**ILORIN
JOURNAL OF SCEINCE**

**ILJS-15-048**

# Efficiency in Double Sampling under a Particular Linear Regression Model

**Adewara*, A. A.**

Department of Statistics, University of Ilorin, Ilorin, Nigeria.

## Abstract

In this study, three regression and mean per unit estimators in double sampling under a particular linear regression model were proposed using Monte Carlo Method through simulation and from the estimated mean square errors obtained on these estimators, we observed that one of the newly proposed regression estimator, $\overline{y}_{dl2}$, performed better and hence, preferred.

**Keywords:** Double sampling, Estimators, Linear Regression, Monte Carlo Method, Simulation

## 1.     Introduction

Monte Carlo technique according to Kendall and Buckland (1982) is a way of obtaining the solution of mathematical problems in a stochastic context through sampling experiments. It is the solution of any mathematical problem by sample methods. The procedure is to conduct an artificial stochastic model of the mathematical process and then to perform sampling experiments upon it. It is also used to solve physical problems and mostly useful when it is difficult or impossible to use other approaches. It is a broad class of computational algorithm that relies on repeated random sampling to obtain numerical results. The essential idea is using randomness to solve problems that might be deterministic in principle. Mainly used in three distinct problem classes: generating draws from a probability distribution, optimization and numerical integration. In principle, Monte Carlo technique can be used to solve any problem having a probabilistic interpretation.

This paper consists of four sections: introduction, materials and methods, results and discussion and conclusion.

*Corresponding Author:  Adewara, A.A.
Email: aaadewara@gmail.com, aadewara@unilorin.edu.ng

## 2.      Materials and Methods

**On the model used:** The linear regression model used in this study is of the form: $y_i = \beta x_i + \mu_i$ , $i = 1,2,..., N$   where   $x_i$   follows   a   gamma   distribution,   $\mu_i$ is   usually distributed with mean 0 and variance $(\mu/x) = \sigma^2 x_i^{\ t}$ and t is said to be the variance function usually between 0 and 2. The gamma distribution with parameter $\alpha, \beta$ and probability density   function   $\dfrac{\beta^\alpha}{\Gamma(\alpha)} X^{\alpha-1} e^{-\beta x}$   was   used,   $\alpha > 0, \beta > 0, E(X) = \dfrac{\alpha}{\beta}, E(Y) = \alpha$ and $Var(X) = \dfrac{\alpha}{\beta^2}$ where, $\alpha$ and $\beta$ are defined as a "scale" and "shape" parameter, respectively while E(X) and E(Y) are the means of X and Y populations.

In sample survey, gamma distribution is a distribution that has skewed population in which a small proportion of the sampling units in the population may account for a high proportion of an aggregate or average being measure. Examples of these skewed populations are sample of farms by size, distribution of sales of retail trade in establishment and family income distribution of the population.

Here, both the sigma $(\sigma)$, alpha $(\alpha)$, beta $(\beta)$ and t values are specified, where t controls the dependence of the variance of y on the value of x. it is always regarded as a non-negative constant. Also, statistical transformation was done from a normal distribution to uniform distribution which made $x_i$ to be distributed as gamma since the normal distribution is the most important distribution in statistics which is likely to provide a good model for a variate when:

- There is a strong tendency for the variate to take a central value.
- Positive and negative deviations from the central value are equally likely.
- The frequency of deviations falls off rapidly as the deviations becomes larger Cooke *et al.* (1982).

The following estimators are considered for comparisons:-

- $\bar{y}$, mean per unit
- $\bar{y}_{dlo} = \bar{y} - \beta(\bar{x} - \bar{x}')$
- $\bar{y}_{dl1} = \bar{y}^* - \beta(\bar{x}^* - \bar{x}')$ and

- $\bar{y}_{dl2} = \bar{y}^{**} - \beta(\bar{x}^{**} - \bar{x}')$

**Mean square error of the conventional estimator, $\bar{y}_{dl0}$ :**

Let N be the population, $n'$ and $n$ be the first sample and subsample drawn, $\bar{x}'$ and $\bar{x}$ be the sample mean of the first sample and subsample drawn, respectively.

$\bar{x} = \bar{X}(1+\Delta_{\bar{x}})$ and $\bar{y} = \bar{Y}(1+\Delta_{\bar{y}})$. Then, $Mse(\bar{y}) = (\frac{1}{n} - \frac{1}{N})s^2{}_y \approx (\frac{1}{n})s^2{}_y$, as

$N \to \infty, \frac{1}{N} \to 0$.

