}_p) = N^2 \frac{(1-f)}{n} (S_y^2 + R^2 S_x^2 + 2R S_{xy})$ , and

$\hat{M}(\hat{Y}_p) = N^2 \frac{(1-f)}{n} (s_y^2 + \hat{R}^2 s_x^2 + 2\hat{R} s_{xy})$ .

### 9.4.5 Unbiased Product Type Estimators

Chaudhuri and Arnab (1982a,b) proposed the unbiased product type estimator and it is constructed as follows:

Let  $\hat{X} = \sum_{i \in s} b_{si} x_i$  and  $\hat{Y} = \sum_{i \in s} b_{si} y_i$  be the unbiased estimators of the population totals  $X$  and  $Y$  as defined in Section 9.3.3. The sample  $s$  of size  $n$  is divided at random into  $k$  groups so that the  $i$ th group  $s_i$  consists of  $n_i$  units  $\left( \sum_{i=1}^k n_i = n \right)$ . Let the product estimator of  $Y$  based on the sample  $s_j$  is denoted by

$$\hat{Y}_p(j) = \hat{X}_j \hat{Y}_j / X \quad (9.4.7)$$

where  $\hat{Y}_j = n \sum_{i \in s_j} \xi_i / n_j$ ,  $\hat{X}_j = n \sum_{i \in s_j} \eta_i / n_j$ ,  $\xi_i = b_{si} y_i$ , and  $\eta_i = b_{si} x_i$ .

The proposed modified product estimator is defined as follows:

$$\hat{Y}_p(m) = \sum_{i=1}^k w_i \hat{Y}_p(j) + \left( 1 - \sum_{i=1}^k w_i \right) \hat{Y}_p \quad (9.4.8)$$

where  $w_i$ 's are suitably chosen weights.

Denoting  $\varepsilon_{jx} = \{\hat{X}_j - X\} / X$  and  $\varepsilon_{jy} = (\hat{Y}_j - Y) / Y$ , we have the following theorem due to Chaudhuri and Arnab (1982a,b).

#### Theorem 9.4.1

The modified product estimator  $\hat{Y}_p(m)$  is unbiased for  $Y$  and has the same efficiency as  $\hat{Y}_p$  upto the first order of approximation if

$$(i) \sum_{i=1}^k c_i w_i + \left( 1 - \sum_{i=1}^k w_i \right) = 0 \quad \text{and} \quad (ii) \quad w_i = \frac{n_i}{n} \left( \sum_{i=1}^k w_i \right)$$

where  $c_i = E(\varepsilon_{ix} \varepsilon_{iy}) / E(\varepsilon_x \varepsilon_y)$ .

**Proof**

$$\begin{aligned}
 E\{\hat{Y}_p(m)\} &= \sum_{i=1}^k w_i E\{\hat{Y}_p(j)\} + \left(1 - \sum_{i=1}^k w_i\right) E(\hat{Y}_p) \\
 &= Y \left\{ \sum_{i=1}^k w_i E(1 + \varepsilon_{ix} + \varepsilon_{iy} + \varepsilon_{ix}\varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right) \right. \\
 &\quad \left. E(1 + \varepsilon_x + \varepsilon_y + \varepsilon_x\varepsilon_y) \right\} \\
 &= Y \left[ \left\{ \sum_{i=1}^k w_i c_i + \left(1 - \sum_{i=1}^k w_i\right) \right\} E(\varepsilon_x\varepsilon_y) + 1 \right]
 \end{aligned} \tag{9.4.9}$$

From Eq. (9.4.9), we note that  $E\{\hat{Y}_p(m)\}$  is unbiased if and only if

$$\sum_{i=1}^k c_i w_i + \left(1 - \sum_{i=1}^k w_i\right) = 0 \tag{9.4.10}$$

The MSE of  $\hat{Y}_p(m)$  is

$$\begin{aligned}
 M\{\hat{Y}_p(m)\} &= E\{\hat{Y}_p(m) - Y\}^2 \\
 &= Y^2 \left[ \sum_{i=1}^k w_i (1 + \varepsilon_{ix} + \varepsilon_{iy} + \varepsilon_{ix}\varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right) \right. \\
 &\quad \left. (1 + \varepsilon_x + \varepsilon_y + \varepsilon_x\varepsilon_y) - 1 \right]^2
 \end{aligned}$$

Now neglecting terms with  $E(\varepsilon_x^r \varepsilon_y^l)$  for  $r + l > 2$ , we get upto the first order of approximation:

$$\begin{aligned}
 M\{\hat{Y}_p(m)\} &\cong Y^2 E \left[ \sum_{i=1}^k w_i (\varepsilon_{ix} + \varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right) (\varepsilon_x + \varepsilon_y) \right]^2 \\
 &= Y^2 E \left[ E \left\{ \sum_{i=1}^k w_i (\varepsilon_{ix} + \varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right) (\varepsilon_x + \varepsilon_y) \right\}^2 \middle| s \right] \\
 &= Y^2 E \left[ V \left\{ \sum_{i=1}^k w_i (\varepsilon_{ix} + \varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right) (\varepsilon_x + \varepsilon_y) \middle| s \right\} \right] \\
 &\quad + Y^2 E \left[ E \left\{ \sum_{i=1}^k w_i (\varepsilon_{ix} + \varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right) (\varepsilon_x + \varepsilon_y) \middle| s \right\}^2 \right]
 \end{aligned} \tag{9.4.11}$$

