

Ratio Estimators.

Definition and Notation:

y_i - the value of the characteristics under study for the i^{th} unit of the population.

x_i - the value of the auxiliary characteristics on the same unit.

\bar{y} - the total of y characteristics of the population

\bar{x} - the total of x characteristics of the population.

$R = \frac{\bar{y}}{\bar{x}} = \frac{\bar{Y}}{\bar{x}}$ = the ratio of the population totals
or mean of characteristics y and x .

ρ = the correlation coefficient between x and y in the population

The ratio estimators of the population ratio $\bar{y}/\bar{x} = R$,
the total \bar{Y} , and the mean \bar{Y} .

$$\hat{R} = \frac{\bar{y}}{\bar{x}} = \frac{\bar{Y}}{\bar{x}}$$

$$\hat{Y}_R = \frac{\bar{y}}{\bar{x}} x = \frac{\bar{Y}}{\bar{x}} x = \hat{R} x.$$

$$\hat{\bar{Y}}_R = \frac{\bar{y}}{\bar{x}} \bar{x} = \hat{R} \bar{x}$$

Bias of Ratio estimators:

Biased Estimators are to be used when

these are comparatively more efficient. There can be many illustrations where one would like to use such estimators in order to obtain higher precision without changing the cost

Theorem: 1

In Simple Random Sampling, the bias of the ratio estimator \hat{R} is given by

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

Proof:

$$\text{Cov}(u, v) = E(uv) - E(u)E(v)$$

$$\begin{aligned}\text{Cov}(\hat{R}, \bar{x}) &= E(\hat{R}\bar{x}) - E(\bar{x})E(\hat{R}) \\ &= E\left(\frac{\bar{y}}{\bar{x}} \cdot \bar{x}\right) - E(\bar{x})E(\hat{R}) \\ &= E(\bar{y}) - E(\bar{x})E(\hat{R}).\end{aligned}$$

$$= \bar{Y} - \bar{x} E(\hat{R})$$

$$\text{Cov}(\hat{R}, \bar{x}) = \bar{x} (\bar{Y}/\bar{x} - E(\hat{R}))$$

$$\frac{\text{Cov}(\hat{R}, \bar{x})}{\bar{x}} = (\bar{Y}/\bar{x} - E(\hat{R}))$$

$$= (R - E(\hat{R})) = -B(\hat{R}).$$

Corollary: 1

$$\frac{|B(\hat{R})|}{\sigma_{\hat{R}}} \leq CV(\bar{x}).$$

CV - Correlation coefficient of Variation.

Corollary: 2

$$B(\hat{Y}_R) = \bar{x} B(\hat{R}) \text{ and } B(\bar{Y}_R) = \bar{x} B(\hat{R}).$$

Theorem:- 2

S.T the first approximation to the relative bias of the ratio estimator in simple random sampling, w.r.t. is given by.

$$\frac{B(\hat{R})}{R} \cong \frac{1-\frac{f}{n}}{n\bar{x}\bar{y}} (R s_x^2 - ps_n s_y) = B_1$$

$$(\hat{R} - R) \cong \frac{1-\frac{f}{n}}{n} (C_n^2 - R C_n C_y).$$

Where $C_n = \frac{s_n}{\bar{x}}$; $C_y = \frac{s_y}{\bar{y}}$ are the correlation coefficient of variation of x and y respectively

Proof:

$$\hat{R} - R = \frac{\bar{y}}{\bar{x}} - R = \frac{\bar{y} - R\bar{x}}{\bar{x}} = \frac{\bar{y} - R\bar{x}}{\bar{x} + (\bar{x} - \bar{x})}$$

$$= \frac{1}{\bar{x}} (\bar{y} - R\bar{x}) \left(1 + \frac{\bar{x} - \bar{x}}{\bar{x}} \right)^{-1}$$

Expanding by a Taylor's series we get

$$\hat{R} - R \cong \frac{1}{\bar{x}} (\bar{y} - R\bar{x}) \left(1 - \frac{\bar{x} - \bar{x}}{\bar{x}} + \dots \right)$$

Ignoring the second and higher orders.

$$\hat{R} - R \cong \frac{1}{\bar{x}} [E(\bar{y} - R\bar{x}) - \frac{1}{\bar{x}} E(\bar{y} - R\bar{x})(\bar{x} - \bar{x})]$$

$$E(\bar{y} - R\bar{x}) = \bar{Y} - R\bar{x} = 0 \text{ respectively.}$$

$$E(\bar{y} - R\bar{x})(\bar{x} - \bar{x}) = E[\bar{y}(\bar{x} - \bar{x}) - RE\{\bar{x}(\bar{x} - \bar{x})\}] \\ = E[(\bar{y} - \bar{Y})(\bar{x} - \bar{x}) - RE(\bar{x} - \bar{x})^2]$$

