# Estimation of Population Variance Using a Generalized Double Sampling Estimator

## Peeyush Misra[1*] and R. Karan Singh[2]

[1]Department of Statistics, D.A.V.(P.G.) College, Dehradun- 248001, Uttarakhand, India.

[2]Department of Statistics, Lucknow University, Lucknow- 226007 Uttar Pradesh, India.

[*]Corresponding Author: dr.pmisra.dav@gmail.com

## ABSTRACT

*For the estimation of finite population variance a generalized double sampling estimator is proposed. The bias and mean square error (MSE) of the proposed estimator are found. Theoretical comparison with the traditional estimator is carried out and it is shown that the proposed estimator is more efficient than the traditional one.*

## 1. Introduction

In sampling theory, auxiliary information is widely used at the stages of selection and estimation, at the selection stage the auxiliary information is used by designing various sampling schemes and at the estimation stage it is used in formulating various types of estimators of different population parameters with a view of getting increased efficiency. Estimators like ratio, product, difference, regression and the classes of ratio and product type estimators for population parameters mainly population mean and variance are studied by many authors and are available in the literature. To cite some references in this context, one may see estimation procedures and their properties by Das and Tripathi (1978), Liu (1974) and Srivastava and Jhajj (1980). But when parameters of one or more auxiliary variables are not available in advance then the alternative is to use double sampling or two phase sampling technique where we first take a preliminary large sample of

---

size $n'$ (called first phase sample) on which only the auxiliary variable is observed and then from $n'$ taking a sub-sample of size $n$ (called second phase sample) on which both the variables are observed. In such situations the different estimators known as double sampling ratio, product, difference and regression estimators were developed. For more details regarding other double sampling estimation procedures, see Sukhatme et. al. (1984), Murthy (1967) and Cochran (1977). This present paper too contributes to this area.

Let a large preliminary sample of size $n'$ is drawn from a population of size $N$ and then a subsample of size $n$ from $n'$ is drawn by using simple random sampling without replacement scheme for both the phases. At first phase sample of size $n'$, only the auxiliary variable X is observed and at the second phase sample of size $n$, both the study variable Y and the auxiliary variable X are observed.

Let $(\bar{y}, \bar{x})$ be the sample means of $(y, x)$ based on second phase sample of size n and $\bar{x}'$ be the sample mean of first phase $n'$ sample values on the auxiliary character X. Let $\rho$ be the population correlation coefficient between Y and X,

$$S_Y^2 = \frac{1}{N-1}\sum_{i=1}^{N}\left(Y_i - \bar{Y}\right)^2 \ , \quad S_X^2 = \frac{1}{N-1}\sum_{i=1}^{N}\left(X_i - \bar{X}\right)^2 \ \text{ and } \ \mu_{rs} = \frac{1}{N}\sum_{i=1}^{N}\left(Y_i - \bar{Y}\right)^r\left(X_i - \bar{X}\right)^s \ ,$$

where $(Y_i, X_i)$ are the population values of $(Y, X)$ respectively for the $i^{th}$ (i = 1, 2, . . . , N) unit of the population.

For estimating finite population variance, a generalized double sampling estimator is proposed as

$$d_g = \frac{1}{n}\sum_{i=1}^{n} y_i^{\ 2} - \bar{y}\, f(\bar{y}, \bar{x}, \bar{x}'), \tag{1.1}$$

where $f(\bar{y}, \bar{x}, \bar{x}')$ satisfying the validity conditions of Taylor's series expansion is a bounded function of $(\bar{y}, \bar{x}, \bar{x}')$ such that

(i) $$f(\bar{Y}, \bar{X}, \bar{X}) = \bar{Y} \tag{1.2}$$

(ii) first order partial differential coefficient of $f(\bar{y}, \bar{x}, \bar{x}')$ with respect to $\bar{y}$ at $T = \left(\bar{Y}, \bar{X}, \bar{X}\right)$ is unity, that is

$$f_0 = \left(\frac{\partial}{\partial \bar{y}} f(\bar{y}, \bar{x}, \bar{x}')\right)_T = 1 \tag{1.3}$$

(iii) second order partial differential coefficient of $f(\bar{y}, \bar{x}, \bar{x}')$ with respect to $\bar{y}$ at $T = (\bar{Y}, \bar{X}, \bar{X})$ is zero, that is

$$f_{00} = \left( \frac{\partial^2}{\partial \bar{y}^2} f(\bar{y}, \bar{x}, \bar{x}') \right)_T = 0 \tag{1.4}$$

