



Ranked Set Sampling from Finite Population under Randomization Framework

Anil Rai and Praveen Krishana

Indian Agricultural Statistics Research Institute, New Delhi

Received 12 February 2010; Revised 18 May 2013; Accepted 19 August 2013

SUMMARY

Ranked Set Sampling (RSS) provides efficient estimation of population mean as compared to Simple Random Sampling (SRS). The procedure of RSS generates virtual stratification of the population, as a consequence of this, it provides better representation of the population in a selected sample as compared to SRS. Published literatures related to RSS are either based on infinite population theory or super population framework. In this article, an attempt has been made to examine the RSS procedure in randomization framework of survey sampling. It has been observed that, this procedure pertains to the category of equal probability sampling method, *i.e.*, the probability of including every unit of the population in a sample is equal. The estimator for estimating population mean has been proved to be unbiased and an expression of its variance has been derived in terms of variability among the sampling units of individual ranks. Statistical properties of this sampling strategy have been studied using simulation under two different cases, *i.e.*, (i) usual case when $N = mn^2$, where N is the population size, n denotes the RSS sample size from each cycle and m is the number of cycle. (ii) two phase sampling when $N > mn^2$. It was found that the proposed RSS estimator is always better than SRS estimator for both the cases. Gain of 20 to 40 percent is achieved in RSS over SRS for equivalent sample size.

Keywords: Ranked set sampling, Randomization approach, Simulation, Finite population.

1. INTRODUCTION

Ranked set sampling (RSS) was first introduced by McIntyre (1952) to improve upon Simple Random Sampling (SRS) for situation where some preliminary ranking of sampled units is possible. The basic idea was to randomly partition the sampled units into small group and each member of individual group is ranked relative to other members of the group. Based on this ranking one member from each group is selected for quantification. Name of this procedure has been coined by Halls and Dell (1966) whereas mathematical foundation was provided by Takahasi and Wakimoto (1968) and Dell (1969). Dell and Clutter (1972), David and Levine (1972), Stokes (1976, 1977, 1980) studied RSS procedure when ranks are subjected to errors. It

was observed that estimators of population mean and population variance are asymptotically unbiased under some realistic assumptions. Bohu and Dougals (1994) and Stokes and Sagar (1988) studied inferential aspects of RSS through empirical investigation. Yu and Lam (1997) proposed regression type of estimator for estimating population mean. Neweihi and Sinha (2000) and Zhua Chen (2000) further studied this procedure from inferential point of view.

RSS has great potential for observational economy over its competitive techniques and it has been applied by McIntyre (1952) for estimation of mean pasture and forage yield. Due to its cost efficiency, RSS has been extensively used in forestry (Evans 1967, Cobby *et al.* 1985, Halls and Dell 1966, Jewiss 1981,

Martin *et al.* 1980). Since, in case of RSS, ranking of units can be done on the basis of visual inspection or expert opinion, therefore, it can be applied in visual ranking of forest trees for estimation of timber volume, area soils can be ranked with respect to its physical and chemical properties including toxicity based on surface properties, presence of vegetation etc. Geographical variables, which are spatial in nature can be ranked on the basis of remote sensing images of electromagnetic reflectance and this has numerous applications in estimation of parameters across different fields including agriculture, environment, hydrology, geology etc.

Most of the literature available on ranked set sampling (RSS) procedure is based on properties of statistical distributions and based on infinite population theory as it provides mathematical simplicity. In this paper an attempt has been made to estimate the finite population mean using RSS procedure by following randomization approach based on sampling from finite population. In order to achieve this objective, second order inclusion probabilities of sampling units were obtained. The effective sub-population size of each rank were also derived as sampling unit of each rank selected in the sample are implicitly drawn in the sample from sub-set of finite population units under consideration. Estimators of population mean and its variance were derived for two different cases *i.e.* (i) when finite population size is equal to product of square of sample size and number of cycles and (ii) when finite population size is greater than product of square of sample size and number of cycles. Since, it was not possible to compare statistical properties of these derived estimators of RSS with corresponding estimators SRS theoretically, therefore, simulation study has been conducted to compare them empirically to demonstrate efficiency of RSS over SRS in this framework.