$$= \frac{1-f}{n} [R s_n s_y / R s_x^2]$$

$$B(\hat{R}) = \left\{ \frac{1-f}{n\bar{x}^2} \right\} (R s_{xy}^2 - f s_n s_y)$$

$$= \frac{1-f}{n\bar{x}\bar{y}} (R s_{xy}^2 - f s_n s_y)$$

$$= E[(\bar{y} - \hat{y})(\bar{x} - \hat{x})] - R E(\bar{x} - \hat{x})^2$$

$$= \frac{\sum (\bar{y} - \hat{y})(\bar{x} - \hat{x})}{n} - R \frac{\sum (\bar{x} - \hat{x})^2}{n}$$

$$\hat{f} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n s_n s_y}$$

$$f \text{ } \cancel{\times} \text{ } s_n s_y = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

$$B(\hat{R}) = \frac{1-f}{n\bar{x}\bar{y}} [R s_x^2 + s_n s_y - R s_n^2]$$

$$B(\hat{R}) = \frac{(1-f)}{n\bar{x}\bar{y}} [R s_x^2 - f s_n s_y] \quad [\because R = \frac{\bar{y}}{\bar{x}}]$$

$$= \frac{1-f}{n} \left[\frac{s_x^2}{\frac{\bar{x}}{\bar{y}} \times (\bar{x}\bar{y})} - f \frac{s_n s_y}{\bar{x}\bar{y}} \right]$$

$$= \frac{1-f}{n} \left[\frac{s_x^2}{\bar{x}^2} - f \frac{s_n s_y}{\bar{x}\bar{y}} \right]$$

$$= \frac{1-f}{n} [C_x^2 - f C_n C_y]$$

Approximate Variance of Ratio Estimator

Theorem: 3

$$\text{In SRSWOR for large } n \quad V(\hat{R}) = \frac{1-f}{n\bar{x}^2} \sum_{i=1}^N \frac{(y_i - Rx_i)^2}{N-1}$$

Pf:

$$V(\hat{R}) = E[(\hat{R} - E(\hat{R}))^2]$$

$$= E\left(\frac{\bar{y}}{\bar{x}} - R\right)^2$$

$$= E\left[\frac{\bar{y} - R\bar{x}}{\bar{x}}\right]^2$$

$$= \frac{E(\bar{y} - R\bar{x})^2}{E(\bar{x})^2} \quad E(\bar{x})^2 = \bar{x}^2$$

$$V(\hat{R}) = \frac{E(\bar{y} - R\bar{x})^2}{\bar{x}^2}$$

$$\text{let } u_i = y_i - Rx_i$$

$$\bar{y} = \bar{y} - R\bar{x} \quad \text{for sample}$$

$$\bar{v} = \bar{y} - R\bar{x} \quad \text{for population}$$

$$\bar{v} = \bar{y} - \frac{\bar{Y}}{\bar{x}}\bar{x} = 0.$$

$$V(\hat{R}) = \frac{E(u - v)^2}{\bar{x}^2}$$

$$\text{In SRSWOR, } V(\bar{Y}) = (1-f) \frac{s^2}{n}$$

$$\text{Here } s^2 = \sum \frac{(y_i - \bar{y})^2}{N-1}$$

$$V(\hat{R}) = \frac{1}{\bar{x}^2} \frac{(1-f)}{n} \sum_{i=1}^N \frac{(\bar{y}_p - \bar{y})^2}{N-1}$$

$$= \frac{1-f}{n\bar{x}^2} \sum_{i=1}^N \frac{(y_i - Rx_i)^2}{N-1}$$

COROLLARY:

$$V(\hat{y}) = V(\hat{R}\bar{x})$$

$$= \bar{x}^2 V(\hat{R})$$

$$= \bar{x}^2 \frac{1-f}{n\bar{x}^2} \sum_{i=1}^N \frac{(y_i - Rx_i)^2}{N-1}$$

$$V(\hat{y}) = \frac{1-f}{n} \sum_{i=1}^N \frac{(y_i - Rx_i)^2}{N-1}$$

Theorem. 4

Show that 1st order of approximate variance of \hat{R} can be expressed.

$$V(\hat{R}) = \frac{1-f}{n\bar{x}^2} (S_y^2 + R^2 S_x^2 - 2f RS_x S_y)$$

$$= \frac{(1-f)R^2}{n} [C_y^2 + C_x^2 - 2f C_x C_y].$$

PF:

$$\text{Since } \bar{Y} = R\bar{x}$$

$$y_i = y_i - Rx_i = (y_i - \bar{y}) - R(x_i - \bar{x})$$

Then,

$$V(\hat{R}) = \frac{1-f}{n\bar{x}^2(N-1)} \sum_{i=1}^N [(y_i - \bar{y}) - R(x_i - \bar{x})]^2$$

