

USE OF ANCILLARY INFORMATION IN COLLAPSING OF STRATA

By

S. MOHANTY AND PADAM SINGH

Institute of Agricultural Research Statistics, New Delhi

1. INTRODUCTION

If the population is highly variable, resort may be made for deep stratification so that only one unit is sampled from each stratum. In this case usual formula for estimating $V(\bar{y}_{st})$ cannot be used. Cochran (1963) suggested a technique called the method of "Collapsed strata" for estimating $V(\bar{y}_{st})$ which consists in grouping the strata in pairs such that two strata in pair should have equal mean and size. As these assumptions are too restrictive, Seth (1966) suggested a modification to Cochran's procedure.

In the present paper a general method of collapsing any number of strata and the use of an ancillary character for collapsing have been suggested. Further, the same ancillary character is used in improving the estimate of the parameter under study.

2. THE METHOD

Let L be the number of strata and N_i, \bar{Y}_i, S_i^2 denote the number of units in the i th stratum, i th stratum mean and variance of the character Y under study respectively. Let y_1, y_2, \dots, y_L be the sample drawn by selecting one unit from each stratum. The stratified random sampling estimate is given by

$$\bar{y}_{st} = \sum_{i=1}^L \frac{N_i}{N} y_i \quad \dots(2.1)$$

and its variance is given by

$$V(\bar{y}_{st}) = \sum_{i=1}^L \frac{N_i^2}{N^2} S_i^2 \quad \dots(2.2)$$

The estimate of (2.2) as proposed by Cochran (1963) is

$$\sum_{i=1}^{L/2} \frac{N_i^2}{N^2} (y_{i1} - y_{i2})^2 \quad \dots(2.3)$$

where $N_{i_1} = N_{i_2} = N_i$, i_1 th, i_2 th strata are collapsed to form the i th stratum and L is assumed to be an even integer. The above estimate is unbiased when $\bar{Y}_{i_1} = \bar{Y}_{i_2}$. Seth (1966) suggested another estimate of (2.2) given by

$$\frac{1}{N^2} \sum_{i=1}^{L/2} (N_{i_1}^p y_{i_1} - N_{i_2}^p y_{i_2})(N_{i_1}^{2-p} y_{i_1} - N_{i_2}^{2-p} y_{i_2}), \quad \dots(2.4)$$

where the strata are collapsed in such a way that the strata means are proportional to some power ' p ' of the number of units in the stratum. The bias of the estimate (2.4) is small when

$$N_{i_1}^p \bar{Y}_{i_1} - N_{i_2}^p \bar{Y}_{i_2} \text{ is negligible.}$$

Now the general method of collapsing consists of collapsing k_i strata to form a new i th stratum ($i=1, 2, \dots, l$) such that the k_i strata to be collapsed must satisfy the relation

$$N_{ij}^{p_i} \bar{Y}_{ij} = N_{ij}'^{p_i} \bar{Y}_{ij}' \quad \dots(2.5)$$

for all $j \neq j' = 1, 2, \dots, k_i$. When $p_i = 1$, the above relation means that the k_i strata to be collapsed must have same total value of the character Y under study.

The estimate of (2.1) for the above case is given by

$$\frac{1}{N^2} \sum_{i=1}^l \sum_{j \neq j'}^{k_i} (N_{ij}^{p_i} y_{ij} - N_{ij}'^{p_i} y_{ij}')(N_{ij}^{2-p_i} y_{ij} - N_{ij}'^{2-p_i} y_{ij}') \quad \dots(2.6)$$

The expected value of (2.6) is given by

$$V(\bar{y}_{st}) + \sum_{i=1}^l \sum_{j \neq j'}^{k_i} (N_{ij}^{p_i} \bar{Y}_{ij} - N_{ij}'^{p_i} \bar{Y}_{ij}')(N_{ij}^{2-p_i} \bar{Y}_{ij} - N_{ij}'^{2-p_i} \bar{Y}_{ij}') \quad \dots(2.7)$$

The estimate (2.6) is unbiased when the conditions given in (2.5) are satisfied.

The proposed estimate (2.6) over-estimate the variance as it appears from (2.7) and is likely to be less efficient compared to the sampling designs where strata are formed in such a manner that more than one unit are selected per stratum to form the sample.