Let a sample of size $n' = nm$ is drawn through RSS procedure from a finite population of size N . Here, n and m are integers, representing sub-sample size and number of cycles respectively. For simplicity, let us consider $N = mn^2$. Let the parameter of interest is population mean, which is linear in functions of population values. In a particular cycle let n random sets of n units are selected by Simple Random Sampling Without Replacement (SRSWOR). The units within each set are ranked relative to each other on some

criterion which is easily observable, say, size measure. From the first set, sampling unit having 1st rank is selected, and then from the set, unit of 2nd rank is selected and so on until the n^{th} ranked unit is selected from the n^{th} set. The above procedure is repeated m times (cycles) without replacing selected units of the previous cycles in the population. Hence, a total of mn ultimate sampling units are selected from m cycles. Here, a sample of mn sampling units consists of m units of each rank. It can be seen that the sets are formed at random and the ranking in different sets are independent from each other. Hence, the units in the ultimate samples are independent of each other. This ranking of sampling units is based on some auxiliary character that may or may not be quantifiable. It is presumed that in any case ranking is perfect. Now, observations on the study character, *i.e.* Y may be taken on ultimate sampling units.

2. INCLUSION PROBABILITY OF SAMPLING UNITS

The inclusion probability of any sampling unit is defined as the probability of including the unit in the sample. Let, units of population be denoted by $\{U_1, U_2, \dots, U_N\}$ where N is population size. It may be presumed that complete population is ranked on the basis of size measure and the ranked unit of the population is denoted by $\{U_{(1)}, U_{(2)}, \dots, U_{(N)}\}$, where $U_{(k)}$ denotes the unit with k^{th} rank in the population. A SRSWOR sample of size n is selected from the population. The units of this ordered sample are denoted by $\{y_{(i)} : (i = 1, 2, \dots, n)\}$, where $y_{(i)}$ denotes the i^{th} ranked of sampled unit. It is assumed that there is no error in ranking in the population as well as in the sample. Probability that U_k [k^{th} ranked unit of the population] has i^{th} rank in the SRSWOR sample of size n is given by

$$P[y_{(i)} = U_{(k)}] = \frac{\binom{k-1}{i-1} \binom{N-k}{n-i}}{\binom{N}{n}} \quad (2.1)$$

In case of RSS procedure as described above, n such sub-samples are selected by SRSWOR without replacing the previous selected sub-samples into the population. Unit of 1st rank is selected from 1st set, unit of second rank from 2nd set and so on. Hence, i^{th} unit of the population may be included in the ranked set sample if it has either 1st rank in first set, or 2nd rank

in the second set, or 3rd rank in the 3rd set and so on. Further, any sub-sample may be selected at any order. Therefore, the inclusion probability of i^{th} unit can be obtained by summing the probabilities of selecting i^{th} unit having ranks 1 to n . Thus it can be easily seen that

$$\pi_i = \frac{n}{N} \tag{2.2}$$

Therefore, in case RSS the inclusion probability of every unit of the population is same. With the above described procedure of RSS for selecting a sample of size n from n^2 units, total number of all possible ranked set samples is

$$N_{RSS} = \frac{(n^2)}{n!n! \dots (n \text{ times})} \tag{2.3}$$

All the N_{RSS} ranked set sample are equal probable and probability of drawing any sample at a particular drawn is given by

$$P(S_{RSS}^*) = \frac{1}{N_{RSS}} \tag{2.4}$$

Here, S_{RSS}^* denotes the ranked set sample.

Similarly, the joint probability of inclusion of any two units of the population, *i.e.*, i and j is given by

$$\begin{aligned} \pi_{ij} = & \frac{\left[\sum_{k=1}^n \sum_{s=1}^{n-1} \binom{i-1}{k-1} \sum_{l=0}^{\min(n-k,z)} \binom{n^2-j}{n-k-1} \binom{z}{1} \right. \\ & \times \binom{j-k-1}{s-1} \binom{n^2-j-(n-k-1)}{n-s} \\ & + \sum_{s=1}^n \sum_{k=1}^{n-1} \binom{n^2-j}{n-k} + \sum_{s=1}^n \sum_{k=11}^{n-1} \binom{n^2-j}{n-k} + \sum_{m=0}^{\min(n-k,z)} \binom{i-1}{k-m-1} \\ & \left. \times \binom{z}{m} \binom{i-k+m}{s-1} \binom{n^2-n-i-m+k-1}{n-s} \right] \\ & \tag{2.5} \end{aligned}$$

Since, the inclusion probability of every unit of population in the sample as well as probability of selection of any sample at any draw is also equal, it can be concluded that above described ranked set sampling procedure is equal probability sampling.