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

$$S_y^2 = \frac{\sum (y_i - \bar{y})^2}{N-1} \quad S_x^2 = \frac{\sum (x_i - \bar{x})^2}{N-1}$$

$$V(\hat{R}) = \frac{1-f}{n\bar{x}} [S_y^2 + R^2 S_x^2 - 2C S_x S_y]$$

$$= \frac{(1-f)R^2}{n} [C_y^2 + C_x^2 - \frac{2f}{n} C_x C_y]$$

Then,
Top row

$$\bar{Y} = R\bar{x}$$

$$V(\hat{R}) \leq \frac{(1-f)}{n\bar{x}^2} \leq \frac{(y_i - Rx_i)^2}{N-1}$$

$$= \frac{(1-f)}{n\bar{x}^2(N-1)} \sum (y_i - Rx_i)^2$$

$$\Rightarrow \boxed{u_i = y_i - Rx_i = (y_i - \bar{Y}) - R(x_i - \bar{x})}$$

$$\text{Now } V(\hat{R}) = \frac{(1-f)}{n\bar{x}^2(N-1)} \leq [(y_i - \bar{Y}) - R(x_i - \bar{x})]^2$$

$$= \frac{1-f}{n\bar{x}^2(N-1)} \left[\sum (y_i - \bar{Y})^2 + R^2 \sum (x_i - \bar{x})^2 - 2R \sum (x_i - \bar{x})(y_i - \bar{Y}) \right]$$

$$= \frac{1-f}{n\bar{x}^2} \left[\frac{\sum (y_i - \bar{Y})^2}{N-1} + R^2 \frac{\sum (x_i - \bar{x})^2}{N-1} - 2R \frac{\sum (x_i - \bar{x})(y_i - \bar{Y})}{N-1} \right]$$

Come off

$$P = \frac{\sum (y_i - \bar{Y})(x_i - \bar{x})}{N-1} S_x S_y$$

$$= \frac{1-f}{n\bar{x}^2} [S_y^2 + R^2 S_x^2 - 2R P S_x S_y]$$

$$= \frac{(1-f)}{n\bar{n}^2} R^2 \left[\frac{s_y^2}{R^2} + \frac{s_x^2}{\bar{n}^2} - 2f \frac{s_n s_y}{R} \right]$$

$$= \frac{1-f}{n\bar{n}^2} R^2 \left[\frac{s_y^2}{\bar{y}^2} + \frac{s_n^2}{\bar{x}^2} - 2f \frac{s_n s_y}{\bar{y} \cdot \bar{x}} \right]$$

$$= \frac{(1-f)}{n\bar{n}^2} R^2 \left[\frac{s_y^2}{\bar{x}^2 \bar{y}^2} \bar{x}^2 + \frac{s_x^2}{\bar{x}^2} - 2f \frac{s_n s_y}{\bar{x} \cdot \bar{y}} \bar{x} \right]$$

$$= \frac{1-f}{n} R^2 \left[\frac{s_y^2}{\bar{y}^2} + \frac{s_x^2}{\bar{x}^2} - 2f \frac{s_n s_y}{\bar{x} \cdot \bar{y}} \right]$$

$$= \frac{1-f}{n} R^2 [C_y^2 + C_x^2 - 2f C_n C_y].$$

Corollary:

$$V(\bar{Y}_R) = \frac{1-f}{n} [s_y^2 + R^2 s_n^2 - 2R f s_n s_y]$$

$$= \frac{1-f}{n} [C_y^2 \bar{y}^2 + R^2 C_n^2 \bar{x}^2 - 2R f C_n C_y \bar{y}]$$

$$= \frac{1-f}{n} [C_y^2 \bar{y}^2 + \frac{\bar{y}^2}{\bar{x}^2} C_n^2 \bar{x}^2 - 2R f C_n C_y \bar{y}]$$

$$= \frac{1-f}{n} \bar{y}^2 [C_y^2 + C_x^2 - 2f C_n C_y].$$

Corollary:

Ratio Estimators In Stratified Sampling :-

Separate Ratio estimator:

If y_m, x_m are the sample totals in the m^{th} stratum and X_m is the m^{th} stratum total.

We may define the estimator \hat{Y}_{RS} (S for separate)

as,

$$\hat{Y}_{RS} = \sum_m \frac{y_m}{x_m} \cdot X_m = \sum_m \bar{y}_m / \bar{x}_m \cdot X_m.$$

Theorem: 5 If the sample size n_m 's are large in all strata & SRS, w.r.t., & done independently within each stratum, show that \hat{Y}_{RS} is a biased estimator with negligible bias and its sampling variance.

$$V(\hat{Y}_{RS}) = \sum_m N_m^2 \frac{(1-f_m)}{n_m} (S_{ym}^2 + R_m^2 S_{xm}^2 - 2 R_m f_m S_{ym} S_{xm})$$

where $R_m = Y_m/X_m$ and f_m are true ratio and coefficient of correlation, resp. in the m^{th} stratum.

