

VIVEKANANDA COLLEGE  
THAKURPUKUR  
KOLKATA-700063

NAAC ACCREDITED 'A' GRADE



Topic: Theory of Ratio and Regression Estimator

Course Title: Sample Survey

Paper: CC-8

Unit: 2

Semester: 4

Name of the Teacher: Riddhi Das Majumder

Name of the Department: Statistics

# Ratio and Regression Estimators: Use of auxiliary information

## (A) Ratio Estimation:

Frequently we come across situations in which the ratio of  $y$  to another character  $x$  is believed to be less variable than the  $y$ 's themselves. In that case it would be better to estimate  $R$ , the ratio of  $y$  to  $x$  in the population, from the sample and then multiply it by the known total of  $x$  to estimate the total for  $y$ . This procedure is called ratio estimation.

Frequently we wish to estimate a ratio rather than a total or mean, for example, it is desired to estimate the total agricultural area in a region containing  $N$  Communes. There are very big Communes and very small Communes and this makes the character  $y$  vary tremendously over the region. But the ratio of agricultural area and the popl. size of the Commune, which is the per capita agricultural area, would be less variable.

Let  $Y$ , and  $X$  be the total agricultural area and the total popl. in the region. Then the per capita agricultural area in the region is  $R = \frac{Y}{X}$ . If a simple random sample of  $n$

Communes gives  $\sum_{i=1}^n y_i$  and  $\sum_{i=1}^n x_i$  as the total for  $y$  and  $x$ , respectively, it is natural to estimate  $R$  by  $\hat{R} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i} = \frac{\bar{y}}{\bar{x}}$

and the total of  $y$  (i.e.  $Y$ ) is estimated by  $\hat{Y}_R = \hat{R} \cdot X = \frac{\bar{y}}{\bar{x}} X$ .

where  $X$  is known total of  $x$ . It should be noted that the two problems are different, although they are connected. For estimating  $Y$  we could have used information on any character  $x$ ; this information need not be recent, but must be known for the entire population. On the otherhand, information on a sample basis is required for  $y$  as well as for  $x$  (the denominator of the ratio) if the purpose is to estimate the ratio  $R = \frac{Y}{X}$  in population.

Since the theory is the same in either case, most of the subsequent results will relate to the problem of estimating a ratio.

Bias of the ratio estimator: The following theorem gives the exact bias associated with  $\hat{R}$ . (28)

Theorem: In simple random sampling, bias of the ratio estimator  $\hat{R} = \frac{\bar{y}}{\bar{x}}$  is given by  $B(\hat{R}) = -\frac{\text{Cov}(\hat{R}, \bar{x})}{E(\bar{x})}$

Proof: As  $\text{Cov}(\frac{\bar{y}}{\bar{x}}, \bar{x}) = E(\bar{y}) - E(\frac{\bar{y}}{\bar{x}})E(\bar{x})$ , we have

$$\bar{x} E(\frac{\bar{y}}{\bar{x}}) = \bar{y} - \text{Cov}(\frac{\bar{y}}{\bar{x}}, \bar{x}), \Rightarrow E(\hat{R}) = R - \frac{1}{\bar{x}} \text{Cov}(\hat{R}, \bar{x})$$

Therefore  $B(\hat{R}) = E(\hat{R}) - R = -\frac{1}{\bar{x}} \text{Cov}(\hat{R}, \bar{x})$

Corollary: Denoting the standard deviation of  $\hat{R}$  by  $\sigma(\hat{R})$ , we have  $B(\hat{R}) = -\frac{1}{\bar{x}} \cdot \sigma(\hat{R}) \sigma(\bar{x}) \rho(\hat{R}, \bar{x})$ .

or  $\frac{B(\hat{R})}{\sigma(\hat{R})} = -\rho(\hat{R}, \bar{x}) \cdot \frac{\sigma(\bar{x})}{\bar{x}} = -\rho(\hat{R}, \bar{x}) \cdot \text{C.V.}(\bar{x})$

Hence,  $|\frac{B(\hat{R})}{\sigma(\hat{R})}| \leq \text{C.V.}(\bar{x})$ , since  $|\rho(\hat{R}, \bar{x})| \leq 1$ .

where C.V. stands for the coefficient of variation. The same bound applies, of course, to the bias in  $\hat{Y}_R$  and  $\hat{y}_R$ .