(iv) $$f_1 = -\,f_2, \tag{1.5}$$

for $f_1$ and $f_2$ being the first order partial derivatives of $f(\bar{y}, \bar{x}, \bar{x}')$ with respect to $\bar{x}$ and $\bar{x}'$ respectively at the point $T = (\bar{Y}, \bar{X}, \bar{X})$ and

(v) $$f_{01} = -\,f_{02}\,, \tag{1.6}$$

for $f_{01} = \left( \dfrac{\partial^2}{\partial \bar{y} \partial \bar{x}} f(\bar{y}, \bar{x}, \bar{x}') \right)_T$ and $f_{02} = \left( \dfrac{\partial^2}{\partial \bar{y} \partial \bar{x}'} f(\bar{y}, \bar{x}, \bar{x}') \right)_T$.

## 2. Some Particular Members Belonging to the Proposed Estimator

Some particular members belonging to this proposed generalized double sampling estimator $d_g$ are:

(a) $t_1 = \hat{\theta} - \bar{y}\left\{ \bar{y} \dfrac{\bar{x}}{\bar{x}'} \right\}$

(b) $t_2 = \hat{\theta} - \bar{y}\left\{ \bar{y} \left( \dfrac{\bar{x}}{\bar{x}'} \right)^k \right\}$

(c) $t_3 = \hat{\theta} - \bar{y}\left\{ \bar{y} + k(\bar{x} - \bar{x}') \right\}$

(d) $t_4 = \hat{\theta} - \bar{y}\left\{ \bar{y} \left( \dfrac{\bar{x}}{\bar{x}'} \right)^{k_1} + k_2(\bar{x} - \bar{x}') \right\}$

(e) $t_5 = \hat{\theta} - \bar{y}\left\{ \bar{y} \left( \dfrac{\bar{x}}{\bar{x}'} \right)^{k_1} + k_2\left( \bar{x}^{k_3} - \bar{x}'^{k_3} \right) \right\}$,

where $k$, $k_1$, $k_2$ and $k_3$ are the characterizing scalars to be chosen suitably and $\hat{\theta} = \dfrac{1}{n}\sum_{i=1}^{n} y_i^{\,2}$.

## 3. Bias and Mean Square Error of Proposed Estimator

The proposed generalized double sampling estimator as in equation (1.1) is

$$d_g = \frac{1}{n}\sum_{i=1}^{n} y_i^2 - \bar{y}\ f(\bar{y},\bar{x},\bar{x}')$$

$$= \hat{\theta} - \bar{y}\ f(\bar{y},\bar{x},\bar{x}')\ ,\ \text{where}\ \hat{\theta} = \frac{1}{n}\sum_{i=1}^{n} y_i^2\ .$$

Let $\bar{y} = \bar{Y}+e_0$, $\bar{x} = \bar{X}+e_1$, $\bar{x}' = \bar{X}+e_1'$ and

$$\frac{1}{n}\sum_{i=1}^{n} y_i^2 = \frac{1}{N}\sum_{i=1}^{N} Y_i^2 + e_2$$

or $\quad \dfrac{1}{n}\sum_{i=1}^{n} z_i = \dfrac{1}{N}\sum_{i=1}^{N} Z_i + e_2 \quad,\ \text{where}\ z_i = y_i^2\ \&\ Z_i = Y_i^2$

or $\quad \bar{z} = \bar{Z} + e_2$

or $\quad \hat{\theta} = \theta + e_2\ ,\ \text{where}\ \hat{\theta} = \bar{z}\ \text{and}\ \theta = \bar{Z}\,.$

For simplicity, we assume that the population size $N$ is large enough as compared to the sample size $n$ so that the finite population correction terms may be ignored.