Proof:-

$$\hat{Y}_{RM} = \frac{y_m}{x_m} \cdot X_m$$

$$\hat{Y}_{RS} = \sum_m \hat{Y}_{RM}$$

$$\frac{B(\hat{Y})}{Y} = \frac{1-f}{n} (C_n^2 - \rho C_n C_y) \quad [\text{To refer Theorem: 2}]$$

$$= \sum_m \frac{(1-f)}{N_m} (C_{xm}^2 - \rho C_{xm} C_{ym}).$$

$$V(\hat{Y}_{RS}) = V\left(\sum_m \hat{Y}_{RM}\right)$$

$$= \sum_m V(\hat{Y}_{RM})$$

$$= \sum_m N_m^2 \frac{(1-f)}{n_m} (S_{ym}^2 + R_m^2 S_{xm}^2)$$

$$= \sum_m N_m^2 \frac{(1-f_m)}{n_m} (S_{ym}^2 + R_m^2 S_{xm}^2 - 2 R_m f_m S_{ym} S_{xm})$$

COROLLARY

In stratified random sampling, ~~we~~,
almost unbiased estimator of $V(\hat{Y}_{RS})$ is given by

$$V(\hat{Y}_{RS}) = \sum_m N_m^2 \frac{(1-f_m)}{n_m} (S_{ym}^2 + R_m^2 S_{xm}^2 - 2 R_m f_m S_{ym} S_{xm}).$$

Combined Ratio estimator:

It was assumed, in case of separate estimator, that n_m 's were large in each stratum.

However, it may not hold good always in practice.

To overcome this difficulty, Hansen, Hurwitz and Grunow (1946) suggested a combined ratio estimator \hat{Y}_{RC} (c for combined) for a sample from a stratified population as.

$$\hat{Y}_{RC} = \frac{\hat{Y}_{st}}{\hat{x}_{st}} \cdot X = \frac{\hat{Y}_{st}}{\bar{n}_{st}} \cdot X$$

$$\text{where } \hat{Y}_{st} = \frac{\hat{Y}_{st}}{N} = \sum_m \frac{N_m}{N} \hat{y}_m$$

$$\hat{x}_{st} = \frac{\bar{x}_{st}}{N} = \sum_m \frac{N_m}{N} \bar{x}_m.$$

\hat{x}_{st} are estimated population means from a stratified sample and X is overall total of x .

Theorem: 6

If the total sample size n is large and simple random sampling, w.r.t., is done in each stratum independently, then \hat{Y}_{RC} is a consistent estimator and its sampling variance is given by

$$V(\hat{Y}_{RC}) = \sum_m N_m^2 \frac{(1-f_m)}{n_m} (S_{ym}^2 + R^2 S_{xm}^2 - 2R f_m S_{ym} S_{xm})$$

Proof:

Since R is a consistent estimator, \hat{Y}_{RC} is a ratio estimator and therefore, also a consistent estimator of the population total, to derive its bias.

$$\hat{Y}_{RC} - Y = \frac{N\bar{X}}{\bar{n}_{st}} (\bar{y}_{st} - R\bar{x}_{st}) \cong N(\bar{y}_{st} - R\bar{x}_{st})$$

Let us define a variate $U_{mj} = y_{mj} - R x_{mj}$

$$\bar{u}_{st} = \bar{y}_{st} - R\bar{x}_{st} \text{ and.}$$

$$U = \bar{Y} - R\bar{X} = 0.$$

$$\begin{aligned} E(\bar{Y} - R\bar{X})(\bar{x} - \bar{X}) &= E[\bar{y} \cdot (\bar{x} - \bar{X})] - R E[\bar{x} \cdot (\bar{x} - \bar{X})] \\ &= E(\bar{y} - \bar{Y})(\bar{x} - \bar{X}) - R E(\bar{x} - \bar{X})(\bar{x} - \bar{X}) \\ &= E(\bar{y} - \bar{Y})(\bar{x} - \bar{X}) - R E(\bar{x} - \bar{X})^2 \\ &= \frac{1-f}{n} [f S_n S_y - \frac{R}{R} S_n^2] \end{aligned}$$

$$\begin{aligned} B(\hat{Y}_{RC}) &= \sum_m \frac{(1-f)}{n \bar{x}^2} \left[\frac{R S_n^2}{R} - f S_n S_y \right] \\ &= \sum_m \frac{1-f}{n \bar{x}^2} \left[S_x^2 - \frac{f S_n S_y}{R} \right] \\ &\leq \sum_m \frac{1-f}{n} \left[\frac{S_x^2}{\bar{x}^2} - \frac{f S_n S_y}{\frac{\bar{x}}{\bar{x}} \cdot \bar{x}^2} \right] \end{aligned}$$