Remarks (1)  $\hat{R}$  is consistent for  $R$  in the sense that  $\hat{R} \rightarrow R$  when the sample size is  $n$ .

(2) The bias associated with  $\hat{Y}_R = \hat{R}X$  is  $X B(\hat{R})$ .

(3)  $\hat{R}$  is unbiased if  $\rho(\hat{R}, \bar{x}) = 0$ .

Theorem The approximate bias and mean square error (MSE) of the ratio estimator  $\hat{R}$  are  $B(\hat{R}) = \frac{(\frac{1}{n} - \frac{1}{N})}{\bar{x}^2} (R S_y^2 - \rho S_y S_x)$

and  $\text{MSE}(\hat{R}) = \frac{(\frac{1}{n} - \frac{1}{N})}{\bar{x}^2} (S_y^2 + R^2 S_x^2 - 2R\rho S_x S_y)$

Proof: Define  $e_0 = \frac{\bar{y} - \bar{Y}}{\bar{y}}$ , and  $e_1 = \frac{\bar{x} - \bar{X}}{\bar{x}}$ .

It may be noted that

(i)  $E(e_0) = E(\frac{\bar{y} - \bar{Y}}{\bar{y}}) = 0$  (ii)  $E(e_1) = 0$

(iii)  $E(e_0^2) = E(\frac{(\bar{y} - \bar{Y})^2}{\bar{y}^2}) = \frac{V(\bar{y})}{\bar{y}^2}$  (iv)  $E(e_1^2) = \frac{V(\bar{x})}{\bar{x}^2}$

(v)  $E(e_0 e_1) = E\left\{\frac{(\bar{x} - \bar{X})(\bar{y} - \bar{Y})}{\bar{x} \bar{y}}\right\} = \frac{\text{Cov}(\bar{x}, \bar{y})}{\bar{x} \bar{y}}$

Assume that the sample size is large enough so that  $|e_0| < 1$  and  $|e_1| < 1 \Leftrightarrow 0 < \bar{x} < 2\bar{x}, 0 < \bar{y} < 2\bar{y}$ .

Since  $\bar{y} = \bar{Y} (1+e_0)$ ,  $\bar{x} = \bar{X} (1+e_1)$  the estimator  

$$= \frac{\bar{y}}{\bar{x}} \text{ can be written as } \hat{R} = \frac{\bar{Y} (1+e_0)}{\bar{X} (1+e_1)} = R (1+e_0) (1+e_1)^{-1}$$

$= R \{1 + e_0 - e_1 + e_1^2 + e_0 e_1 + \dots\}$ . Hence,  $E(\hat{R}) - R = B(\hat{R})$   

$$B(\hat{R}) \approx R \{E(e_1^2) - E(e_0 e_1)\} = R \left\{ \frac{V(\bar{x})}{\bar{x}^2} - \frac{\text{Cov}(\bar{x}, \bar{y})}{\bar{x} \bar{y}} \right\}$$
  

$$= \left( \frac{1}{n} - \frac{1}{N} \right) \left\{ R \frac{S_x^2}{\bar{x}^2} - \rho S_x S_y \right\}$$

[ In SRSWOR,  $V(\bar{x}) = \left( \frac{1}{n} - \frac{1}{N} \right) S_x^2$ ,  $V(\bar{y}) = \left( \frac{1}{n} - \frac{1}{N} \right) S_y^2$   
 and  $\text{Cov}(\bar{x}, \bar{y}) = \left( \frac{1}{n} - \frac{1}{N} \right) S_{xy} = \left( \frac{1}{n} - \frac{1}{N} \right) \rho S_x S_y$  ]

Again,  $\text{MSE}(\hat{R}) = E(\hat{R} - R)^2 \approx R^2 E[e_0^2 + e_1^2 - 2e_0 e_1]$   
 ignoring terms of degree greater than two.