Because the first term of Eq. (9.4.11) is nonnegative and  $E\left\{\sum_{i=1}^k w_i(\varepsilon_{ix} + \varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right)(\varepsilon_x + \varepsilon_y)|s\right\} = \varepsilon_x + \varepsilon_y$ , we find

$$M\{\widehat{Y}_p(m)\} \geq Y^2 E(\varepsilon_x + \varepsilon_y)^2 = M(\widehat{Y}_p) \quad (9.4.12)$$

Equality in Eq. (9.4.12) holds if and only if

$$\begin{aligned} & \sum_{i=1}^k w_i(\varepsilon_{ix} + \varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right)(\varepsilon_x + \varepsilon_y) \\ &= E\left\{\sum_{i=1}^k w_i(\varepsilon_{ix} + \varepsilon_{iy}) + \left(1 - \sum_{i=1}^k w_i\right)(\varepsilon_x + \varepsilon_y)|s\right\} = \varepsilon_x + \varepsilon_y \end{aligned} \quad (9.4.13)$$

for every  $s$  with  $p(s) > 0$ .

Now noting  $\varepsilon_x + \varepsilon_y = \sum_{i=1}^k \frac{n_i}{n}(\varepsilon_{ix} + \varepsilon_{iy})$ , the above condition (Eq. 9.4.13) reduces to

$$\sum_{i=1}^k w_i(\varepsilon_{ix} + \varepsilon_{iy}) = \left(\sum_{i=1}^k w_i\right) \sum_{i=1}^k \frac{n_i}{n}(\varepsilon_{ix} + \varepsilon_{iy}) \quad (9.4.14)$$

The condition (Eq. 9.4.14) is realized if we choose

$$w_i = \frac{n_i}{n} \left( \sum_{i=1}^k w_i \right) \quad \text{for } i = 1, \dots, k \quad (9.4.15)$$

The theorem follows from Eqs. (9.4.10) and (9.4.15).

#### 9.4.5.1 Simple Random Sampling Without Replacement

Let  $b_{si} = N/n$ , then  $\widehat{X} = N\bar{x}(s)$ ,  $\widehat{Y} = N\bar{y}(s)$ ,  $\widehat{X}_j = N \sum_{i \in s_j} x_i/n_j$ ,  $\widehat{Y}_j = N \sum_{i \in s_j} y_i/n_j$ , and  $c_j = 1 + \left(\frac{1}{n_j} - \frac{1}{n}\right) \left(\frac{1}{n} - \frac{1}{N}\right)^{-1}$ . In this case Eq. (9.4.10) reduces to

$$\left(\frac{1}{n} - \frac{1}{N}\right)^{-1} \sum_{j=1}^k w_j \left(\frac{1}{n_j} - \frac{1}{n}\right) + 1 = 0 \quad (9.4.16)$$

In particular if we choose  $n_j = n/k$ , Eqs. (9.4.15) and (9.4.16) yield a very simple solution of  $w_j$  as



$$w_j = -\frac{1}{k(k-1)} \frac{N-n}{N} \quad \text{for } j = 1, \dots, k \quad (9.4.17)$$

Shukla (1976) derived Eq. (9.4.17) for  $k = 2$ .

## 9.5 COMPARISON BETWEEN THE RATIO, REGRESSION, PRODUCT, AND CONVENTIONAL ESTIMATORS

The regression estimator is always superior to the ratio, product, and the conventional estimator provided the estimate  $\hat{\beta}$  becomes very close to the true value  $\beta$ . In practice, the reliability of the estimator of the regression coefficient  $\beta$  is in question, especially if the sample size is small. Hence, in this situation, the regression estimator may not be efficient. Noting the difference  $M(\hat{Y}_p) - M(\hat{Y}_R) = 4RCov(\hat{X}, \hat{Y})$ , we conclude (i) for  $R > 0$ , the ratio estimator is superior to the product estimator if  $\rho_{\hat{X}, \hat{Y}} > 0$  and the product estimator is superior to the ratio estimator if  $\rho_{\hat{X}, \hat{Y}} < 0$ ; (ii) for  $R < 0$ , the ratio estimator is superior to the product estimator if  $\rho_{\hat{X}, \hat{Y}} < 0$  but for  $\rho_{\hat{X}, \hat{Y}} > 0$ , the product estimator is superior. In case  $\left| \rho_{\hat{X}, \hat{Y}} \right| < \frac{1}{2} \frac{C_{\hat{X}}}{C_{\hat{Y}}}$ , the conventional estimator  $\hat{Y}$  is always better than the product and ratio estimator. In particular, if  $C_{\hat{X}} = C_{\hat{Y}}$ , then the conventional estimator becomes more efficient than both the ratio and product estimators when  $\left| \rho_{\hat{X}, \hat{Y}} \right| < 1/2$ .