3. EFFECTIVE POPULATION SIZE FOR THE i^{th} RANK IN THE RANKED SET SAMPLE

In the above procedure of RSS it can be seen that certain units of the population have zero probability of being selected as i^{th} rank unit from sets generated through randomization procedure. It can be seen that $(n - 1)$ units of the population will never get i^{th} rank, $i = 1, 2, \dots, n$, in any of the sets. So for a particular i^{th} ranked unit in a sample is not selected from the population size N but a reduced population of size $N^* = N - (n - 1) = N - n + 1$. Therefore, it may be visualized that a ranked set sample of size mn consists of m units which are i^{th} rank selected from a population of size N^* through SRSWOR.

In view of the above fact, let, $Y_{\{i\}j}$ = value of the i^{th} unit from i^{th} sub-sample of size n in the j^{th} cycle. $S = \{U_1, U_2, \dots, U_N\}$, be the set of all the units present in the population of size N . $S\{i\} = \{U_1, U_2, \dots, U_N\} \subset S$, be the set of all the units from the population of size N^* which can be assigned i^{th} rank in random sample of size n and $s\{i\} = \{y_{\{i\}j} : j = 1, 2, \dots, m\} \subset S\{i\}$, be the set of units having the i^{th} rank; $i = 1, 2, \dots, n$ selected in the ranked set sample of size mn .

3.1 Estimation of Population Mean

Consider the situation of estimation of population mean when sample has been selected through RSS.

Case 1: When $N = mn^2$

Sample mean is the usual estimator for population mean. In sampling by RSS procedure as described above for the case $N = mn^2$, sample mean is given by

$$\bar{y}_{RSS} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m y_{\{i\}j} \tag{3.1}$$

The above estimator is unbiased for estimation of population mean *i.e.* $E(\bar{y}_{RSS}) = \bar{Y}$.

The variance of \bar{y}_{RSS} can be obtained as

$$V(\bar{y}_{RSS}) = V \left[\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m y_{\{i\}j} \right]$$

$$= \frac{1}{n^2} \sum_{i=1}^n \left[V \left(\frac{1}{m} \sum_{j=1}^m y_{\{i\}j} \right) \right]$$

As $y_{\{i\}j}$'s are independent, hence the covariance term resulting from the expansion of above expression vanishes. Therefore,

$$V(\hat{y}_{RSS}) = \frac{1}{n^2} \sum_{i=1}^n [V(\bar{y}_{\{i\}})] \tag{3.2}$$

where, $[V(\bar{y}_{\{i\}})]$ is the variance of the units having i^{th} rank in the ranked set sample. Here m numbers of i^{th} ranks are selected in the sample by SRSWOR from N^* ($= N - n + 1$) units of the population that can take i^{th} rank in sample. Hence,

$$V(\bar{y}_i) = \left(\frac{1}{m} - \frac{1}{N'} \right) S_{Y_{\{i\}}}^2$$

Therefore, the above variance expression can be written as

$$V(\bar{y}_{RSS}) = \frac{1}{n^2} \left(\frac{1}{m} - \frac{1}{N'} \right) \sum_{i=1}^n S_{Y_{\{i\}}}^2 \tag{3.3}$$

where,

$$S_{Y_{\{i\}}}^2 = \frac{1}{N'-1} \sum_{s_{\{i\}}} (Y_{\{i\}j} - \bar{Y}_{\{i\}})^2 \tag{3.4}$$

Thus it can be seen that variance expression is expressed in terms of variability of units within ranks. The expression for estimate of variance can be obtained as usual.