Therefore  $\text{MSE}(\hat{R}) \approx R^2 \left\{ \frac{V(\bar{x})}{\bar{x}^2} + \frac{V(\bar{y})}{\bar{y}^2} - 2 \frac{\text{Cov}(\bar{x}, \bar{y})}{\bar{x} \bar{y}} \right\}$   

$$\approx \left( \frac{1}{n} - \frac{1}{N} \right) \left\{ R^2 \frac{S_x^2}{\bar{x}^2} + S_y^2 - 2R \rho S_x S_y \right\}$$

Remark:

(1)  $B(\hat{Y}_R) = B(N\bar{x} \hat{R}) = N \left( \frac{N-n}{Nn} \right) \frac{1}{\bar{x}} \left\{ R S_x^2 - \rho S_x S_y \right\}$   
 and  $\text{MSE}(\hat{Y}_R) = N^2 \frac{N-n}{nN} \left\{ R^2 S_x^2 + S_y^2 - 2R \rho S_x S_y \right\}$ .

(2) The quantity  $\frac{\text{Bias}}{\text{S.E}}$ , which is the same for  $\hat{R}$ ,  $\hat{Y}_R$ ,  $\hat{\bar{y}}_R$   
 may be expressed as  $\frac{\text{Bias}}{\text{S.E}} = \frac{\text{C.V}(\bar{x}) \cdot (R S_x - \rho S_y)}{\sqrt{\left\{ R^2 S_x^2 + S_y^2 - 2R \rho S_x S_y \right\}}}$

where  $\text{C.V}(\bar{x}) = \frac{\sqrt{V(\bar{x})}}{\bar{x}} = \frac{\left( \frac{1}{n} - \frac{1}{N} \right)^{1/2} S_x}{\bar{x}}$

The following theorem gives the Condition under which the Ratio estimator will be more efficient than the Conventional expansion estimator, or estimator based on the mean per unit ( $\bar{y}$ )

Theorem: The ratio estimator  $\hat{Y}_R = \frac{\bar{y}}{\bar{x}} X = \hat{R} X$  is more efficient than the expansion estimator  $\hat{y}$ , with simple random sample, if  $\rho > \frac{1}{2} \frac{\text{C.V}(x)}{\text{C.V}(y)} = \frac{1}{2} \left( \frac{S_x}{\bar{x}} \right) \left( \frac{S_y}{\bar{y}} \right)$ .

$$V(\hat{Y}_R) > \text{MSE}(\hat{Y}_R), \text{ under SRS} \quad (29)$$

$$\Rightarrow N^2 \frac{N-n}{nN} S_y^2 > N^2 \frac{N-n}{nN} \{ R^2 S_x^2 + S_y^2 - 2R\rho S_x S_y \}$$

= approximately.

$$\Rightarrow S_y^2 > \{ S_y^2 + R^2 S_x^2 - 2R\rho S_x S_y \}$$

$$\Rightarrow \rho > \frac{1}{2} \cdot R \cdot \frac{S_x}{S_y} = \frac{1}{2} \cdot \left( \frac{S_x}{\bar{x}} \right) / \left( \frac{S_y}{\bar{y}} \right) = \frac{1}{2} \cdot \frac{C.V.(x)}{C.V.(y)}, \text{ if } R > 0$$

Hence the proof. [The theorem holds for large samples]

Estimated MSE under Simple Random Sampling:

$$\text{Note that } \sum_{i=1}^N [Y_i - R X_i]^2 = \sum_{i=1}^N [Y_i - \bar{y} + \bar{y} - R X_i]^2$$

$$= \sum_{i=1}^N [Y_i - \bar{y} + R\bar{x} - R X_i]^2, \text{ since } R = \frac{\bar{y}}{\bar{x}}$$

$$= \sum_{i=1}^N [Y_i - \bar{y}]^2 + R^2 \sum_{i=1}^N [X_i - \bar{x}]^2 - 2R \sum_{i=1}^N (X_i - \bar{x})(Y_i - \bar{y})$$

$$\Rightarrow \frac{1}{N-1} \sum_{i=1}^N [Y_i - R X_i]^2 = S_y^2 + R^2 S_x^2 - 2R\rho S_x S_y$$

$$\text{Then, we have } \text{MSE}(\hat{R}) = \left( \frac{1}{n} - \frac{1}{N} \right) \frac{1}{\bar{x}^2} \cdot \frac{1}{N-1} \sum_{i=1}^N (Y_i - R X_i)^2$$