## 9.6 DUAL TO RATIO ESTIMATOR

Srivenkataramana (1980) proposed the following product type estimator, known as a dual to ratio estimator

$$\hat{Y}_{dR} = \frac{\hat{Y} \hat{X}_{\bar{s}}}{X} \quad (9.6.1)$$

where  $X_{\bar{s}} = \sum_{i \in \bar{s}} \frac{x_i}{1 - \pi_i}$  is an estimator for the total  $X$  based on the set  $\bar{s} (= U - s)$  of  $N - n$  units, which are not selected in the sample  $s$ . So the estimator  $\hat{Y}_{dR}$  is usable when the values of the auxiliary variable  $x$  are known at unit level. For SRSWOR,  $\pi_i = n/N$  and the estimator  $\hat{Y}_{dR}$  reduces to

$$\begin{aligned}\hat{Y}_{dR} &= \frac{\hat{Y}(X - n\bar{x}_s)N}{(N - n)X} \\ &= \frac{N\hat{Y} - n\hat{Y}_p}{N - n}\end{aligned}\quad (9.6.2)$$

where  $\hat{Y}_p$  is the product estimator defined in Eq. (9.4.1).

### 9.6.1 Bias of the Dual Estimator

Now writing

$$\hat{Y}_{dR} - Y = \frac{N(\hat{Y} - Y) - n(\hat{Y}_p - Y)}{N - n} \quad (9.6.3)$$

we find the bias of  $\hat{Y}_{dR}$  as

$$B(\hat{Y}_{dR}) = -\frac{n}{N - n}B(\hat{Y}_p) \quad (9.6.4)$$

Furthermore, from Section 9.4.4.2 we note

$$B(\hat{Y}_{dR}) = -\frac{S_{xy}}{\bar{X}} \quad (9.6.5)$$

The expression (Eq. 9.6.4) indicates that the absolute bias of  $\hat{Y}_{dR}$  is less than the absolute bias of the product estimator  $\hat{Y}_p$  if the sampling fraction  $f = n/N < 1/2$ . The expression (Eq. 9.6.5) indicates that  $\hat{Y}_{dR}$  is an inconsistent estimator for  $Y$  because the bias does not decrease with the increase in the sample size.

### 9.6.2 Mean Square Error of the Dual Estimator

Writing  $\varepsilon_x = \frac{\hat{X} - X}{X}$  and  $\varepsilon_y = \frac{\hat{Y} - Y}{Y}$ , we find

$$\hat{Y}_{dR} - Y = Y \frac{(N - n)\varepsilon_y - n\varepsilon_x - n\varepsilon_x\varepsilon_y}{N - n} \quad (9.6.6)$$

The expression of the MSE of  $\hat{Y}_{dR}$  upto the first-order approximation comes out as

$$\begin{aligned}M(\hat{Y}_{dR}) &= E(\hat{Y}_{dR} - Y)^2 \\ &= N^2(1 - f)\left(S_y^2 + g^2R^2S_x^2 - 2gRS_{xy}\right)/n\end{aligned}\quad (9.6.7)$$

where  $g = n/(N - n)$ .

### 9.6.3 Comparison With Other Estimators

#### 9.6.3.1 Conventional Estimator

The  $\hat{Y}_{dR}$  becomes more efficient than the conventional estimator  $\hat{Y} = N\bar{y}(s)$  if

$$(i) \rho_{xy} \geq \frac{g}{2} \frac{C_x}{C_y} \quad \text{with } R > 0, \quad \text{or} \quad (ii) \rho_{xy} \leq -\frac{g}{2} \left| \frac{C_x}{C_y} \right| \quad \text{with } R < 0.$$

#### 9.6.3.2 Ratio Estimator

$$M(\hat{Y}_{dR}) - M(\hat{Y}_R) = (g - 1) \left[ (g + 1)R^2 S_x^2 - 2R\rho_{xy} S_x S_y \right] \quad (9.6.8)$$

The estimator  $\hat{Y}_{dR}$  will be more efficient than the ratio estimator  $\hat{Y}_R$  if the difference (Eq. 9.6.8) is negative. Here we consider the following situations:

- (i)  $g - 1 > 0$ , i.e.,  $f > 1/2$ , in this case the difference (Eq. 9.6.8) is negative if either  $\rho_{xy} > \frac{g+1}{2} \frac{C_x}{C_y}$  with  $R > 0$  or  $\rho_{xy} < -\frac{g+1}{2} \left| \frac{C_x}{C_y} \right|$  with  $R < 0$ . The situation  $f > 1/2$  rarely happens in practice.
- (ii)  $g - 1 < 0$ , i.e.,  $f < 1/2$ , in this case (Eq. 9.6.8) is negative if either  $\rho_{xy} < \frac{g+1}{2} \frac{C_x}{C_y}$  with  $R > 0$  or  $\rho_{xy} > -\frac{g+1}{2} \left| \frac{C_x}{C_y} \right|$  with  $R < 0$ . The situation  $f < 1/2$  most likely happens in practice.