$\bar{y}_{dlo} = \bar{y} - \beta(\bar{x} - \bar{x}')$, then,

$Mse(\bar{y}_{dl0}) \approx v(\bar{y}_{dl0}) = v[\bar{y} - \beta(\bar{x} - \bar{x}')]$

$$= v(\bar{y}) + \beta^2(\bar{x} - \bar{x}') - 2\beta\,\mathrm{cov}(\bar{y}, \bar{x} - \bar{x}')$$

$$= v(\bar{y}) + \beta^2[v(\bar{x}) + v(\bar{x}') - 2\,\mathrm{cov}(\bar{x} - \bar{x}')] - 2\beta[\mathrm{cov}(\bar{x}, \bar{y}) - \mathrm{cov}(\bar{x}', \bar{y})].$$

Using the unconditional expectation $E_1$ and conditional expectation $E_2$ over the first sample and the second subsample given the first sample, unconditional variance $V_1$ and conditional variance $V_2$ of the first sample and the second subsample given the first sample, then

$v(\bar{y}) = (\frac{1}{n})s^2{}_y, v(\bar{x}) = (\frac{1}{n})s^2{}_x, v(\bar{x}') = (\frac{1}{n'})s^2{}_{x'}, \mathrm{cov}(\bar{x}, \bar{x}') = (\frac{1}{n'})s^2{}_x, \mathrm{cov}(\bar{x}, \bar{y}) = (\frac{1}{n})s_{xy}$ and

$\mathrm{cov}(\bar{x}', \bar{y}) = (\frac{1}{n'})s_{x'y}$. Then,

$Mse(\bar{y}_{dl0}) = (\frac{1}{n'} - \frac{1}{N})s^2{}_y + (\frac{1}{n} - \frac{1}{n'})(s^2{}_y - 2\beta s_{xy} + \beta^2 s^2{}_x)$ ( Okafor, 2002; Xuejun *et al,.* 2011;

Agunbiade and Ogunyinka, 2013; Ogunyinka and Sodipo, 2013; Yi-Hau and Hung, 2000) as

$N \to \infty, \frac{1}{N} \to 0$, then, $Mse(\bar{y}_{dl0}) \approx (\frac{1}{n'})s^2{}_y + (\frac{n'-n}{nn'})(s^2{}_y - 2\beta s_{xy} + \beta^2 s^2{}_x) \dots$  (1)

**Mean square error of the proposed estimator, $\bar{y}_{dl1}$ :**

Let $\bar{x}^*$ and $\bar{y}^*$ be the means of the auxiliary variable(x) and for the variable of interest (y) yet to be drawn when the first sample is drawn (Srivenkataramana and Srinath, 1976), that is, the means corresponding to the ($N - n'$) population units and $\rho$ be the coefficient of correlation between X and Y (Srivenkataramana and Srinath, 1976; Adewara, 2005; 2006).

Let, $\bar{y}_{dl1} = \bar{y} - \beta(\bar{x}^* - \bar{x}')$, where

$$\bar{X} = \theta\bar{x} + (1-\theta)\bar{x}^*, \quad \bar{Y} = \theta\bar{y} + (1-\theta)\bar{y}^*, \quad \theta = \frac{1}{n'} - \frac{1}{N} \text{ as } N \to \infty, \frac{1}{N} \to 0 \text{, then,}$$

$$\bar{x}^* = \bar{X}(1 - (\frac{1}{(n'-1)})\Delta_{\bar{x}^*}), \quad \bar{y}^* = \bar{Y}(1 - (\frac{1}{(n'-1)})\Delta_{\bar{y}^*}), \quad |\Delta_{\dot{x}*}| = \left|\frac{\bar{x} - \bar{x}'}{\bar{x}'}\right| < 1 \text{ and}$$

$$|\Delta_{\bar{y}*}| = \left|\frac{\bar{y} - \bar{y}'}{\bar{y}'}\right| < 1.$$

Then

$$Mse(\bar{y}_{dl1}) \approx v(\bar{y}_{dl1}) = v[\bar{y}^* - \beta(\bar{x}^* - \bar{x}')]$$

$$= v(\bar{y}^*) + \beta^2(\bar{x}^* - \bar{x}') - 2\beta\,cov(\bar{y}^*, \bar{x}^* - \bar{x}')$$