3. USE OF ANCILLARY INFORMATION

As it is not possible to examine the conditions given in (2.5) for strata to be collapsed to form a new stratum, we may use an ancillary character X which is highly correlated with Y and which is readily available. Let \bar{X}_i be the mean of the ancillary variate for the i th stratum ($i=1,2,\dots,L$). We collapsed those k_i strata to form a new i th stratum ($i=1,2,\dots,l$) such that for the strata to be collapsed the conditions

$$N_{ij} p_i \bar{X}_{ij} = N_{ij'} p_i \bar{X}_{ij'}, j \neq j' = 1, 2, \dots, k_i \text{ are satisfied.}$$

The above conditions ensure that the conditions in (2.5) are also (almost) satisfied as Y and X are highly correlated and hence the bias in the estimate (2.6) tends to zero.

After collapsing the strata by use of the ancillary variate Y , the same variate X is again used to improve the estimate of \bar{Y} by using ratio method of estimation. Here we will consider both separate and combined methods of ratio estimation.

(a) Separate ratio estimate is defined as

$$\bar{y}_{RS} = \sum_{i=1}^l \frac{N_{i.}}{N} \times \frac{\bar{y}_{i.}}{\bar{x}_{i.}} \times \bar{X}_{N_{i.}}, \quad \dots(3.1)$$

where

$$N_{i.} = \sum_{j=1}^{k_i} N_{ij}, \bar{y}_{i.} = \sum_{j=1}^{k_i} N_{ij} y_{ij} / N_{i.},$$

$$x_{i.} = \sum_{j=1}^{k_i} N_{ij} x_{ij} / N_{i.},$$

and

$$\bar{X}_{N_{i.}} = \sum_{r=1}^{k_i} \sum_{j=1}^{N_{ir}} X_{ij} / N_{i.}.$$

The variance of (3.1) is given by (upto first degree of approximation)

$$V(\bar{y}_{RS}) = \sum_{i=1}^l \frac{N_i^2}{N^2} [V(\bar{y}_i) + R_i^2 V(\bar{x}_i) - 2R_i \text{Cov}(\bar{x}_i, \bar{y}_i)] \quad \dots(3.2)$$

where $R_i = \bar{Y}_{N_i} / \bar{X}_{N_i}$.

Now under the assumptions $N_{ij}\bar{Y}_{ij} = N_{ij}'\bar{Y}_{ij}'$ (considering $p_i=1$ for all i) for all $j \neq j' = 1, 2, \dots, k_i$, we get

$$\left. \begin{aligned} \hat{V}(\bar{y}_i) &= \frac{1}{N_i^2} \sum_{j \neq j'=1}^{k_i} (y_{ij}N_{ij} - y_{ij}'N_{ij}')^2 \\ \hat{V}(\bar{x}_i) &= \frac{1}{N_i^2} \sum_{j \neq j'=1}^{k_i} (x_{ij}N_{ij} - x_{ij}'N_{ij}')^2 \\ \text{and Cov}(\bar{x}_i, \bar{y}_i) &= \frac{1}{N_i^2} \sum_{j \neq j'=1}^{k_i} (y_{ij}N_{ij} - y_{ij}'N_{ij}')(x_{ij}N_{ij} - x_{ij}'N_{ij}') \end{aligned} \right\} \dots(3.3)$$

So the estimate of (3.2) is given by

$$\hat{V}(\bar{y}_{RS}) = \frac{1}{N^2} \sum_{i=1}^l \sum_{j \neq j'=1}^{k_i} [(y_{ij}N_{ij} - y_{ij}'N_{ij}')^2 + r_i^2 (x_{ij}N_{ij} - x_{ij}'N_{ij}')^2 - 2r_i (x_{ij}N_{ij} - x_{ij}'N_{ij}')(y_{ij}N_{ij} - y_{ij}'N_{ij}')] \quad \dots(3.4)$$

where $r_i = \hat{R}_i = \bar{y}_i / \bar{x}_i$, which on simplification becomes

$$\frac{1}{N^2} \sum_{i=1}^l \sum_{j \neq j'=1}^{k_i} [(y_{ij}N_{ij} - y_{ij}'N_{ij}') - r_i (x_{ij}N_{ij} - x_{ij}'N_{ij}')]^2 \quad \dots(3.5)$$