#### 9.6.3.3 Product Estimator

From Eq. (9.6.7) and Section 9.4.4.2, we note that

$$M(\hat{Y}_{dR}) - M(\hat{Y}_P) = (g + 1) \left[ (g - 1)R^2 S_x^2 - 2R\rho_{xy} S_x S_y \right] \quad (9.6.9)$$

The estimator  $\hat{Y}_{dR}$  will be more efficient than  $\hat{Y}_P$  if the difference (Eq. 9.6.9) is negative, i.e., if (i)  $\rho_{xy} > \frac{g-1}{2} \frac{C_x}{C_y}$  with  $R > 0$  or (ii) if  $\rho_{xy} < -\frac{g-1}{2} \left| \frac{C_x}{C_y} \right|$  with  $R < 0$ .

## 9.7 CALIBRATION ESTIMATORS

The calibration estimator was proposed by Deville and Särndal (1992). In calibration estimators, the weights of the conventional unbiased estimators are modified by using auxiliary information. The weights of the

conventional estimator are calibrated in such a way that the average distances between the calibrated and original weights attain a minimum according to a certain distance measure and at the same time, the calibrated weights satisfy some constraints that are known as calibration equations. The calibrated estimators become asymptotically design consistent and more efficient than the corresponding conventional estimator.

Suppose our objective is to estimate the population total  $Y$  on the basis of a sample  $s$  selected with probability  $p(s)$ . Let the conventional estimator for  $Y$  be the HTE

$$\hat{Y}_{ht} = \sum_{i \in s} d_i y_i \quad (9.7.1)$$

where  $d_i = 1/\pi_i$  and  $\pi_i$  is the inclusion probability for the  $i$ th unit, assumed to be positive for every  $i \in U$ .

In case  $X$ , the population total of the auxiliary variable  $x$  is known, we use this known total  $X$  to form a calibrated estimator

$$\hat{Y}_c = \sum_{i \in s} w_i y_i \quad (9.7.2)$$

where  $w_i$ 's are the weights that satisfy the following calibration constraint or calibration equation.

$$\sum_{i \in s} w_i x_i = X \quad (9.7.3)$$

Here it is important to note that in deriving the constraints (Eq. 9.7.3), no model between  $x$  and  $y$  has been assumed. To determine the modified weights  $w_i$ 's, we define the distance function  $D(w, d)$ , which satisfies the following conditions.

For every fixed  $d > 0$ , (i)  $D(w, d) \geq 0$ , differentiable with respect to  $w$ ; (ii) strictly convex on an interval containing  $d$ , and such that  $D(d, d) = 0$ ; and (iii)  $D'(w, d) = \frac{\partial D(w, d)}{\partial w}$  is continuous.

Deville and Särndal (1992) proposed the following distance functions: (i)  $D(w, d) = (w_i - d_i)^2 / (d_i)$ , (ii)  $D(w, d) = w_i \log (w_i / d_i) - w_i + d_i$ , (iii)  $D(w, d) = 2(\sqrt{w_i} - \sqrt{d_i})^2$ , (iv)  $-d_i \log (w_i / d_i) + w_i - d_i$ , and (v)  $D(w, d) = (w_i - d_i)^2 / (2w_i)$

The most popular is the chi-square type distance function

$$D(w_i, d_i) = \frac{(w_i - d_i)^2}{d_i} \quad (9.7.4)$$

Here we minimize the average chi-square distance

$$\sum_{i \in s} \frac{(w_i - d_i)^2}{d_i q_i} \text{ subject to } \sum_{i \in s} w_i x_i = X \text{ for every } s \quad (9.7.5)$$

where  $1/q_i$  is a known positive weight.

For minimization consider

$$\phi = \frac{1}{2} \sum_{i \in s} \frac{(w_i - d_i)^2}{d_i q_i} - \lambda \left( \sum_{i \in s} w_i x_i - X \right) \quad (9.7.6)$$

with  $\lambda$  as an undetermined Lagrange multiplier.

$$\frac{\partial \phi}{\partial w_i} = 0 \text{ implies} \quad (9.7.7)$$

$$w_i = d_i + \lambda d_i q_i x_i$$

On eliminating  $w_i$  from Eqs. (9.7.3) and (9.7.7), we get

$$\lambda = \frac{(X - \sum_{i \in s} d_i x_i)}{\sum_{i \in s} d_i q_i x_i^2} \quad (9.7.8)$$

Eqs. (9.7.7) and (9.7.8) yield

$$w_i = d_i + d_i q_i x_i \frac{(X - \sum_{i \in s} d_i x_i)}{\sum_{i \in s} d_i q_i x_i^2} \quad (9.7.9)$$

Finally substituting Eq. (9.7.9) in Eq. (9.7.2), we get the calibrated estimator as

$$\begin{aligned} \hat{Y}_c &= \sum_{i \in s} d_i y_i - B_s \left( \sum_{i \in s} d_i x_i - X \right) \\ &= \hat{Y}_{ht} - B_s (\hat{X}_{ht} - X) \end{aligned} \quad (9.7.10)$$

where

$$\hat{X}_{ht} = \sum_{i \in s} d_i x_i \text{ and } B_s = \frac{\sum_{i \in s} d_i q_i x_i y_i}{\sum_{i \in s} d_i q_i x_i^2} \quad (9.7.11)$$

The resulting calibrated estimator (Eq. 9.7.10) is known as generalized regression estimator (Cassel et al., 1976 and Särndal, 1982).