$$= v(\bar{y}^*) + \beta^2[v(\bar{x}^*) + v(\bar{x}') - 2cov(\bar{x}^* - \bar{x}')] - 2\beta[cov(\bar{x}^*, \bar{y}^*) - cov(\bar{x}', \bar{y}^*)].$$

$$v(\bar{y}^*) = \frac{1}{n}(\frac{1}{n'-1})^2 s^2_y, v(\bar{x}^*) = \frac{1}{n}(\frac{1}{n'-1})^2 s^2_x, v(\bar{x}') = (\frac{1}{n'})s^2_{x'}, cov(\bar{x}^*, \bar{x}') = \frac{1}{n'}(\frac{1}{n'-1})s^2_x$$

$$cov(\bar{x}^*, \bar{y}^*) = \frac{1}{n}(\frac{1}{n'-1})^2 s_{xy} \text{ and } cov(\bar{x}', \bar{y}^*) = \frac{1}{n'}(\frac{1}{n'-1})s_{xy}.$$

Then

$$Mse(\bar{y}_{dl1}) \approx (\frac{1}{n'})(\frac{1}{n'-1})^2 s^2_y + (\frac{n'-n}{n'n})(\frac{1}{n'-1})^2(s^2_y - 2\beta s_{xy} + \beta^2 s^2_x)\dots \qquad (2)$$

**Mean square error of the proposed estimator, $\bar{y}_{dl2}$:**

Let $\bar{x}^{**}$ and $\bar{y}^{**}$ be the means of the auxiliary variable(x) and for the variable of interest (y) yet to be drawn when the second sample is drawn (Srivenkataramana and Srinath, 1976), that is, the means corresponding to the $(n'-n)$ population units and $\rho$ be the coefficient of correlation between X and Y (Srivenkataramana and Srinath, 1976).

Let,

$$\bar{X} = (1+\theta)\bar{x}^{**} - \theta\bar{x}, \quad \bar{Y} = (1+\theta)\bar{y}^{**} - \theta\bar{y}, \quad \theta = \frac{1}{n'} - \frac{1}{N} \text{ as } N \to \infty, \frac{1}{N} \to 0 \text{, then,}$$

$$\theta = \frac{1}{n'}, \bar{x}^{**} = \bar{X}(1 + (\frac{1}{(n'+1)})\Delta_{\bar{x}^{**}}) \text{ and } \bar{y}^{**} = \bar{Y}(1 + (\frac{1}{(n'+1)})\Delta_{\bar{y}^{**}}).$$

$$\bar{y}_{dl2} = \bar{y}^{**} - \beta(\bar{x}^{**} - \bar{x}'), \text{ then,}$$

$$Mse(\bar{y}_{dl2}) \approx v(\bar{y}_{dl2}) = v[\bar{y}^{**} - \beta(\bar{x}^{**} - \bar{x}')]$$

$$= v(\bar{y}^{**}) + \beta^2(\bar{x}^{**} - \bar{x}') - 2\beta\,cov(\bar{y}^{**}, \bar{x}^{**} - \bar{x}')$$

$$= v(\bar{y}^{**}) + \beta^2[v(\bar{x}^{**}) + v(\bar{x}') - 2cov(\bar{x}^{**} - \bar{x}')] - 2\beta[cov(\bar{x}^{**}, \bar{y}^{**}) - cov(\bar{x}', \bar{y}^{**})],$$

where

$$v(\bar{y}^{**}) = \frac{1}{n}(\frac{1}{n'+1})^2 s^2_y, v(\bar{x}^{**}) = \frac{1}{n}(\frac{1}{n'+1})^2 s^2_x, v(\bar{x}') = (\frac{1}{n'})s^2_{x'}, cov(\bar{x}^{**}, \bar{x}') = \frac{1}{n'}(\frac{1}{n'+1})s^2_x$$

$$cov(\bar{x}^{**}, \bar{y}^{**}) = \frac{1}{n}(\frac{1}{n'+1})^2 s_{xy} \text{ and } cov(\bar{x}', \bar{y}^{**}) = \frac{1}{n'}(\frac{1}{n'+1})s_{xy}.$$

Then

$$Mse(\bar{y}_{dl2}) \approx (\frac{1}{n'})(\frac{1}{n'+1})^2 s^2_y + (\frac{n'-n}{n'n})(\frac{1}{n'+1})^2 (s^2_y - 2\beta s_{xy} + \beta^2 s^2_x) \dots . \qquad (3)$$