So we have, $\mathrm{E}\,(e_0) = \mathrm{E}\,(e_1) = \mathrm{E}\,(e'_1) = \mathrm{E}\,(e_2) = 0$  (3.1)

$$\left. \begin{aligned} E(e_0^{\,2}) &= \frac{1}{n}S_Y^2 = \frac{\mu_{20}}{n}, & E(e_1^{\,2}) &= \frac{1}{n}S_X^2 = \frac{\mu_{02}}{n}\\[2mm] E(e_1'^{\,2}) &= \frac{\mu_{02}}{n'}, & E(e_0 e_1) &= \frac{1}{n}\rho\ S_Y S_X = \frac{\mu_{11}}{n}\\[2mm] E(e_0 e_1') &= \frac{\mu_{11}}{n'}, & E(e_1 e_1') &= \frac{\mu_{02}}{n'} \end{aligned} \right\}$$

(3.2),

$$E(e_2{}^2) = \frac{1}{n}(\mu_{40} + 4\bar{Y}\mu_{30} + 4\bar{Y}^2\mu_{20} - \mu_{20}{}^2)$$

$$E(e_0 e_2) = \frac{1}{n}(\mu_{30} + 2\bar{Y}\mu_{20})$$

$$E(e_1 e_2) = \frac{1}{n}(\mu_{21} + 2\bar{Y}\mu_{11})$$

$$E(e_1' e_2) = \frac{1}{n'}(\mu_{21} + 2\bar{Y}\mu_{11})$$

(3.3)

Now similar to Singh (1982) and Senapati and Sahoo (2006), expanding (say) $t = f(\bar{y}, \bar{x}, \bar{x}')$ in the third order Taylor's series about the point $T = (\bar{Y}, \bar{X}, \bar{X})$ and noting that $f(\bar{Y}, \bar{X}, \bar{X}) = \bar{Y}$, we have

$$t = f(\bar{Y}, \bar{X}, \bar{X}) + (\bar{y} - \bar{Y})f_0 + (\bar{x} - \bar{X})f_1 + (\bar{x}' - \bar{X})f_2 + \frac{1}{2!}\{(\bar{y} - \bar{Y})^2 f_{00}$$

$$+ (\bar{x} - \bar{X})^2 f_{11} + (\bar{x}' - \bar{X})^2 f_{22} + 2(\bar{y} - \bar{Y})(\bar{x} - \bar{X})f_{01} + 2(\bar{y} - \bar{Y})(\bar{x}' - \bar{X})f_{02}$$

$$+ 2(\bar{x} - \bar{X})(\bar{x}' - \bar{X})f_{12}\} + \frac{1}{3!}\left\{(\bar{y} - \bar{Y})\frac{\partial}{\partial \bar{y}} + (\bar{x} - \bar{X})\frac{\partial}{\partial \bar{x}} + (\bar{x}' - \bar{X})\frac{\partial}{\partial \bar{x}'}\right\}^3 f(\bar{y}* \bar{x}*, \bar{x}'*)$$

On employing the conditions from (1.2) to (1.6), we have

$$t = \bar{Y} + (\bar{y} - \bar{Y}) + (\bar{x} - \bar{X})f_1 + (\bar{x}' - \bar{X})f_2 + \frac{1}{2!}\{(\bar{x} - \bar{X})^2 f_{11} + (\bar{x}' - \bar{X})^2 f_{22}$$

$$+ 2(\bar{y} - \bar{Y})(\bar{x} - \bar{X})f_{01} + 2(\bar{y} - \bar{Y})(\bar{x}' - \bar{X})f_{02} + 2(\bar{x} - \bar{X})(\bar{x}' - \bar{X})f_{12}\}$$

$$+ \frac{1}{3!}\left\{(\bar{y} - \bar{Y})\frac{\partial}{\partial \bar{y}} + (\bar{x} - \bar{X})\frac{\partial}{\partial \bar{x}} + (\bar{x}' - \bar{X})\frac{\partial}{\partial \bar{x}'}\right\}^3 f(\bar{y}*, \bar{x}*, \bar{x}'*) \;,$$

(3.4)

where $f_0$, $f_{00}$, $f_1$, $f_2$ and $f_{02}$ are already defined from (1.3) to (1.6), second order partial derivatives $f_{11}$, $f_{22}$ and $f_{12}$ are given by

$$f_{11} = \left(\frac{\partial^2}{\partial \bar{x}^2} f(\bar{y}, \bar{x}, \bar{x}')\right)_T , \; f_{22} = \left(\frac{\partial^2}{\partial \bar{x}'^2} f(\bar{y}, \bar{x}, \bar{x}')\right)_T , \; f_{12} = \left(\frac{\partial^2}{\partial \bar{x}\partial \bar{x}'} f(\bar{y}, \bar{x}, \bar{x}')\right)_T$$

and

$$\bar{y}* = \bar{Y} + h(\bar{y} + \bar{Y}), \; \bar{x}* = \bar{X} + h(\bar{x} - \bar{X}), \; \bar{x}'* = \bar{X} + h(\bar{x}' - \bar{X}) \;, \text{ for } 0 < h < 1.$$