(b) The combined ratio estimate, with usual notation is defined as

$$\bar{y}_{RC} = \bar{y}_{st} \bar{X}_N / \bar{x}_{st} \quad \dots(3.6)$$

The variance of (3.6) upto first degree of approximation,

when
$$\sum_{i=1}^l k_i$$

is sufficiently large is given by

$$V(\bar{y}_{RC}) = \sum_{i=1}^l (N_i^2/N^2) [V(\bar{y}_i) + R^2 V(x_i) - 2R \text{Cov}(\bar{x}_i, \bar{y}_i)] \quad \dots(3.7)$$

where $R = \bar{Y}_N / \bar{X}_N$. After substituting the values of the estimate from (3.3), the estimate of (3.7) is given by

$$\hat{V}(\bar{y}_{RC}) = \frac{1}{N^2} \sum_{i=1}^l \sum_{j \neq j'=1}^{k_i} [(N_{ij}y_{ij} - N_{ij'}y_{ij'}) - r(N_{ij}x_{ij} - N_{ij'}x_{ij'})]^2 \quad \dots(3.8)$$

where $r = \bar{y}_{st} / \bar{x}_{st}$.

The estimates (3.1) and (3.6) will be more efficient than the estimate given in (2.1), if

$$\rho_i > \frac{1}{2} R_i \cdot \sqrt{\frac{V(\bar{x}_i)}{V(\bar{y}_i)}}, \quad \rho_i > \frac{1}{2} R \sqrt{\frac{V(\bar{x}_i)}{V(\bar{y}_i)}} \text{ respectively,}$$

where ρ_i is the population correlation coefficient between y and x in the i th stratum, $i=1, 2, \dots, l$.

The proposed estimates (3.5) and (3.8) are, of course, subject to the same objection as estimate (2.7), that, they overestimate the variance if the conditions in (2.5) are not satisfied and are likely to be less efficient compared to the sampling designs where strata are so formed that at least two samples are drawn per stratum to form the sample.

4. DOUBLE SAMPLING

So far we have considered the case when \bar{X} is known for each stratum. But this is not the case always. So, at times it pays to devote a part of the resources in estimating the population mean of the ancillary variate for each stratum by taking a large sample. Let a sample of size n'_{ij} from (ij) th stratum is first drawn to estimate $\bar{X}_{N_{ij}}$, where (ij) th ($j=1, 2, \dots, k_i$) strata are collapsed to form the new i th stratum.

Let $n'_{i.} = n'_{i1} + n'_{i2} + \dots + n'_{ik_i}$ and

$$\bar{x}_{n'_{i.}} = \frac{1}{n'_{i.}} \sum_{j=1}^{k_i} n'_{ij} \bar{x}_{ij}.$$

Then from n'_{ij} units of each stratum one unit is selected for which the Y character is measured.

(a) The separate ratio estimate under double sampling is defined as

$$\bar{y}_{\text{DRS}} = \sum_{i=1}^l \frac{N_i}{N} \frac{\bar{y}_i}{\bar{x}_i} \bar{x}_{n'_i} \quad \dots(4.1)$$

The variance of the above estimate is given by

$$V(\bar{y}_{\text{DRS}}) = \sum_{i=1}^l \frac{N_i^2}{N^2} \left[V(\bar{y}_i) + R_i^2 \left\{ V(\bar{x}_i) - V(\bar{x}_{n'_i}) \right\} - 2R_i \{ \text{Cov}(\bar{x}_i, \bar{y}_i) - \text{Cov}(\bar{x}_{n'_i}, \bar{y}_{n'_i}) \} \right] \quad \dots(4.2)$$