Therefore

$$Mse(\bar{y}_{dl0}) \approx (\frac{1}{n'})s^2_y + (\frac{n'-n}{nn'})(s^2_y - 2\beta s_{xy} + \beta^2 s^2_x)$$

$$Mse(\bar{y}_{dl1}) \approx (\frac{1}{n'})(\frac{1}{n'-1})^2 s^2_y + (\frac{n'-n}{n'n})(\frac{1}{n'-1})^2 (s^2_y - 2\beta s_{xy} + \beta^2 s^2_x) \text{ and}$$

$$Mse(\bar{y}_{dl2}) \approx (\frac{1}{n'})(\frac{1}{n'+1})^2 s^2_y + (\frac{n'-n}{n'n})(\frac{1}{n'+1})^2 (s^2_y - 2\beta s_{xy} + \beta^2 s^2_x).$$

Here, the coefficient for $(\frac{1}{n'})s^2_y$ and $(\frac{n'-n}{nn'})(s^2_y - 2\beta s_{xy} + \beta^2 s^2_x)$ are:

- +1 and +1 for $\bar{y}_{dl0}$,

- $(\frac{1}{n'-1})^2$ and $(\frac{1}{n'-1})^2$ for $\bar{y}_{dl1}$,

- $(\frac{1}{n'+1})^2$ and $(\frac{1}{n'+1})^2$ for $\bar{y}_{dl2}$.

One could see that their difference lies on these coefficients. Under what condition would an estimator of this kind be preferred? An estimator of this kind among others being considered here would be preferred if it has least estimated mean square error.


3.      **Results and Discussion**

Using the gamma distribution described above, where $n' = 140, n = 2, 20, 40, 80$ and $100$, $\alpha = 1, 2, 3$ and $\beta = 1, 2$, and 3. Then, the estimates obtained using these information are shown in Tables 1 – 3 below.

**Tab.1:** Estimated mean square error of $\bar{y}, \bar{y}_{dlo}, \bar{y}_{dl1}$ and $\bar{y}_{dl2}$ when $n' = 140, n = 2,20,40,80$ and 100, $\alpha = 1$ and $\beta = 1,2$ and 3.

| $n'$ | $n$ | $\alpha$ | $\beta$ | $mse(\bar{y})$ | $mse(\bar{y}_{dl0})$ | $mse(\bar{y}_{dl1})$ | $mse(\bar{y}_{dl2})$ |
|------|-----|----------|---------|----------------|----------------------|----------------------|----------------------|
| 140 | 2 | 1 | 1 | 1.8242 | 1.8552 | $9.6x10^{-5}$ | $9.3x10^{-5}$ |
|     |   |   | 2 | 27.5348 | 27.8381 | 0.00144 | 0.00140 |
|     |   |   | 3 | 123.502 | 123.805 | $6.41x10^{-3}$ | $6.23x10^{-3}$ |
| 140 | 20 | 1 | 1 | 0.0333 | 0.0419 | $2.2x10^{-6}$ | $2.1x10^{-6}$ |
|     |   |   | 2 | 0.5649 | 0.5743 | $2.97x10^{-5}$ | $2.89x10^{-5}$ |
|     |   |   | 3 | 2.9739 | 2.9833 | $1.54x10^{-4}$ | $1.50x10^{-4}$ |
| 140 | 40 | 1 | 1 | 0.0165 | 0.0220 | $1.14x10^{-6}$ | $1.11x10^{-6}$ |
|     |   |   | 2 | 0.2804 | 0.2869 | $1.49x10^{-5}$ | $1.44x10^{-5}$ |
|     |   |   | 3 | 1.4759 | 1.4824 | $7.7x10^{-5}$ | $7.5x10^{-5}$ |
| 140 | 80 | 1 | 1 | 0.0108 | 0.0118 | $6.1x10^{-6}$ | $5.9x10^{-6}$ |
|     |   |   | 2 | 0.1460 | 0.1472 | $7.6x10^{-6}$ | $7.4x10^{-6}$ |
|     |   |   | 3 | 0.7181 | 0.7193 | $3.7x10^{-5}$ | $3.6x10^{-5}$ |
| 140 | 100 | 1 | 1 | 31.3132 | 34.1866 | $1.77x10^{-3}$ | $1.72x10^{-3}$ |
|     |   |   | 2 | 0.1200 | 0.1209 | $6.26x10^{-6}$ | $6.08x10^{-6}$ |
|     |   |   | 3 | 0.6009 | 0.6015 | $3.11x10^{-5}$ | $3.03x10^{-5}$ |

Here, $mse(\bar{y}_{dl2}) < mse(\bar{y}_{dl1}) < mse(\bar{y}) < mse(\bar{y}_{dl0})$, $\bar{y}_{dl2}$ has the least estimated mean square error, hence preferred.