$$B(\hat{Y}_{RC}) = \sum_{i=1}^m \frac{1-f_m}{n_m} \left(\frac{s_x^2}{n_i} + f \frac{s_{nx} s_y}{\bar{x} \bar{y}} \right)$$

$$\frac{B(\hat{Y}_{RC})}{\gamma} = \sum_{i=1}^m \frac{(1-f_m)}{n_m} W_m^2 \left(\frac{s_{nxm}^2}{\bar{x}^2} - \frac{f_m s_{nxm} s_{ym}}{\bar{x} \bar{y}} \right)$$

$$\therefore C_x = \frac{s_{xy}}{\bar{x}} ; C_y = \frac{s_y}{\bar{y}}$$

$$B = \frac{1-f}{n} (C_x^2 - PC_x C_y)$$

$$V(\hat{Y}_{RC}) = V(N \bar{Y}_{st}) = N^2 V(\bar{Y}_{st})$$

$$= \sum_{i=1}^m N_m^2 \underbrace{\frac{(1-f_m)}{n_m} s_{nx}^2}_{V(\hat{R})}$$

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

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

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

$$P = \frac{\sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x})}{(N-1)s_{nx}s_y}$$

$$[P(N-1)s_{nx}s_y = \sum_{i=1}^N (y_i - \bar{y})(x_i - \bar{x})]$$

$$V(\hat{R}) = \frac{1-f}{n \bar{x}^2} [s_y^2 + R^2 s_x^2 - 2R P s_{nx} s_y]$$

$$= \frac{1-f}{n \bar{x}^2} R^2 \left[\frac{s_y^2}{R^2} + s_x^2 - 2 \frac{P s_{nx} s_y}{R} \right]$$

$$= \frac{1-f}{n} R^2 \left[\frac{S_y^2}{\bar{Y}_m^2 - \bar{x}^2} + \frac{S_x^2}{\bar{x}^2} - \frac{2f S_{xy} S_y}{\bar{Y}_m \cdot \bar{x}^2} \right]$$

$$= \frac{1-f}{n} R^2 \left[\frac{S_y^2}{\bar{Y}^2} + \frac{S_x^2}{\bar{x}^2} - \frac{2f S_{xy} S_y}{\bar{Y} \bar{x}} \right]$$

$$V(\hat{Y}_{RC}) = \sum_m N_m^2 \frac{(1-f_m)}{n_m} \left[\frac{S_{ym}^2}{\bar{Y}_m^2} + \frac{S_{xm}^2}{\bar{x}_m^2} - \frac{2f S_{ym} S_{xm}}{\bar{x}_m \bar{Y}_m} \right] \quad (\text{or})$$

$$V(\hat{Y}_{RS}) = \sum_m N_m^2 \frac{(1-f_m)}{n_m} \left[S_{ym}^2 + R^2 S_{xm}^2 - 2R f_m S_{ym} S_{xm} \right]$$

Thus proved.

Comparison of separate and combined Ratio estimators:

$$V(\hat{Y}_{RC}) - V(\hat{Y}_{RS}) = \sum_m N_m^2 \frac{(1-f_m)}{n_m} \left[(R^2 - R_m^2) S_{xm}^2 - 2(R - R_m) f_m S_{ym} S_{xm} \right]$$

$$= \sum_m N_m^2 \frac{(1-f_m)}{n_m} \left[(R - R_m)^2 S_{xm}^2 + 2(R - R_m) (f_m S_{ym} S_{xm} - R_m S_{ym})^2 \right]$$

Since the last term

Since the last term is small, the right side is likely to be positive. Thus, the separate ratio estimator is expected to be more precise provided the sample in each stratum is large enough for the approximate variance formula to be applied. The combined estimator will have a large variance as compared to the separate estimator, but the bias in the former is expected to be smaller than in the latter. Unless the population ratio of y to x in different strata varies considerably, the combined estimator will usually have negligible bias and precision will be as high as the separate estimator. In cases when the line of regression of y on x passes through the origin within

each stratum, the separate ratio estimator will be more precise than the combined one. Hence the rules for choice between the two methods should be

1. If the sample taken from each stratum is small the combined ratio estimator should be used unless there is a wide difference between the strata ratios R_m . When there are such difference and regrouping of the strata is possible so that each group will not differ much and has a large sample size, the separate ratio estimator should be used.
2. If the sample size in each stratum is large so that the approximate variance formula can be applied, then it is better to use the separate ratio estimator unless it involves extra calculation work. If it does, then one should make sure of the gain before actually using it.
3. If X_m are known independently from stratum to stratum the separate ratio estimator may be used as a method of estimation. When $X_m = N_m$, ~~are known independently~~ ^{the estimator becomes} equivalent to the stratified estimator and there is no problem of bias from these estimates.
4. If the sampling units are the elementary units and the denominator of the ratio is simply the number of elementary units in the sample, both the estimators are unbiased. In such cases, the estimator should be chosen by consideration of the variance and gain to be obtained by using it.