Case 2. When $N > mn^2$

In case, population size $N > mn^2$, the sampling is done in two phases to select the sample by RSS procedure described above. In first phase, from the population of size $N (> mn^2)$ a sample of size mn^2 is selected by SRSWOR. Then from this first phase sample of mn^2 units a sample of size mn is selected by following the method of RSS described above.

Sample mean for this case is given by

$$\bar{y}_{RSS,2} = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m y_{\{i\}j,k} \tag{3.5}$$

were $y_{\{i\}j,k}$ denotes the value for the k^{th} first phase unit and i^{th} unit of the i^{th} set in the j^{th} cycle, ($i = 1, 2, \dots, n$; $j = 1, 2, \dots, m$; $k = 1, 2, \dots, N$), and $\bar{y}_{RSS,2}$ is the ranked set sample mean when sample is selected in two phases.

Now

$$E(\bar{y}_{RSS,2}) = E_1 E_2 \left[\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m y_{\{i\}j,k} \right]$$

In the preceding section it was proved that ranked set sample mean is an unbiased estimator for population mean. Hence, if we denote the first phase SRSWOR sample mean by \bar{y}_{SRS} , where

$$\bar{y}_{SRS} = \frac{1}{mn^2} \sum_{j=1}^{mn^2} y_{\{i\}j,k}$$

Hence,

$$E_2 \left[\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m y_{\{i\}j} \right] = \bar{y}_{SRS}$$

It is known that SRSWOR sample mean is an unbiased estimator of population mean.

$$E(\bar{y}_{RSS,2}) = E_1 E_2 \left[\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m y_{\{i\}j,k} \right] = \bar{Y} \tag{3.6}$$

Hence, the ranked set sample mean is an unbiased estimator of population mean even when the population size $N > mn^2$ and sampling is carried out in two phases.

3.2 Variance of Ranked Set Sample Mean

Since, the sample has been selected in two phases, in the first phase a SRSWOR sample of size mn^2 is selected from population of size N , and in the second phase ranked set sample of size mn is selected from the first phase sample.

$$V(\bar{y}_{RSS,2}) = V_1 E_2(\bar{y}_{RSS,2}) + E_1 V_2(\bar{y}_{RSS,2})$$

It has already been proved that \bar{y}_{RSS} is an unbiased estimator of population mean.

Therefore,

$$V_1 E_2(\bar{y}_{RSS,2}) = V_1(\bar{y}_{SRS,1}) = \left(\frac{1}{mn^2} - \frac{1}{N} \right) S_Y^2$$

$$E_1V_2(\bar{y}_{RSS}) = E_1\left[\frac{1}{n^2}\left(\frac{1}{m} - \frac{1}{N'_2}\right)\sum s'^2_{Y\{i\}j}\right]$$

$$= \left[\frac{1}{n}\left(\frac{1}{m} - \frac{1}{N'_2}\right)\right]S^2_Y$$

where

$$N'_2 = mn^2 - n + 1$$

Therefore, we have,

$$V(\bar{y}_{RSS,2}) = \left[\left(\frac{1}{mn^2} - \frac{1}{N}\right) + \frac{1}{n}\left(\frac{1}{m} - \frac{1}{N'_2}\right)\right]S^2_Y \quad (3.7)$$

where

$$s'^2_{Y\{i\}j} = \frac{1}{N'_2 - 1} \sum_{j \in s_2\{i\}} (y_{\{i\}j} - y_{\{i\}}) \quad (3.8)$$

$$S^2_Y = \frac{1}{N - 1} \sum_{k=1}^N (Y_{\{i\}jk} - \bar{Y}) \quad (3.9)$$

Hence above equation gives the variance of the sample mean when sampling is done in two phases following the procedure described above. The expression of unbiased estimator of this variance can be obtained as usual.