Therefore a reasonable estimator for the MSE of the ratio estimate is

$$V(\hat{R}) = \left( \frac{1}{n} - \frac{1}{N} \right) \frac{1}{\bar{x}^2} \cdot \frac{\sum_{i=1}^n (y_i - \hat{R} x_i)^2}{n-1} \text{ where } \hat{R} = \frac{\bar{y}}{\bar{x}}$$

The estimator is biased. ~~\_\_\_\_\_~~

$$\text{Now } \frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{R} x_i)^2 = \frac{1}{n-1} \left\{ \sum_{i=1}^n y_i^2 + \hat{R}^2 \sum_{i=1}^n x_i^2 - 2\hat{R} \sum_{i=1}^n x_i y_i \right\}$$

$$\text{or } = \frac{1}{n-1} \sum_{i=1}^n \{ y_i - \bar{y} - \hat{R} (x_i - \bar{x}) \}^2 = \{ S_y^2 + \hat{R}^2 S_x^2 - 2\hat{R} S_{xy} \}$$

$$\text{where } S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{Hence, } V(\hat{R}) = \left( \frac{1}{n} - \frac{1}{N} \right) \frac{1}{\bar{x}^2} \{ S_y^2 + \hat{R}^2 S_x^2 - 2\hat{R} S_{xy} \}$$

$$\text{Since } \hat{Y}_R = \hat{R} \bar{X}_N, \quad V(\hat{Y}_R) = N^2 \left( \frac{1}{n} - \frac{1}{N} \right) \{ S_y^2 + \hat{R}^2 S_x^2 - 2\hat{R} S_{xy} \}$$

## Unbiased Ratio-type estimators

To modify the usual ratio estimator itself ~~modification~~ so that a ratio-type estimator is obtained that is unbiased under a simple random sampling. Actually, the estimator  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n \frac{y_i}{x_i} x_i$  is corrected for its bias to obtain an unbiased estimator.

Theorem: In simple random sampling, an unbiased estimator of  $R = \bar{y}/\bar{x}$  is given by

$$\hat{R} = \bar{r} + \frac{(N-1)n}{N(n-1)} \cdot \frac{\bar{y} - \bar{r}\bar{x}}{\bar{x}}$$

Proof:  $\frac{1}{N} \sum_{i=1}^N R_i (x_i - \bar{x})$ , where  $R_i = \frac{y_i}{x_i}$ ,  $i = 1, \dots, N$ .

$$= \frac{1}{N} \sum_{i=1}^N (y_i - \bar{x} R_i) = \bar{y} - \bar{x} \cdot \frac{1}{N} \sum_{i=1}^N R_i = \bar{y} - \bar{x} \cdot E(R_i)$$

But in simple random sampling  $E(\bar{r}) = E(R_i)$ .

$$\text{Hence, bias in } \bar{r} = E(\bar{r}) - R = -\frac{1}{N} \sum_{i=1}^N R_i (x_i - \bar{x}) \quad (*)$$

Again, an unbiased estimator of  $\frac{1}{N-1} \sum_{i=1}^N R_i (x_i - \bar{x})$

$$\text{is } \frac{1}{n-1} \sum_{i=1}^n x_i (x_i - \bar{x}) = \frac{1}{n-1} \sum_{i=1}^n (\bar{y} - r_i \bar{x})$$

$$= \frac{n}{n-1} (\bar{y} - \bar{r}\bar{x})$$

$$\text{From } (*), \text{ bias in } \bar{r} = E(\bar{r}) - R = -\frac{(N-1)n}{N(n-1)} \cdot \frac{(\bar{y} - \bar{r}\bar{x})}{\bar{x}}$$

$$\Rightarrow E\left\{ \bar{r} + \frac{(N-1)n}{N(n-1)} \cdot \frac{\bar{y} - \bar{r}\bar{x}}{\bar{x}} \right\} = R$$

Hence,  $\hat{R}_* = \bar{r} + \frac{(N-1)n}{N(n-1)} \cdot \frac{(\bar{y} - \bar{r}\bar{x})}{\bar{x}}$  is an unbiased estimator

of  $R = \frac{\bar{y}}{\bar{x}}$ .

Remark:

(i) The corresponding unbiased estimator of the popl. total  $Y$  is  $\hat{Y}_R = \hat{R}_* X = \bar{r} X + \frac{(N-1)n}{n-1} (\bar{y} - \bar{r}\bar{x})$ .