#### Remark 9.7.1

The calibrated weights  $w_i$  may take negative values. Distance functions (ii), (iii), (iv) and (v) guarantee positive weights. Details of properties are given by Deville and Särndal (1992), Singh and Mohl (1996), and Anderson and Thorburn (2005), among others.

#### Example 9.7.1: Ratio Estimator

If we choose  $q_i = 1/x_i$ , then we get  $B_s = \frac{\sum_{i \in s} d_i y_i}{\sum_{i \in s} d_i x_i}$  and the estimator  $\hat{Y}_c$  in Eq. (9.7.10) reduces to

$$\hat{Y}_c = \frac{\hat{Y}_{ht}}{\hat{X}_{ht}} X = \hat{Y}_R \quad (9.7.12)$$

the well known ratio estimator.

#### Example 9.7.2

If we choice of  $q_i = 1$  and  $d_i = N/n$ , then the estimator (Eq. 9.7.10) reduces to

$$\hat{Y}_c = N \left[ \bar{y}(s) - \frac{\sum_{i \in s} y_i x_i}{\sum_{i \in s} x_i^2} (\bar{x}(s) - \bar{X}) \right] \quad (9.7.13)$$

The estimator (Eq. 9.7.13) is quite different from the regression estimator given in Eq. (9.7.12).

### 9.7.1 Efficiency of Calibrated Estimator

The calibrated estimator  $\hat{Y}_c$  is asymptotically design unbiased. The asymptotic variance (AV) of  $\hat{Y}_c$  given by Deville and Särndal (1992) is

$$AV(\hat{Y}_c) = \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} (d_i E_i - d_j E_j)^2 \quad (9.7.14)$$

where  $\Delta_{ij} = \pi_i \pi_j - \pi_{ij}$ ,  $E_i = y_i - Bx_i$ , and  $B = \frac{\sum_{i \in U} x_i y_i q_i}{\sum_{i \in U} x_i^2 q_i}$ .

From the expression (Eq. 9.7.14), it is evident that the estimator  $\hat{Y}_c$  is more efficient than  $\hat{Y}_{ht}$  in most situations since  $AV(\hat{Y}_c)$  is expected to be smaller than  $V(\hat{Y}_{ht})$ .

A design-consistent estimator of  $AV(\hat{Y}_c)$  is given by

$$\hat{V}(\hat{Y}_c) = \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} (w_i \hat{E}_i - w_j \hat{E}_j)^2 \quad (9.7.15)$$

where  $\hat{E}_i = y_i - \hat{B}_s x_i$  and  $\hat{B}_s = \left( \sum_{i \in s} w_i q_i x_i y_i \right) / \left( \sum_{i \in s} w_i q_i x_i^2 \right)$ .

### 9.7.2 Calibration Estimator for Several Auxiliary Variables

Suppose we have  $p$  auxiliary variables  $x_1, \dots, x_p$  with known totals  $X_1, \dots, X_p$ , respectively. In this situation the weights  $w_i$  are obtained by minimizing

$$\sum_{i \in s} \frac{(w_i - d_i)^2}{d_i q_i}$$

Subject to

$$\sum_{i \in s} w_i x_{1i} = X_1, \dots, \sum_{i \in s} w_i x_{ji} = X_j, \dots, \sum_{i \in s} w_i x_{pi} = X_p$$

where  $x_{ji}$  be the value of the  $i$ th unit of the  $j$ th auxiliary variable  $x_j$ ,  $j = 1, \dots, p$ .

For this minimization problem, we consider

$$\begin{aligned} \phi = & \frac{1}{2} \sum_{i \in s} \frac{(w_i - d_i)^2}{d_i q_i} - \lambda_1 \left( \sum_{i \in s} w_i x_{1i} - X_1 \right) - \dots - \lambda_j \left( \sum_{i \in s} w_i x_{ji} - X_j \right) \\ & - \dots - \lambda_p \left( \sum_{i \in s} w_i x_{pi} - X_p \right) \end{aligned}$$

with  $\lambda_1, \dots, \lambda_p$  as Lagrange multipliers.