Now using (3.4) in (1.1) and taking approximation, we have

$$d_g = \frac{1}{n}\sum_{i=1}^{n} y_i^2 - \bar{y}\left[\bar{Y} + \left(\bar{y} - \bar{Y}\right) + \left(\bar{x} - \bar{X}\right)f_1 + \left(\bar{x}' - \bar{X}\right)f_2 + \frac{1}{2!}\left\{\left(\bar{x} - \bar{X}\right)^2 f_{11}\right.\right.$$

$$\left.\left. + \left(\bar{x}' - \bar{X}\right)^2 f_{22} + 2\left(\bar{y} - \bar{Y}\right)\left(\bar{x} - \bar{X}\right)f_{01} + 2\left(\bar{y} - \bar{Y}\right)\left(\bar{x}' - \bar{X}\right)f_{02} + 2\left(\bar{x} - \bar{X}\right)\left(\bar{x}' - \bar{X}\right)f_{12}\right\}\right]$$

$$= (\theta + e_2) - (\bar{Y} + e_0)\left[(\bar{Y} + e_0) + e_1 f_1 + e_1' f_2 + \frac{1}{2!}\left\{e_1^2 f_{11} + e_1'^2 f_{22}\right.\right.$$

$$\left.\left. + 2e_0 e_1 f_{01} + 2e_0 e_1' f_{02} + 2e_1 e_1' f_{12}\right\}\right]$$

$$= (\theta + e_2) - \bar{Y}^2 - \bar{Y}e_0 - \bar{Y}e_1 f_1 - \bar{Y}e_1' f_2 - \frac{\bar{Y}}{2}\left\{e_1^2 f_{11} + e_1'^2 f_{22} + 2e_0 e_1 f_{01} + 2e_0 e_1' f_{02}\right.$$

$$\left. + 2e_1 e_1' f_{12}\right\} - \bar{Y}e_0 - e_0^2 - e_0 e_1 f_1 - e_0 e_1' f_2$$

$$= \left(\frac{1}{N}\sum_{i=1}^{N} Y_i^2 - \bar{Y}^2\right) + e_2 - 2\bar{Y}e_0 - \bar{Y}e_1 f_1 - \bar{Y}e_1' f_2 - \frac{\bar{Y}}{2}e_1^2 f_{11} - \frac{\bar{Y}}{2}e_1'^2 f_{22}$$

$$- \bar{Y}e_0 e_1 f_{01} - \bar{Y}e_0 e_1' f_{02} - \bar{Y}e_1 e_1' f_{12} - e_0^2 - e_0 e_1 f_1 - e_0 e_1' f_2$$

$$d_g - \sigma_Y^2 = \left(e_2 - 2\bar{Y}e_0 - \bar{Y}e_1 f_1 - \bar{Y}e_1' f_2\right) - \left(e_0^2 + e_0 e_1 f_1 + e_0 e_1' f_2 + \bar{Y}e_0 e_1 f_{01}\right.$$

$$\left. + \bar{Y}e_0 e_1' f_{02} + \bar{Y}e_1 e_1' f_{12} + \frac{\bar{Y}}{2}e_1^2 f_{11} + \frac{\bar{Y}}{2}e_1'^2 f_{22}\right) \cdot \tag{3.5}$$

Taking expectation on both sides of (3.5) and ignoring terms in $(e_i, e_i')$, $i = 0,1,2$ Murthy (1967), the bias in $d_g (= E(d_g) - \sigma_Y^2)$ up to terms of order $\left(\frac{1}{n}\right)$ is given by

$$E(d_g) - \sigma_Y^2 = E(e_2) - 2\bar{Y}E(e_0) - \bar{Y}f_1 E(e_1) - \bar{Y}f_2(e_1') - E(e_0^2) - f_1 E(e_0 e_1)$$

$$- f_2 E(e_0 e_1') - \bar{Y}f_{01} E(e_0 e_1) - \bar{Y}f_{02} E(e_0 e_1') - \bar{Y}f_{12} E(e_1 e_1')$$

$$- \frac{\bar{Y}}{2} f_{11} E(e_1^2) - \frac{\bar{Y}}{2} f_{22} E(e_1'^2) \cdot$$