(ii) An unbiased estimator of the popl. mean  $\bar{y}$  is  $\hat{Y}_R = \hat{R}_* \bar{x} = \bar{r}\bar{x} + \frac{(N-1)n}{N(n-1)} (\bar{y} - \bar{r}\bar{x})$ .

## Regression Estimator

Like the ratio estimator, the linear regression estimator is designed to increase precision by the use of an auxiliary variable  $x_i$  that is ~~assumed~~ correlated with  $y_i$ . The ratio estimator is at its best when the relation between  $y$  and  $x$  is a straight line through the origin, that is,  $y - kx = 0 \Leftrightarrow y/x = k$ . When the relation between  $y_i$  and  $x_i$  is examined, it may be found that although the relation is (approximately) linear, the line does not go through the origin. This suggests an estimator based on the linear regression of  $y$  on  $x$  rather than on the ratio of the variables.

We suppose that  $y_i$  and  $x_i$  are each obtained for every unit in the sample and that the population mean  $\bar{x}$  of the  $x_i$  is known. The linear regression estimator of  $\bar{y}$ , the pop'n mean of  $y_i$ , is

$$\bar{y}_{lr} = \bar{y} + b(\bar{x} - \bar{x}_i)$$

where  $b$  is an estimator of the change in  $y$  when  $x$  is increased by unit. The rationale behind this estimator is that if  $\bar{x}$  is below average we should expect  $\bar{y}$  also to be below average by an amount  $b(\bar{x} - \bar{x})$  because of the regression of  $y$  on  $x$ . For an estimator of the pop'n  $\bar{y}$ , we take  $\hat{y}_{lr} = N \bar{y}_{lr}$ .

Suppose that we can take a rapid estimate  $x_i$  of some characteristic for every unit and can also, by some more costly method, determine the correct value  $y_i$  of the characteristic for a simple random sample of the units. For an example, an eye estimate of the volume of timber was made on each of a population of  $\frac{1}{10}$ -acre plots, and the actual timber volume was measured for a simple random sample of the plots. The regression estimator

$$\bar{y} + b(\bar{x} - \bar{x}_i)$$

adjusts the sample mean of the actual measurements by the regression of the actual measurements on the rapid estimates.

By a suitable choice of  $b$ , the regression estimator includes as particular cases both the mean per unit and the ratio estimate. Obviously if  $b$  is taken as zero,  $\bar{y}_{lr}$  reduces to  $\bar{y}$ . If  $b = \frac{\bar{y}}{\bar{x}}$ ,  $\bar{y}_{lr} = \bar{y} + \frac{\bar{y}}{\bar{x}}(\bar{x} - \bar{x}_i) = \frac{\bar{y}}{\bar{x}} \bar{x} = \hat{y}_R$

Regression Estimator When b is computed from the sample

Let  $y = \bar{y} + B(x - \bar{x})$  be the pop'n regression line of y on x, where  $B = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$  is the pop'n regression coeff.

Here, 'b' must be (the sample estimate of B, that is) the least squares estimate of B, that is,  $b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$

(i) Bias of the Linear regression Estimator

Introduce the variate  $e_i = y_i - \bar{y} - B(x_i - \bar{x})$ .

The properties of  $e_i$  are:  $\sum_{i=1}^n e_i = 0$

and  $\sum_{i=1}^n e_i (x_i - \bar{x}) = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) - B \sum_{i=1}^n (x_i - \bar{x})^2 = 0$

, by def'n. of B.

Now,  $b = \frac{\sum_{i=1}^n y_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$

$= \frac{\sum_{i=1}^n [e_i + \bar{y} + B(x_i - \bar{x})] (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$

$= B + \left\{ \frac{\sum_{i=1}^n e_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right\}$

We have  $E(\bar{y}_1) = \bar{y} - E\{b(x - \bar{x})\}$ . Thus one expression for bias is  $-E\{b(x - \bar{x})\} = -\text{Cov}(b, \bar{x})$ .