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

$$w_i = d_i + \lambda_1 d_i q_i x_{1i} + \dots + \lambda_j d_i q_i x_{ji} + \dots + \lambda_p d_i q_i x_{pi} \quad \text{for } i = 1, \dots, n \quad (9.7.16)$$

On multiplying both sides of Eq. (9.7.16) by  $x_{ki}$  then summing over  $i \in s$  and writing  $\sum_{i \in s} w_i x_{ki} = X_k$ , for  $k = 1, \dots, p$ , we get estimated value of  $\lambda' = (\lambda_1, \dots, \lambda_p)$  as

$$\hat{\lambda} = \mathbf{A}_s^{-1} \mathbf{T}_s \quad (9.7.17)$$

where

$$\mathbf{T}_s = \begin{pmatrix} X_1 - \sum_{i \in s} d_i x_{1i} \\ \dots \\ X_j - \sum_{i \in s} d_i x_{ji} \\ \dots \\ X_p - \sum_{i \in s} d_i x_{pi} \end{pmatrix},$$

$$\mathbf{A}_s = \begin{pmatrix} \sum_{i \in s} d_i q_i x_{1i}^2 & \dots & \sum_{i \in s} d_i q_i x_{ji} x_{1i} & \dots & \sum_{i \in s} d_i q_i x_{pi} x_{1i} \\ \dots & \dots & \dots & \dots & \dots \\ \sum_{i \in s} d_i q_i x_{1i} x_{ji} & \dots & \sum_{i \in s} d_i q_i x_{ji}^2 & \dots & \sum_{i \in s} d_i q_i x_{pi} x_{ji} \\ \dots & \dots & \dots & \dots & \dots \\ \sum_{i \in s} d_i q_i x_{1i} x_{pi} & \dots & \sum_{i \in s} d_i q_i x_{pi} x_{ji} & \dots & \sum_{i \in s} d_i q_i x_{pi}^2 \end{pmatrix} \text{ and}$$

$$\hat{\lambda} = \begin{pmatrix} \hat{\lambda}_1 \\ \dots \\ \hat{\lambda}_j \\ \dots \\ \hat{\lambda}_p \end{pmatrix}$$

Eq. (9.7.17) yields

$$\begin{aligned} \hat{Y}_c &= \sum_{i \in s} w_i y_i \\ &= \hat{Y}_{ht} + \hat{\lambda}_1 \sum_{i \in s} d_i q_i x_{1i} y_i + \dots + \hat{\lambda}_j \sum_{i \in s} d_i q_i x_{ji} y_i + \dots + \hat{\lambda}_p \sum_{i \in s} d_i q_i x_{pi} y_i \\ &= \hat{Y}_{ht} - B_{1s}(\hat{X}_{ht1} - X_1) - \dots - B_{js}(\hat{X}_{htj} - X_j) - \dots - B_{ps}(\hat{X}_{htp} - X_p) \end{aligned} \quad (9.7.18)$$



where  $\hat{X}_{l_{ij}} = \sum_{i \in s} d_i x_{ji}$ ,  $\mathbf{B}'_s \mathbf{A}_s = \left( \sum_{i \in s} d_i q_i x_{1i} y_i, \dots, \sum_{i \in s} d_i q_i x_{ji} y_i, \dots, \sum_{i \in s} d_i q_i x_{pi} y_i \right)$ , and  $\mathbf{B}'_s = (B_{1s}, \dots, B_{js}, \dots, B_{ps})$ .

Following Deville and Särndal (1992), the asymptotic variance (AV) of  $\hat{Y}_c$  and its approximate unbiased estimators are obtained as

$$AV(\hat{Y}_c) = \sum_{i \neq j} \sum_{j \in U} \Delta_{ij} \left( d_i E_i^* - d_j E_j^* \right)^2 \quad (9.7.19)$$

and

$$\hat{V}(\hat{Y}_c) = \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( d_i \hat{E}_i^* - d_j \hat{E}_j^* \right)^2 \quad (9.7.20)$$

where  $E_i^* = y_i - B_1 x_{1i} - \dots - B_j x_{ji} - \dots - B_p x_{pi}$ ,  $\hat{E}_i^* = y_i - B_{1s} x_{1i} - \dots - B_{js} x_{ji} - \dots - B_{ps} x_{pi}$ , and

$$\begin{pmatrix} B_1 \\ \dots \\ B_j \\ \dots \\ B_p \end{pmatrix} = \begin{pmatrix} \sum_{i \in U} q_i x_{1i}^2 & \dots & \sum_{i \in U} q_i x_{ji} x_{1i} & \dots & \sum_{i \in U} q_i x_{pi} x_{1i} \\ \dots & \dots & \dots & \dots & \dots \\ \sum_{i \in U} q_i x_{1i} x_{ji} & \dots & \sum_{i \in U} q_i x_{ji}^2 & \dots & \sum_{i \in U} q_i x_{pi} x_{ji} \\ \dots & \dots & \dots & \dots & \dots \\ \sum_{i \in U} q_i x_{1i} x_{pi} & \dots & \sum_{i \in U} q_i x_{pi} x_{ji} & \dots & \sum_{i \in U} q_i x_{pi}^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i \in U} q_i x_{1i} y_i \\ \dots \\ \sum_{i \in U} q_i x_{ji} y_i \\ \dots \\ \sum_{i \in U} q_i x_{pi} y_i \end{pmatrix}$$

The variance estimator (Eq. 9.7.20) provides strictly design-consistent variance estimator. When design-based as well as model-based properties are considered, the following alternative variance estimator is proposed.

$$\hat{V}_w(\hat{Y}_c) = \sum_{i \neq j} \sum_{j \in s} \frac{\Delta_{ij}}{\pi_{ij}} \left( w_i \hat{E}_i^* - w_j \hat{E}_j^* \right)^2$$

## 9.8 EXERCISES

- 9.8.1** Describe the regression methods of estimation. Derive approximate expressions of bias and mean square errors of the regression estimator of the population mean. Show that the regression estimator is more efficient than the ratio estimator for estimating the population mean.
- 9.8.2** Consider the sampling scheme where the first unit is selected with probability proportional to  $(x_i - \bar{X})^2$  and the rest  $n - 1$  units are

selected from the remaining  $(N - 1)$  units by SRSWOR method.