Using values of the expectations given in (3.1) to (3.3), we have

$$Bias(d_g) = E(d_g) - \sigma_Y^2 = -\frac{\mu_{20}}{n} - f_1\frac{\mu_{11}}{n} - f_2\frac{\mu_{11}}{n'} - \bar{Y}f_{01}\frac{\mu_{11}}{n} - \bar{Y}f_{02}\frac{\mu_{11}}{n'} - \bar{Y}f_{12}\frac{\mu_{02}}{n'}$$

$$-\frac{\bar{Y}}{2}f_{11}\frac{\mu_{02}}{n} - \frac{\bar{Y}}{2}f_{22}\frac{\mu_{02}}{n'} \quad (3.6)$$

Now squaring (3.5) on both the sides and then taking expectation, the mean square error to the first degree of approximation is given by

$$E(d_g - \sigma_Y^2)^2 = E\Big\{(e_2 - 2\bar{Y}e_0 - \bar{Y}f_1e_1 - \bar{Y}f_2e_1') - (e_0^2 + e_0e_1f_1 + e_0e_1'f_2 + \bar{Y}e_0e_1f_{01}$$

$$+ \bar{Y}f_{02}e_0e_1' + \bar{Y}e_1e_1'f_{12} + \frac{\bar{Y}}{2!}f_{11}e_1^2 + \frac{\bar{Y}}{2}f_{22}e_1'^2)\Big\}^2$$

$$= E(e_2^2) + 4\bar{Y}^2E(e_0^2) + \bar{Y}^2f_1^2E(e_1^2) + \bar{Y}^2f_2^2(e_1'^2) - 4\bar{Y}E(e_0e_2) - 2\bar{Y}f_1E(e_1e_2)$$

$$- 2\bar{Y}f_2E(e_1'e_2) + 4\bar{Y}^2f_1E(e_0e_1) + 4\bar{Y}^2f_2E(e_0e_1') + 2\bar{Y}^2f_1f_2E(e_1e_1').$$

Using values of the expectations given in (3.1) to (3.3), we have

$$\text{MSE}(d_g) = E(d_g - \sigma_Y^2)^2$$

$$= \frac{1}{n}(\mu_{40} + 4\bar{Y}\mu_{30} + 4\bar{Y}^2\mu_{20} - \mu_{20}^2) + 4\bar{Y}^2\frac{\mu_{20}}{n} - 4\bar{Y}\frac{1}{n}(\mu_{30} + 2\bar{Y}\mu_{20}) + \bar{Y}^2f_1^2\frac{\mu_{02}}{n}$$

$$+ \bar{Y}^2f_2^2\frac{\mu_{02}}{n'} - 2\bar{Y}f_1\frac{1}{n}(\mu_{21} + 2\bar{Y}\mu_{11}) - 2\bar{Y}f_2\frac{1}{n'}(\mu_{21} + 2\bar{Y}\mu_{11}) + 4\bar{Y}^2f_1\frac{\mu_{11}}{n}$$

$$+ 4\bar{Y}^2f_2\frac{\mu_{11}}{n'} + 2\bar{Y}^2f_1f_2\frac{\mu_{02}}{n'}$$

$$= \frac{1}{n}\Big(\mu_{40} - \mu_{20}^2\Big) + \bar{Y}^2f_1^2\frac{\mu_{02}}{n} + \bar{Y}^2f_2^2\frac{\mu_{02}}{n'} - 2\bar{Y}f_1\frac{\mu_{21}}{n} - 2\bar{Y}f_2\frac{\mu_{21}}{n'} + 2\bar{Y}^2f_1f_2\frac{\mu_{02}}{n'}$$

$$. (3.7)$$

On employing the condition $f_2 = -f_1$, the mean square error of $d_g$ becomes

$$\text{MSE}(d_g)$$

$$= \frac{1}{n}(\mu_{40} - \mu_{20}^2) + \bar{Y}^2f_1^2\frac{\mu_{02}}{n} + \bar{Y}^2f_1^2\frac{\mu_{02}}{n'} - 2\bar{Y}f_1\frac{\mu_{21}}{n} + 2\bar{Y}f_1\frac{\mu_{21}}{n'} - 2\bar{Y}^2f_1^2\frac{\mu_{02}}{n'}$$