Now,  $-\text{Cov}(b, \bar{x}) = -\text{Cov}\left(B + \frac{\sum_{i=1}^n e_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \bar{x}\right)$

$= -\frac{\text{Cov}\left(\sum_{i=1}^n e_i (x_i - \bar{x}), \bar{x}\right)}{\sum_{i=1}^n (x_i - \bar{x})^2} = -\frac{\text{Cov}\left\{\sum_{i=1}^n e_i (x_i - \bar{x}) + n\bar{e}(\bar{x} - \bar{x}), \bar{x}\right\}}{\sum_{i=1}^n (x_i - \bar{x})^2}$

$= -\frac{E\left\{\sum_{i=1}^n e_i (x_i - \bar{x})(\bar{x} - \bar{x})\right\}}{nS_x^2} + \frac{E\left\{\bar{e}(\bar{x} - \bar{x})\right\}}{S_x^2}$  replacing  $\sum_{i=1}^n (x_i - \bar{x})$  by its leading term  $n\bar{e}$

$= -\frac{E(\bar{U} - \bar{U})(\bar{x} - \bar{x})}{S_x^2}$ , where  $U_i = e_i (x_i - \bar{x})$  and  $\bar{U} = 0$ .

$= -\left(\frac{1}{n} - \frac{1}{N}\right) \cdot \frac{E(U_i - \bar{U})(x_i - \bar{x})}{S_x^2} = -\frac{1-f}{n} \cdot \frac{E\{e_i (x_i - \bar{x})^2\}}{S_x^2}$

which is, the bias turns out to be in the estimator  $\bar{y}_1$ .

(23)  
 The approximate MSE of the regression estimator under SRSWOR  
 Theorem: If  $b$  is the least square estimate of  $B$  and  
 $\bar{y}_{lr} = \bar{y} + b(\bar{x} - \bar{x})$ , then in SRSWOR of size  $n$ , with  $n$  large.

MSE or  $V(\bar{y}_{lr}) \approx \frac{1-f}{n} S_y^2 (1 - \rho^2)$ , where  $\rho = \frac{S_{yx}}{S_x S_y}$  is the population correlation between  $y$  and  $x$ .

Proof: The sampling error of  $\bar{y}_{lr}$  arises from the quantity  
 $\bar{y}_{lr} - \bar{Y} = \bar{y} - \bar{Y} + b(\bar{x} - \bar{x})$ .

As an approximation, replace  $\bar{y}_{lr}$  by  $\bar{y}_{lr}^* = \bar{y} + B(\bar{x} - \bar{x})$   
 where  $B$  is the population linear regression coefficient of  $y$  on  $x$ .

The error committed in this approximation is  $(B-b)(\bar{x} - \bar{x})$

Note that  $(b-B) = O(\frac{1}{\sqrt{n}})$  and  $(\bar{x} - \bar{x}) = O(\frac{1}{\sqrt{n}})$ , hence  $(b-B)(\bar{x} - \bar{x})$   
 is of order  $\frac{1}{n}$  in SRS. Again  $V(\bar{y}_{lr}^*)$  is of order  $\frac{1}{n}$ , since it  
 is the variance of the sample mean of the variate  $(y - Bx)$ .

Hence,  $E(\bar{y}_{lr} - \bar{Y})^2 = E\{(\bar{y}_{lr}^* - \bar{Y} + (B-b)(\bar{x} - \bar{x}))^2\}$   
 $= V(\bar{y}_{lr}^*) + E[(b-B)^2(\bar{x} - \bar{x})^2] + 2E[(\bar{y}_{lr}^* - \bar{Y})(b-B)(\bar{x} - \bar{x})]$

Now,  $E[(b-B)^2(\bar{x} - \bar{x})^2] \leq \{E(b-B)^2 E(\bar{x} - \bar{x})^2\}^{1/2}$ , which is  
 of order  $1/n^2$ . Similarly,  $E[(b-B)(\bar{y}_{lr}^* - \bar{Y})(\bar{x} - \bar{x})]$   
 $\leq \{E(b-B)^2\}^{1/2} \{E(\bar{y}_{lr}^* - \bar{Y})^2 E(\bar{x} - \bar{x})^2\}^{1/2}$ , which is of order  $1/n^{3/2}$

Thus, the large sample variance of the regression estimator  
 $\bar{y}_{lr}$  is  $V(\bar{y}_{lr}) \approx V(\bar{y}_{lr}^*) = \left(\frac{1-f}{n} - \frac{1}{N}\right) S_y^2 (1 - \rho^2) = \frac{1-f}{n} S_y^2 (1 - \rho^2)$

Sample estimate of the MSE or Variance:

Note that  $V(\bar{y}_{lr}) = \frac{1-f}{n} s_c^2$  where  $s_c^2 = S_y^2 (1 - \rho^2)$ .