**Tab. 2:** Estimated mean square error of $\bar{y}, \bar{y}_{dlo}, \bar{y}_{dl1}$ and $\bar{y}_{dl2}$ when $n' = 140, n = 2,20,40,80$ and 100, $\alpha = 2$ and $\beta = 1,2$ and 3.

| $n'$ | $n$ | $\alpha$ | $\beta$ | $mse(\bar{y})$ | $mse(\bar{y}_{dl0})$ | $mse(\bar{y}_{dl1})$ | $mse(\bar{y}_{dl2})$ |
|------|-----|----------|---------|----------------|----------------------|----------------------|----------------------|
| 140 | 2 | 1 | 1 | 1.0562 | 1.2387 | $6.4x10^{-5}$ | $6.23x10^{-5}$ |
|     |   |   | 2 | 7.9533 | 8.1357 | 0.00044 | 0.000409 |
|     |   |   | 3 | 33.7278 | 33.9102 | $1.76x10^{-3}$ | $1.71x10^{-3}$ |
| 140 | 20 | 1 | 1 | 0.0703 | 0.0849 | $4.39x10^{-6}$ | $4.23x10^{-6}$ |
|     |   |   | 2 | 1.3694 | 1.3854 | $7.17x10^{-5}$ | $6.97x10^{-5}$ |
|     |   |   | 3 | 7.1739 | 7.1894 | $3.72x10^{-4}$ | $3.62x10^{-4}$ |
| 140 | 40 | 1 | 1 | 1.7629 | 2.2500 | $1.2x10^{-6}$ | $1.1x10^{-6}$ |
|     |   |   | 2 | 0.2495 | 0.2543 | $1.32x10^{-5}$ | $1.28x10^{-5}$ |
|     |   |   | 3 | 1.2757 | 1.2807 | $6.63x10^{-5}$ | $6.44x10^{-5}$ |
| 140 | 80 | 1 | 1 | 0.0248 | 0.0268 | $1.38x10^{-6}$ | $1.34x10^{-6}$ |
|     |   |   | 2 | 0.3354 | 0.3373 | $1.75x10^{-6}$ | $1.70x10^{-6}$ |
|     |   |   | 3 | 1.6836 | 1.6856 | $8.7x10^{-5}$ | $8.5x10^{-5}$ |
| 140 | 100 | 1 | 1 | 0.0241 | 0.0247 | $1.28x10^{-3}$ | $1.24x10^{-3}$ |
|     |   |   | 2 | 0.3590 | 0.3597 | $1.86x10^{-6}$ | $1.81x10^{-6}$ |
|     |   |   | 3 | 1.8133 | 1.8140 | $9.4x10^{-5}$ | $9.1x10^{-5}$ |

Here, $mse(\bar{y}_{dl2}) < mse(\bar{y}_{dl1}) < mse(\bar{y}) < mse(\bar{y}_{dl0})$, $\bar{y}_{dl2}$ has the least estimated mean square error, hence preferred.

**Tab.3:** Estimated mean square error of $\bar{y}, \bar{y}_{dlo}, \bar{y}_{dl1}$ and $\bar{y}_{dl2}$ when $n' = 140, n = 2,20,40,80$ and 100, $\alpha = 3$ and $\beta = 1,2$ and 3.