Regression Estimators:-

Linear regression estimators also make use of auxiliary information for increasing precision. It was seen that the ratio estimator provides a precise estimate of the population mean of regression. As linear and the line does not go through the origin, it is better to use estimator based on linear regression. In other words, if the study variable (y) is approximately a constant and a multiple of the auxiliary variate, it is more precise to estimate the population mean or total by fitting a linear regression. Such an estimator is called a Regression Estimator.

Difference Estimator:-

Let y and x be correlated characteristics, to estimate \bar{Y} , if from a SRS, we obtain the unbiased estimator \bar{y} and \bar{x} of \bar{Y} and \bar{x} respectively, then we can improve upon the estimator \bar{y} by introducing a difference function. Thus a simple difference estimator is defined by,

$$\hat{y}_D = \bar{y} + (\bar{x} - \bar{x})$$

It is assumed that there is a unit change in y when a unit change is made in x . x and y having equal variances. This assumption not be valid if the relationship is of the type $y = cx + k$, where c & k are constants. In Difference Estimator is defined as,

$$\hat{y}_D = \bar{y} + c(\bar{x} - \bar{x})$$

c & k \Rightarrow one known Quantities.

Theorem: In SRS, we have, The sampling variance of \bar{y}_D is obtained by.

$$V(\bar{y}_D) = \frac{1-f}{n} [s_y^2 + c^2 s_x^2 - 2fc s_x s_y]$$

where c is some specified quantity.

Proof:

Since c is a given constant and $E(\bar{x} - \bar{\pi}) = 0$ the unbiasedness of \bar{y}_D is proved.

$$u_i = y_i - c(x_i - \bar{x})$$

$$\bar{u} = \bar{y} - c(\bar{x} - \bar{x})$$

$$\bar{U} = \bar{y} - c(\bar{x} - \bar{x}) = \bar{y}$$

$$V(\bar{y}_D) = V(\bar{u}) = \frac{1-f}{n} s_u^2$$

$$\text{where } s_u^2 = \frac{\sum_{i=1}^N (u_i - \bar{U})^2}{(N-1)}$$

$$= \frac{\sum_{i=1}^N [(\bar{y} - c(x_i - \bar{x})) - \bar{y}]^2}{N-1}$$

$$= \frac{\sum_{i=1}^N [(\bar{y} - \bar{y}) - c(\bar{x} - \bar{x})]^2}{N-1}$$

$$= \frac{1}{N-1} \left[\sum_{i=1}^N (\bar{y} - \bar{y})^2 + c^2 \sum_{i=1}^N (\bar{x} - \bar{x})^2 - 2c \sum_{i=1}^N (\bar{x} - \bar{x})(\bar{y} - \bar{y}) \right]$$

$$= \sum_{i=1}^N \frac{(y_i - \bar{y})^2}{N-1} + c^2 \sum_{i=1}^N \frac{(x_i - \bar{x})^2}{N-1} - 2c \rho S_n S_y$$

$$S_n^2 = S_y^2 + c^2 S_x^2 - 2c \rho S_n S_y$$

$$S_n^2 = S_y^2 + c^2 S_x^2 - 2c \rho S_n S_y$$

$$V(\hat{Y}_D) = \frac{1-f}{n} [S_y^2 + c^2 S_x^2 - 2c \rho S_n S_y]$$

Hence Proved.

Corollary:-

- For the case $c = R (= Y_{1X})$, the variance of the difference estimator \bar{Y}_D is exactly the same as the first order approximation of $V(\hat{Y}_R)$

$$V(\hat{Y}_R) = \frac{N^2(1-f)}{n} [S_y^2 + R^2 S_x^2 - 2R \rho S_n S_y]$$

$$= \frac{(1-f) \Delta Y^2}{n} [c_y^2 + c_x^2 - 2f c_n c_y]$$

- In SRS, an unbiased estimator of $V(\bar{Y}_D)$ is obtained by,

$$V(\bar{Y}_D) = V(\bar{Y}) + c^2 V(\bar{x}) - 2c \text{cov}(\bar{Y}, \bar{x})$$

$$= \frac{1-f}{n} (S_y^2 + c^2 S_x^2 - 2c S_n S_y)$$

Regression Estimator:

While discussing the difference estimator, it was seen that the optimum value to be given to c & β , where β is the regression coefficient of y and x . Generally, β is not known in advance and its value is estimated from the sample. Suppose y_i and x_i are obtained for each unit in the sample, then a least square estimate of β is,

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

Thus Linear regression estimator of the population mean \bar{Y} , one population mean \bar{Y} are given by,

$$\hat{Y}_1 = \bar{Y} + b(\bar{X} - \bar{x})$$

$$\hat{Y}_2 = N[\bar{Y} + b(\bar{X} - \bar{x})]$$

$\hat{y}_i = N(\bar{y}, \sigma^2)$

Therefore, we shall present the theory of linear regression estimator \hat{y}_i only. Since b is random variate exact expression of the mean and variance of the regression estimator are difficult to derive. Large sample approximation to its bias and sampling variance will be given.