The estimate of the above variance, after substituting the estimate of each term is given by

$$\begin{aligned} \hat{V}(\bar{y}_{\text{DRS}}) &= \frac{1}{N^2} \sum_{i=1}^l \sum_{j \neq i=1}^{k_i} \left[\frac{N_i}{(N_i - k_i)} \frac{(n'_i - k_i)}{n'_i} \left\{ (y_{ij} N_{ij} - y_{ij}' N_{ij}') - r_i (x_{ij} N_{ij} - x_{ij}' N_{ij}') \right\}^2 \right] + \frac{1}{N^2} \sum_{i=1}^l \sum_{j \neq j'=1}^{k_i} \left(\frac{1}{n'_i} - \frac{1}{N_i} \right) \\ &\quad \frac{N_i \cdot k_i}{(N_i - k_i)} \times \left(N_{ij} y_{ij} N_{-ij}' - y_{ij}' \right)^2 \quad \dots(4.3) \end{aligned}$$

(b) The double sampling combined ratio estimate is defined as

$$\bar{y}_{\text{DRC}} = \frac{\bar{y}_{st}}{\bar{x}_{st}} \bar{x}_{n'} \quad \dots(4.4)$$

where

$$\bar{x}_{n'} = \frac{1}{N} \sum_{i=1}^l N_i \bar{x}_{n'_i}.$$

The variance of the above estimate is given by

$$V(\bar{y}_{\text{DRC}}) = V(\bar{y}_{st}) + R^2 [V(\bar{x}_{st}) - V(\bar{x}_{n'})] - 2R [\text{Cov}(\bar{x}_{st}, \bar{y}_{st}) - \text{Cov}(\bar{x}_{n'}, \bar{y}_{n'})]$$

and the estimate is given by

$$\begin{aligned} \hat{V}(\bar{y}_{\text{DRC}}) &= \frac{1}{N^2} \sum_{i=1}^l \sum_{j \neq j'=1}^{k_i} \frac{N_i}{(N_i - k_i)} \frac{(n'_i - k_i)}{n'_i} \left\{ (y_{ij} N_{ij} - y_{ij}' N_{ij}') - r \left(x_{ij} N_{ij} - x_{ij}' N_{ij}' \right) \right\}^2 + \frac{1}{N^2} \sum_{i=1}^l \sum_{j \neq j'=1}^{k_i} \frac{N_i \cdot k_i}{(N_i - k_i)} \\ &\quad \left(\frac{1}{n'_i} - \frac{1}{N_i} \right) \left(y_{ij} N_{ij} - y_{ij}' N_{ij}' \right)^2 \quad \dots(4.5) \end{aligned}$$

5. ILLUSTRATION

The data on the block level estimate for wheat crop in Patna district of Bihar for the year 1967-68 have been taken for illustrating the use of ancillary information in collapsing strata. For this purpose two blocks consisting of 15 V.L.W. circles as strata have been taken. The character under study is the crop cutting yield per plot of size $1/80$ of an acre in kg and the ancillary information is the eye estimate per acre in kg. For our purpose we have selected one village per V.L.W. circle. Our sample consists of one field for crop cutting experiment and four fields for eye estimates in each selected village. The collapsing of strata have been done on the basis of the auxiliary variate such that total of auxiliary variate remains more or less same for the strata to be collapsed. Following the double sampling procedure given in section 4, the results are as follows :

<i>Method of estimation</i>	<i>Mean in kg./plot</i>	<i>Variance</i>	<i>% gain in efficiency over stratified random sampling</i>
Stratified random sampling	5.797	0.8770	
Double sampling separate ratio	5.872	0.2870	23.87
Double sampling combined ratio	5.846	0.2885	23.47

with parameters $n'_{ij}=4$ for all i and j , $k_1=k_2=k_3=3$, $k_4=2$, $k_5=4$, $l=5$ and $r_{i.}=(0.0135, 0.0133, 0.0174, 0.0128, 0.0136)$, $r=0.0141$.

SUMMARY

Here an attempt has been made to utilize ancillary information in collapsing any number of strata in case of deep stratification and the same ancillary information is used in estimating the parameter under study more efficiently.

ACKNOWLEDGEMENT

We are thankful to the referee for his valuable comments to improve the paper.

REFERENCES

1. Cochran, W.G. (1963) Sampling Technique, John Wiley and Sons, Inc. New York.
2. Seth, G.R. (1966) "On collapsing of strata," Journal of Indian Society of Agricultural Statistics, Vol. XVIII, 1966.