Show that for this sampling scheme

- (i) probability of selection of an unordered samples of  $n$  units is

$$p(s) = \frac{\sum_{i \in s} (x_i - \bar{X})^2}{\binom{N-1}{n-1} \sum_{i \in U} (x_i - \bar{X})^2} \text{ and}$$

- (ii)  $\frac{n(N-1)}{N(n-1)} \left[ \bar{y}(s) - \hat{\beta}^* \{ \bar{x}(s) - \bar{X} \} \right]$  is an unbiased estimator  $\bar{Y}$

$$\text{where } \hat{\beta}^* = \frac{\sum_{i \in s} y_i (x_i - \bar{X})}{\sum_{i \in s} (x_i - \bar{X})^2} \text{ (Singh and Srivastava, 1980).}$$

- 9.8.3** Let  $s$  be sample selected from a population by SRSWR method.

Show that the regression estimator is more efficient than the ratio estimator under the following superpopulation model where  $x$  is used as an auxiliary variable.

$$E_m(y|x_i) = \alpha + \beta x_i, V_m(y_i|x_i) = \sigma^2 x_i^g \text{ and } Cov(y_i, y_j|x_i, x_j) = 0 \text{ for } i \neq j$$

with  $\sigma^2, g > 0$  (Raj, 1958).

- 9.8.4** To determine the average height of 50 plants in a certain field, a sample of 10 plants was selected by SRSWOR method. The following table gives the exact height and eye estimate heights.

Serial number of plants	Eye estimate height (in m)	Exact height (in m)
1	15	16.8
2	10	12.2
3	20	18.5
4	8	9.0
5	5	4.8
6	10	10.5
7	15	16.2
8	10	10.0
9	15	16.3
10	14	12.2
11	12	12.0
12	8	9.5
13	12	12.8
14	9	8.5
15	6	5.8

Given that the average eye estimate of heights for the 50 plants is 16.5 m.

- (i) Estimate the average height of the plant by ratio, regression, and product method of estimation.
- (ii) Estimate the biases of the estimators used in (i).
- (iii) Estimate the relative efficiencies of the above estimators with respect to the sample mean.

**9.8.5** A sample of 10 households is selected from a population of 50 households by PPSWR method of sampling using household size as a size variable. The information of household income, household expenditure, and amount of income tax paid were obtained and presented in the following table.

Sampled households	1	2	3	4	5	6	7	8	9	10
Size	5	4	3	2	4	5	6	4	4	3
Expenditure on food (\$)	2000	2500	1500	1000	2500	3000	4000	2000	1500	1200
Income (\$)	5500	4500	5000	2500	7500	6000	8500	6000	2500	5000
Income tax paid (\$)	1000	850	950	400	1400	1100	1450	1050	400	950

From the past survey it is known that the average family size and average income of the population are 4.2 and \$5000, respectively. The total tax collected from the 50 households was \$45,000. Use this information estimate the following:

- (i) Average expenditure on food by using the regression estimator taking tax paid and income as auxiliary variables.
- (ii) Give two separate estimates of the average expenditure of food by using the regression method using income and income tax paid as auxiliary variables. Give a pooled estimate of these two regression estimates.
- (iii) Estimate relative efficiency of the regression estimators used in “i” with respect to the pooled estimator used in “ii.”

**9.8.6** Using the data given in Exercise 9.8.5: (i) Calibrate the Hansen–Hurwitz estimator for the population mean consumption on food using total tax collection as \$45,000. Estimate the standard error of the calibrated estimator you used. (ii) Obtain the calibrated estimator

for the average consumption on food using average income and total tax collection, which were \$5000 and \$45,000, respectively. (iii) Compare precision of the calibrate estimator used in “i” and “ii.”

**9.8.7** To estimate the average production of milk per farm, the farms were classified into three strata: small, medium, and big with respect to cattle population. From each of the strata, samples were selected by SRSWOR method and information on production of milk and the number of cattle was collected.

Sample							
Strata	Number of farms	Average number of cattle	Size	Mean no. of cattle	Mean yield of milk (kg)	SD of no. of cattle	SD of yield of milk (kg)
Small	550	10	50	12	22,000	15	200
Medium	400	25	30	28	15,000	25	325
Large	250	45	20	50	10,000	10	250

- Estimate the average production of milk per farm using (a) separate and (b) combined regression methods taking cattle population as an auxiliary variable.
- Estimate the relative efficiency of the separate regression estimation with respect to the combined regression estimator and comment on your findings.