Bias of Regression Estimator:-

The regression estimator \hat{y}_i is biased since,

- β is generally estimated by taking the ratio of the estimate $\text{cov}(\hat{y}, \bar{x})$ to that $\text{v}(\bar{x})$.
- It involves the product of 2 variables viz. $b\bar{x}$. The bias regression estimator will usually be trivial and will decrease as sample size increase.

Theorem: In SRS, the bias of \hat{y}_i is approximately by. $B(\hat{y}_i) \cong -\text{cov}(\bar{x}, b)$

which will be negligible if the sample size is ~~large~~ large.

Proof:-

$$\text{Suppose } \bar{y} = \bar{Y}(1+e)$$

$$\bar{x} = \bar{X}(1+e_1)$$

$$b = \beta(1+e_2)$$

Where

$$E(e) = E(e_1) = E(e_2) = 0.$$

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

$$\begin{aligned}
 \bar{y}_1 &= \bar{y}(1+e) + \beta(1+e_2)(\bar{x} - \bar{x}(1+e_1)) \\
 &= \bar{y} + e\bar{y} + (\beta + \beta e_2)(\bar{x} - \bar{x} + e_1\bar{x}) \\
 &= \bar{y} + e\bar{y} + (\bar{x}\beta - \bar{x}\beta + \beta e_1\bar{x} + \beta e_2\bar{x} - \beta e_1 e_2 \bar{x}) \\
 &= \bar{y} + e\bar{y} - \beta e_1 \bar{x} - \beta e_1 e_2 \bar{x} \\
 \bar{y}_1 &= \bar{y} + (e\bar{y} - e_1 \beta \bar{x}) - e_1 e_2 \beta \bar{x}
 \end{aligned}$$

\therefore The bias of regression estimator is

$$B(\bar{y}_1) = \bar{y} E(e) + \beta \bar{x} E(e_1) - \beta \bar{x} E(e_1 e_2)$$

$$B(\bar{y}_1) = - \text{cov}(\bar{x}, b) \quad [\because E(e) = E(e_1) = 0.]$$

For large samples, usually $\text{cov}(\bar{x}, b)$ decreases.

It becomes zero if the joint distribution of y and x is a bivariate normal.

Sampling Variance of Regression Estimator:

Theorem: In SRS, the large sample variance of the regression estimator is given by

$$V(\bar{y}_1) = V(\bar{y}) + \beta^2 V(\bar{x}) - 2\beta \text{cov}(\bar{y}, \bar{x}).$$

PP:

$$\begin{aligned}
 \bar{y}_1 &= \bar{y} + \cancel{e\bar{y}} (e\bar{y} - e_1 \beta \bar{x}) + e_1 e_2 \beta \bar{x} \\
 \bar{y}_1 - \bar{y} &= e\bar{y} - e_1 \beta \bar{x} - e_1 e_2 \beta \bar{x}
 \end{aligned}$$

If the term involving $e_1 e_2$ is ignoring, we have

$$y_1 - \bar{y} = e\bar{y} - e_1 \beta \bar{x}$$

Therefore,

$$V(\bar{y}_i) = V(e\bar{Y} - e_i \beta \bar{x})$$

$$= V(e\bar{Y}) + \beta^2 V(e_i \bar{x}) - 2\beta \text{cov}(e\bar{Y}, e_i \bar{x})$$

$$V(\bar{y}_i) = V(\bar{y}) + \beta^2 V(\bar{x}) - 2\beta \text{cov}(\bar{y}, \bar{x}).$$

Combined with the mean per unit and Ratio estimators

for a large numbers sample size n , the variance
for minimum estimators of the population mean \bar{Y} ,

$$V(\bar{y}) = \frac{1-f}{n} s_y^2 \text{ (mean per unit)}$$

$$V(\bar{y}_R) = \frac{1-f}{n} [s_y^2 - 2Rf s_n s_y + R^2 s_n^2] \text{ (ratio).}$$

$$V(\bar{Y}_r) = \frac{1-f}{n} (1-f^2) s_y^2 \text{ (regression).}$$

Finally, comparing with the mean per unit estimator,
we observe that the variance of the regression estimator
always smaller unless $f=0$ in case $f=0$ the variation
for both are equal. The reduction in variance is large
when f is high and small when f is low.

Comparing next with the ratio estimator, it
can be seen that the variances of the regression
estimator is less than that of the ratio estimator if

$$\cancel{\text{if}} (f s_y - R s_x)^2 > 0$$

$$(\beta - R)^2 > 0$$

which is always true unless $\beta = R$. In this
situation, both estimators have the same variance
and this occurs only when the regression of
 y on x is a straight line through the origin.