$$= \frac{1}{n}(\mu_{40} - \mu_{20}^2) + \bar{Y}^2 f_1^2 \frac{\mu_{02}}{n} - \bar{Y}^2 f_1^2 \frac{\mu_{02}}{n'} - 2\bar{Y}f_1 \frac{\mu_{21}}{n} + 2\bar{Y}f_1 \frac{\mu_{21}}{n'}. \qquad (3.8)$$

The optimum value of $f_1$ minimizing the mean square error of $d_g$ is given by

$$f_1^* = \frac{\mu_{21}}{\bar{Y}\mu_{02}}, \qquad (3.9)$$

which when substituted in (3.8) gives the minimum value of mean square error of $d_g$ as

$$\text{MSE}(d_g)_{\min} = \frac{1}{n}(\mu_{40} - \mu_{20}^2) - \left(\frac{1}{n} - \frac{1}{n'}\right)\frac{\mu_{21}^2}{\mu_{02}}. \qquad (3.10)$$

As we know that the mean square error of usual conventional unbiased estimator $s_y^2 = \frac{1}{n-1}\sum_{i=1}^{n}(y_i - \bar{y})^2$ of population variance $\sigma_Y^2$ is $\frac{1}{n}\left(\mu_{40} - \mu_{20}^2\right)$ and from (3.10) the mean square error of the proposed estimator $d_g$ comes out to be $\frac{1}{n}(\mu_{40} - \mu_{20}^2) - \left(\frac{1}{n} - \frac{1}{n'}\right)\frac{\mu_{21}^2}{\mu_{02}}$ showing that the mean square error of the proposed estimator $d_g$ is less than that of the usual conventional unbiased estimator $s_y^2$ of population variance $\sigma_Y^2$.

## 4. Empirical Study

For comparing efficiency of the proposed estimator, let us consider the data given in Cochran (1977) dealing with Paralytic Polio cases 'Placebo' Y group and Paralytic Polio cases in not inoculated group X. We have calculated the required values of $\mu_{rs}$ and a comparison is made.

For $n = 34$ and $n' = 50$ (say), we have

$\mu_{20} = 9.8894,$          $\mu_{02} = 7.1865882 \times 10^7$

$\mu_{40} = 421.96088,$          $\mu_{21} = 93.464705 \times 10^3$

Mean Square Error of usual conventional unbiased estimator is 9.534136697 and Mean Square Error of the proposed estimator is 8.390090538. The percent relative efficiency (PRE) of the proposed estimator over the usual conventional unbiased estimator is 113.63568 which shows that the proposed estimator is more efficient than the usual conventional unbiased estimator.

## Acknowledgment

## References

1. Cochran, W. G. (1977) – 'Sampling Techniques', 3$^{rd}$ Edition, John Willey and Sons, New York.

2. Das, A.K. and Tripathi, T.P. (1978) – Use of auxiliary information in estimating the finite population variance. Sankhya, C 40, 139-148.

3. Liu, T.P. (1974) – A general unbiased estimator for the variance of a finite population. Sankhya, C 36, 23-32.

4. Murthy, M. N. (1967) – 'Sampling Theory and Methods', 1$^{st}$ Edition, Statistical Publishing Society Calcutta (India).

5. Senapati, S. C. and Sahoo, L. N. (2006) – An alternative class of estimators in double sampling, Bulletin of the Malaysian Mathematical Sciences Society 29(1), 89-94.

6. Singh, R. Karan. (1982) – Generalized double sampling estimators for the ratio and product of population parameters, Journal of Indian Statistical Association, 20, 39-49.

7. Srivastava, S.K. and Jhajj, H.S. (1980) – A class of estimators using auxiliary information for estimating finite population variance. Sankhya, C 42, 87-96.

8. Sukhatme, P.V., Sukhatme, B.V., Sukhatme, S. and Asok, C. (1984) – Sampling theory of surveys with applications, 3$^{rd}$ edition, Ames, Iowa (USA) and Indian Society of Agricultural Statistics, New Delhi (India).

Peeyush Misra and R. Karan Singh