Note that, an unbiased estimator of  $s_c^2 = \frac{1}{n-1} \sum_{i=1}^n (e_i)^2$  is  $s_c^2 = \frac{\sum_{i=1}^n (e_i - \bar{e})^2}{n-1}$

Now,  $e_i - \bar{e} = y_i - \bar{y} - B(x_i - \bar{x}) = \{y_i - \bar{y} - b(x_i - \bar{x})\} + (b-B)(x_i - \bar{x})$

The 2nd term on the right, of order  $\frac{1}{\sqrt{n}}$ , may be neglected in  
 relation to the 1st term, which is of order unity.

Hence, in large sample  $s_c^2 = \frac{1}{(n-1)S_y^2} \sum_{i=1}^n \{y_i - \bar{y} - b(x_i - \bar{x})\}^2$  as  $\frac{(S_y^2 - b^2 S_x^2)}{S_y^2} = S_y^2 (1 - \rho^2)$   
 estimate of  $s_c^2$ . The estimator  $\frac{1}{n-2} \sum_{i=1}^n \{y_i - \bar{y} - b(x_i - \bar{x})\}^2$  is  
 as suggested since it is used in regression theory

is an estimate of  $V(\bar{y}_{LR})$  or  $MSE(\bar{y}_{LR})$ , valid in large samples, we may use

$$V(\bar{y}_{LR}) = \frac{1-f}{n(n-2)} \sum_{i=1}^n \{(y_i - \bar{y}) - b(x_i - \bar{x})\}^2$$

$$= \frac{1-f}{n(n-2)} \left\{ \sum_{i=1}^n (y_i - \bar{y})^2 - \frac{[\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})]^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right\}$$

Comparison with the Ratio and the Mean per unit, in large Sample:

Theorem: Under simple random sampling, with large sample,  $V(\bar{y}) > MSE(\bar{y}_{LR})$  and  $MSE(\bar{y}_R) > MSE(\bar{y}_{LR})$

Proof:

$$MSE(\bar{y}_{LR}) \text{ or } V(\bar{y}_{LR}) = \frac{N-n}{Nn} S_y^2 (1-\rho^2), \text{ [Regression]}$$

$$MSE(\bar{y}_R) \text{ or } V(\bar{y}_R) = \frac{N-n}{Nn} (S_y^2 + R^2 S_x^2 - 2R\rho S_x S_y), \text{ [Ratio]}$$

$$V(\bar{y}) = \frac{N-n}{Nn} S_y^2 \text{ (Mean per unit).}$$

$$\text{Since } |\rho| < 1, (1-\rho^2) < 1 \Rightarrow V(\bar{y}) > V(\bar{y}_{LR}) \text{ or } MSE(\bar{y}_{LR})$$

$$\text{Now, } MSE(\bar{y}_R) - MSE(\bar{y}_{LR}) = \frac{N-n}{Nn} \{ S_y^2 + R^2 S_x^2 - 2R\rho S_x S_y - S_y^2 + S_y^2 \rho^2 \}$$

$$= \frac{N-n}{Nn} \{ R S_x - S_y \rho \}^2 > 0.$$

The regression estimator is more precise than the ratio estimator unless  $B = R$  (~~is equal to the ratio estimator~~)  $\Rightarrow y = kx$ , that is, the relation between  $y$  and  $x$  is a straight line through the origin.

$$MSE(\hat{y}_R) - MSE(\hat{y}_{LR}) = \left(\frac{1}{n} - \frac{1}{N}\right) (R S_x - S_y \rho)^2 > 0$$

$$'=' \text{ holds iff } R S_x = S_y \rho, \text{ iff } R = \rho \cdot \frac{S_y}{S_x} = B.$$

which is true if  $y = kx$ , i.e. if the relation b/w  $y$  and  $x$  is a st. line passing through origin.