4. SIMULATION STUDY

In this simulation study, both the above situations of sample selection have been considered. A population of tri-variate normal population have been generated using variables X, Y and Z . Here Y is taken to be the variable of interest, X as the auxiliary variable, and Z is the variable which was used for ranking as require in the procedure of RSS described above. The parameters of the hypothetical generated tri-variate normal population are

	Y	X	Z
35	326	205.73	204.39
$\mu = 26$; Var-cov matrix =	205.73	265	122.85
32	204.39	122.85	365

4.1 Case 1 $\{N = mn^2 (n = 5, m = 20)\}$

In this case population of 5000 units (sets of trivariate normal variate, X, Y, Z) was generated. 200 samples of size 100 have been selected following the RSS procedure. To make comparison 200 samples of size 100 units were selected by SRSWOR. Comparison of these two designs were made on the basis of statistics calculated with these samples such as Sample mean, Variance, Skewness, Kurtosis, Percentage bias and C.V. for both the cases.

4.2 Case 2 $\{N > mn^2, \bar{X} \text{ is known}\}$

A population of 5000 units of trivariate normal population with variables as X, Y , and Z were generated. Under this case RSS sample has been selected in two phases. 200 different samples of following sizes with SRS at first phase and RSS at second phase were selected.

- (a) SRS of size 40 and then RSS of size 20 ($n = 2, m = 10$)
- (b) SRS of size 90 and then RSS of size 30 ($n = 3, m = 10$)
- (c) SRS of size 160 and then RSS of size 40 ($n = 4, m = 10$)
- (d) SRS of size 250 and then RSS of size 50 ($n = 5, m = 10$)
- (e) SRS of size 1000 and then RSS of size 100 ($n = 10, m = 10$)

For the purpose of comparison, 200 different SRS sample of following equivalent sizes were taken from the same population.

- (a) SRS of size 20
- (b) SRS of size 30
- (c) SRS of size 40
- (d) SRS of size 50
- (e) SRS of size 100

Both the above sampling strategies were compared on the basis of Sample mean, Variance, Skewness, Kurtosis, Percentage bias and C.V. and the results for all the above cases were given in the table.

Table. Comparison of method of RSS with SRS using the simulation results.

Case	Design	Mean	%Bias	Variance	C.V.	Skewness	Kurtosis	Relative efficiency
1	2	3	4	5	6	7	8	9
1	SRS	35.67	1.99	3.65	5.35	0.126	0.035	121.80
	RSS	35.60	1.71	3.00	4.95	0.369	-0.192	
2(a)	SRS	34.32	1.93	16.43	11.58	0.212	-0.231	115.68
	RSS	34.78	0.61	14.20	10.76	0.285	-0.276	
2(b)	SRS	34.79	0.34	9.85	8.99	0.195	-0.228	123.52
	RSS	34.82	0.59	7.97	8.11	0.179	-0.356	
2(c)	SRS	34.96	0.11	7.02	7.57	0.251	-0.011	113.48
	RSS	35.23	1.75	6.18	6.86	0.311	0.081	
2(d)	SRS	34.99	0.003	7.07	7.60	0.271	-0.558	138.33
	RSS	35.01	0.04	5.11	6.46	0.321	-0.0084	
2(e)	SRS	34.97	0.08	2.94	4.90	0.138	-0.25	123.35
	RSS	34.90	0.28	2.38	4.41	0.380	-0.075	

From the above table of simulation results, the following points can be noted

- ❖ The % bias for RSS and SRS is of the same order. It may reduce to zero when all possible samples are considered. Hence, RSS mean is an unbiased estimator of population mean like SRS.
- ❖ It can be seen from column 5 of the table that the variance of RSS mean is always less than variance of SRS mean, and it decreases with increase in the ultimate sample size. Hence, it can be concluded that RSS sampling strategy is always better than SRS sampling strategy.
- ❖ From column, 7 and 8 of the table, it can be noted that the values of skewness and kurtosis is very close to zero. Hence, the distribution of RSS means is symmetric and mesokurtic just like the distribution of SRS means.
- ❖ Column 9 of the table shows relative efficiency of RSS with respect to SRS. It can be concluded that RSS mean is always efficient than SRS mean. As it can be observed from the table, gain in efficiency in case of RSS as compared to SRS is ranging between 13% to 39% respectively.