| $n'$ | $n$ | $\alpha$ | $\beta$ | $mse(\bar{y})$ | $mse(\bar{y}_{dl0})$ | $mse(\bar{y}_{dl1})$ | $mse(\bar{y}_{dl2})$ |
|------|-----|----------|---------|----------------|----------------------|----------------------|----------------------|
| 140 | 2 | 1 | 1 | 0.0673 | 0.1045 | $5.4x10^{-6}$ | $5.3x10^{-6}$ |
|     |   |   | 2 | 2.6259 | 2.6631 | $1.38x10^{-6}$ | $1.34x10^{-6}$ |
|     |   |   | 3 | 15.1234 | 15.1605 | $7.8x10^{-4}$ | $7.60x10^{-4}$ |
| 140 | 20 | 1 | 1 | 0.2157 | 0.2313 | $1.7x10^{-5}$ | $1.16x10^{-5}$ |
|     |   |   | 2 | 2.7557 | 2.7724 | $1.43x10^{-4}$ | $1.39x10^{-4}$ |
|     |   |   | 3 | 13.3728 | 13.3895 | 0.00069 | 0.00067 |
| 140 | 40 | 1 | 1 | 0.0570 | 0.0618 | $3.20x10^{-6}$ | $3.11x10^{-6}$ |
|     |   |   | 2 | 0.8651 | 0.8700 | $4.5x10^{-5}$ | $4.38x10^{-5}$ |
|     |   |   | 3 | 4.4034 | 4.4082 | $2.28x10^{-4}$ | $2.22x10^{-4}$ |
| 140 | 80 | 1 | 1 | 0.0514 | 0.0526 | $2.72x10^{-6}$ | $2.65x10^{-6}$ |
|     |   |   | 2 | 0.7721 | 0.7733 | $4.0x10^{-5}$ | $3.89x10^{-5}$ |
|     |   |   | 3 | 3.8821 | 3.8833 | $2.01x10^{-4}$ | $1.95x10^{-4}$ |
| 140 | 100 | 1 | 1 | 0.0245 | 0.0252 | $1.3x10^{-6}$ | $1.27x10^{-6}$ |
|     |   |   | 2 | 0.3563 | 0.3570 | $1.85x10^{-5}$ | $1.80x10^{-5}$ |
|     |   |   | 3 | 1.7941 | 1.7947 | $9.29x10^{-5}$ | $9.03x10^{-5}$ |

Here, $mse(\bar{y}_{dl2}) < mse(\bar{y}_{dl1}) < mse(\bar{y}) < mse(\bar{y}_{dl0})$, $\bar{y}_{dl2}$ has the least estimated mean square error, hence preferred.

## 4.      Conclusion

From Tables 1 -3, irrespective of the values of $\alpha$ and $\beta$, the unique finding is that double sampling regression estimator, $\bar{y}_{dl2}$ has the least estimated mean square error, $mse(\bar{y}_{dl2}) < mse(\bar{y}_{dl1}) < mse(\bar{y}) < mse(\bar{y}_{dl0})$. Hence, double sampling regression estimator, $\bar{y}_{dl2}$ is preferred.

**References**

Adewara, A.A. (2005): *Alternative to Ratio and Product Estimators*, Ph.D thesis, University of Ilorin, Nigeria. (Unpublished).

Adewara, A.A. (2006): Effect of improving both the auxiliary and variable of interest in ratio and product estimators. *Proc. Pakistan Acad. Sci.* **43**(4), 275 – 278

Agunbiade, D.A. and Ogunyinka, P.I. (2013): Effects of Correlation Level on the Use of Auxiliary Variable in Double Sampling for Regression Estimation. *Open J. Stat.* **3**, 312 – 318.

Cooke, D., Craven, A.H. and Clarke, G.M. (1982): *Basic Statistical Computing*. Edward Arnoid.

Kendall, M.G. and Buckland, W.R. (1982): *A dictionary of Statistical terms.* 4th Edn. Longman Group Ltd, New York.

Ogunyinka, P. I. and Sodipo, A.A. (2013): Efficiency of Ratio and Regression Estimators Using Double Sampling. *Journal of Natural Sciences Research.* **3**(7), 201 – 208.

Okafor, F.C. (2002): *Sample Survey Theory With Applications*. First Edition, Afro-Orbit Publications Ltd. University of Nigeria, Nsukka, Nigeria.

Srivenkataramana, T. and Srinath, K.P. (1976): Ratio and Product Method of estimation in Sample Surveys when the two variables are moderately correlated. Vignana Bharathi., **2**, 54 – 58.

Xuejun Jiang, Jiancheng Jiang and Yanling Liu (2011): Nonparametric regression under double sampling designs. *Journal of Systems Science and Complexity.* **24**(1), 167 – 175.

Yi-Hau, C. and Hung, C. (2000): A unified approach to regression analysis under double sampling designs. *Journal of Royal Statistical Society.* **62**(3), 443 – 460.