## APPENDIX 9A

$$\text{Cov}(\bar{y}_s, s_x^2 | \text{SRSWOR}) = \text{Cov}\{s_{xy}, (\bar{x}_s - \bar{X})\} = \frac{1}{n} \frac{N(N-n)}{(N-1)(N-2)} \mu_{21}$$

**Proof:**

$$\begin{aligned} & \text{Cov}(\bar{y}_s, s_x^2 | \text{SRSWOR}) \\ &= E[\bar{y}_s - \bar{Y}, s_x^2] \\ &= A_0 \sum_s \left[ \sum_{i \in s} h_i \left\{ \sum_{i \in s} z_i^2 - \left( \sum_{i \in s} z_i \right)^2 / n \right\} \right], z_i = x_i - \bar{X} \end{aligned}$$

$$\begin{aligned}
& \text{(where } A_0 = \left[ n(n-1) \binom{N}{n} \right]^{-1} \text{)} \\
&= \frac{A_0}{n} \left[ \sum_s \left\{ (n-1) \sum_{i \in s} h_i \sum_{i \in s} z_i^2 - \left( \sum_{i \in s} h_i \right) \left( \sum_{i \neq j \in s} z_i z_j \right) \right\} \right] \\
&= \frac{A_0}{n} \left[ \sum_s (n-1) \left\{ \sum_{i \in s} h_i z_i^2 + \sum_{i \neq j \in s} h_i z_j^2 \right\} - \left\{ 2 \sum_{i \neq j \in s} h_i z_i z_j \right. \right. \\
&\quad \left. \left. + \sum_{i \neq j \neq k \in s} h_i z_j z_k \right\} \right] \\
&= \frac{A_0}{n} \left[ (n-1) \left\{ M_1 \sum_{i \in U} h_i z_i^2 + M_2 \sum_{i \neq j \in U} h_i z_j^2 \right\} \right. \\
&\quad \left. - \left\{ 2M_2 \sum_{i \neq j \in U} h_i z_i z_j + M_3 \sum_{i \neq j \neq k \in U} h_i z_j z_k \right\} \right]; M_i = \binom{N-i}{n-i} \\
&= \frac{A_0}{n} \left[ (n-1)(M_1 - M_2) \sum_{i \in U} h_i z_i^2 + 2 \left\{ M_2 \sum_{i \in U} h_i z_i^2 - M_3 \sum_{i \in U} z_i^2 h_i \right\} \right] \\
&= \frac{A_0}{n} \left[ \left\{ (n-1)(M_1 - M_2) + 2(M_2 - M_3) \right\} \right] \sum_{i \in U} h_i z_i^2 \\
&= \frac{A_0}{n} \left[ \left\{ (n-1)M_1(1 - M_2/M_1) + 2M_2(1 - M_3/M_2) \right\} \right] \sum_{i \in U} h_i z_i^2 \\
&= \frac{A_0}{n} M_1 \left[ (n-1) \frac{N-n}{N-1} + 2 \frac{(n-1)}{N-1} \frac{N-n}{N-2} \right] \sum_{i \in U} h_i z_i^2 \\
&= \frac{1}{n^2} \frac{n}{N} \frac{(N-n)}{N-1} \frac{N}{N-2} \sum_{i \in U} h_i z_i^2 \\
&= \frac{1}{n} \frac{(N-n)}{N-1} \frac{1}{N-2} \sum_{i \in U} h_i z_i^2 \\
&= \frac{1}{n} \frac{(N-n)}{N-1} \frac{1}{N-2} \sum_{i \in U} (y_i - \bar{Y})(x_i - \bar{X})^2 \\
&= \frac{1}{n} \frac{N(N-n)}{(N-1)(N-2)} \mu_{21}
\end{aligned} \tag{A1}$$

Similarly,

$$\begin{aligned}
 & Cov\{s_{xy}, (\bar{x}_s - \bar{X})\} \\
 &= E \left[ \left\{ \frac{1}{n-1} \sum_{i \in s} (y_i - \bar{Y})(x_i - \bar{x}_s) \right\} \left\{ \frac{1}{n} \sum_{i \in s} (x_i - \bar{X}) \right\} \right] \\
 &= E \left[ \frac{1}{(n-1)n} \sum_{i \in s} h_i (z_i - \bar{z}_s) \sum_{i \in s} z_i \right] \\
 &= \frac{1}{(n-1)n^2 M_o} \sum_s \left[ \left\{ (n-1) \sum_{i \in s} h_i z_i - \sum_{i \neq j} \sum_{j \in s} h_i z_j \right\} \sum_{i \in s} z_i \right] \\
 &= \frac{1}{(n-1)n^2 M_o} \left[ (n-1)(M_1 - M_2) + 2(M_2 - M_3) \right] \sum_{i \in U} h_i z_i^2 \\
 &= \frac{1}{n} \frac{N(N-n)}{(N-1)(N-2)} \mu_{21} \tag{A2}
 \end{aligned}$$

From (A1) and (A2), we verify the result.