The comparison suggests that all these estimators belong to the class of estimators.

$$\hat{y}_R = \bar{y} + z(\bar{x} - \hat{x}),$$

where z is a random variate having some values in a finite range.

Some rules for a choice among alternatives are as follows:

1. When advance information on an approximate value of β ($= c$) is available, then with simple computations good results can be obtained with such value of c with the difference estimator.
2. When the correlation coefficient between variables is nearly equal to the ratio of their S.D i.e., $\rho \approx \frac{\sigma_x}{\sigma_y}$, the difference estimator with $c=1$ will obtain equally precise results as the regression estimator.
3. When the correlation coefficient between variates is nearly equal to the ratio of their coefficients of variation i.e., $\rho \approx \frac{c_x}{c_y}$, the ratio estimator will provide equally precise results as the regression estimator.
4. When ρ is different from $\frac{\sigma_x}{\sigma_y}$ the regression estimator should be preferred. In this situation, the difference estimator with $c=1$ will not provide precise results.

5. When f is different form Cx/cy , the regression estimator should be preferred. In this situation, the variance of the ratio estimator will be larger than that of regression estimator.

6. When computations for the regression estimator are heavy, time consuming and expensive, its uses is recommended only if the gains from such computations are such more significant than the additional cases.

Regression Estimators in stratified Sampling:-

Nm units * SRSWOR like the ratio estimator, too regression are possible Separate regression Estimator.

* The combined regression Estimator which is obtained by getting common regression coefficient of the strata.

Separate Regression Estimator:-

Theorem:- If sampling is independent in different strata and sample size is large enough in each stratum, then \bar{Y}_{ls} is ~~bias~~, is an almost unbiased estimator. To the 1st order of approximation its variance is given by,

$$V(\bar{Y}_{ls}) = \sum_m w_m \frac{(1-f_m)}{n_m} [S_{ym}^2 - 2b_m p_m S_{ym} S_{cm} + b_m^2 S_{cm}^2]$$

Proof

$$u_i = y_i - b(x_i - \bar{x}), \bar{u} = \bar{y} - b(\bar{x} - \bar{x})$$

$$U = \bar{y} - b(\bar{x} - \bar{x})$$

$$V(\bar{y}_{lb}) = V(\bar{u}) = \frac{1-\frac{1}{n}}{n} s_u^2$$

$$s_u^2 = \frac{\sum_{i=1}^N (u_i - \bar{u})^2}{N-1}$$

$$= \frac{\sum_{i=1}^N [(\bar{y} - \bar{y})^2 - b(\bar{x} - \bar{x})^2]}{N-1}$$

$$= \frac{1}{N-1} \left[\sum_{i=1}^N (\bar{y} - \bar{y})^2 + b^2 \sum_{i=1}^N (\bar{x} - \bar{x})^2 - \right.$$

$$\left. - 2b \sum_{i=1}^N (\bar{x} - \bar{x})(\bar{y} - \bar{y}) \right]$$

$$= \left[\sum_{i=1}^N \frac{(\bar{y} - \bar{y})^2}{N-1} + b^2 \sum_{i=1}^N \frac{(\bar{x} - \bar{x})^2}{N-1} \right.$$

$$\left. - 2b \frac{\sum (\bar{x} - \bar{x})(\bar{y} - \bar{y})}{N-1} \right]$$

$$s_u^2 = [s_{y_m}^2 + b^2 s_{x_m}^2 - 2b \rho s_x s_y]$$

$$V(\bar{y}_{lb}) = \sum_{m=1}^k w_m^2 \frac{(1-p_m)}{n_m} \left[s_{y_m}^2 + b_m^2 s_{x_m}^2 - 2b_m p_m s_{x_m} s_{y_m} \right].$$

Thus Proved.

Combined Regression Estimator:-

Theorem:- If Sampling is independent in different strata and sample size is large enough in each stratum, the variance of \bar{Y}_{lc} is given by,

$$V(\bar{Y}_{lc}) = \sum_m^k W_m^2 \frac{(1-f_m)}{n_m} (S_{ym}^2 - 2b f_m S_{ym} S_{nm} + b^2 S_{nm}^2)$$

[The proof is previous theorem]

Comparision of Combined and Separate Estimators:-

Some caution is required when deciding whether the separate or combined regression estimator is to be used at any specific situation. The bias in the separate regression estimator is likely to be larger when sample size in each stratum is small. On the other hand, the defect of the combined estimator is that its variance is inflated, if the population regression coeff. differ from stratum to stratum. If the regression are approximately linear and the regression coeff. appear to be the same in all strata, use of the combined regression estimator is preferred. If the regression are linear but the coeff. differ from stratum to stratum, it is advisable to adopt the separate estimator. Though the separate regression estimator is likely to be more efficient than the combined regression estimator, the latter is likely to have smaller bias.