It can be seen that in this case performance of ranked set sampling under randomized framework for finite population is superior than Simple Random

Sampling as in case of infinite population. However, it is well know, that in case of RSS, perfect ranking will lead to optimal performance of RSS. It has been shown (Patil *et al.* 1994) that in any case of imperfect ranking its efficiency i.e. Relative Precision (RP) with respect to SRS is always high. Efficiency of RSS can be equal to SRS only in case of random ranking. Further, in case of highly skewed sample measurements performance of RSS decreases and in extreme situations its performance can be worse than SRS. However, McIntyre (1952) suggested that in this situation allocation of sub-samples of different rank is to be made in proportion the standard deviation of each rank. This allocation strategy will leads to allocation similar to Neyman allocation in case of stratified sampling.

ACKNOWLEDGEMENTS

Authors are grateful to the Associate Editor and referees for their valuable comments.

REFERENCES

- Bohn, L.L. and Wolfe, D.A. (1994). The effect of imperfect judgment rankings on properties of procedures based on the ranked-set samples analog of the Mann-Whitney-Wilcoxon statistic. *J. Amer. Statist. Assoc.*, **89**, 168-176.
- Chui, N.N. and Sunhat, B.K. (1998). On some aspect of ranked set sampling in parametric estimation. In: *Handbook of Statistics*, **17**, 337-377, N. Balakrishna and Rao, C.R. eds., North Holland, New York.
- Cobby, J.M., Ridout, M.S., Bassett, P.J. and Large, R.V. (1985). An investigation into the use of ranked set sampling on grass and grass-clover swards. *Grass Forage Sci.* **40**, 257-263.
- Cochran, W.G. (1997). *Sampling Techniques*. 3rd edition. Wiley, New York.
- David, H.A. and Levine, D.N. (1972). Ranked set sampling in the presence of judgment error. *Biometrics*, **28**, 553-555.
- Doll, T.R. and Cluster, J.L. (1972). Ranked set sampling theory with order statistics background. *Biometrics*, **28**, 545-555.
- Evans, M.J. (1967). Application of ranked set sampling to regeneration surveys in areas direct-seeded to longleaf pine. Master's Thesis, School of Forestry and Wildlife Management, Louisiana State University, Baton Rouge.
- Halls, L.K. and Dell, T.R. (1966) Trial of ranked set sampling for forage yields. *Forest Sci.*, **12**, 22-26.

- Jewiss, O.R. (1981). Shoot development and number. In: *Sward Measurement Handbook*, J. Hodgson, R.D. Baker, A. Davies, A.S. Laidlaw and J.D. Leaver, eds. The British Grassland Society, Hurley, 93-114.
- Kour, A., Patil, G., Susan, J. and Charles, T. (1996). Environmental sampling with a concomitant variable: A comparison between ranked set sampling and stratified simple random sampling. *J. Appl. Statist.*, **23(2 and 3)**, 231-255.
- Martin, W.L., Sharik, T.L., Oderwald, R.G. and Smith, D.W. (1980). Evaluation of ranked set sampling for estimating shrub phytomass in Appalachian oak forests. Publication Number FWS-4-80, School of Forestry and Wildlife Resources, Virginia Polytechnic Institute and State University, Blacksburg.
- McIntyre, G.A. (1952). A method of unbiased selective sampling, using ranked set. *Austr. J. Agric. Res.*, **3**, 385-390.
- Patil, G.P., Sinha, A.K. and Taillie, C. (1994). Ranked set sampling. In: *Handbook of Statistics*, **12**, Patil, G.P. and Rao, C.R. eds., 167-198, North-Holland, Amsterdam.
- Stokes, S.L. (1977). Ranked set sampling with concomitant variables. *Comm. Statist. -Theory Methods*, **A6**, 1207-1211.
- Stokes, S.L. (1980). Estimation of variance using judgment order ranked set samples. *Biometrics*, **36**, 35-42.
- Takahashi, K. (1970). Practical note on estimation of population means based on samples stratified by means of ordering. *Ann. Inst. Statist. Math.*, **22**, 421-428.
- Takahashi, K. and Wakomoto, K. (1968). On biased estimates of the population mean based on sample stratified by means of ordering. *Ann. Inst. Statist. Math.*, **20**, 1-31.
- Yu, Philip, L.H. and Lam, K. (1997). Regression estimator in ranked set sampling. *Biometrics*, **53**, 1070-1080.