

AN ESTIMATOR IN PPS SAMPLING FOR MULTIPLE CHARACTERISTICS*

By P. K. PATHAK

Indian Statistical Institute

SUMMARY. For pps sampling schemes when the characteristic under study and the selection probabilities are poorly correlated, it is shown how customary estimators of population total can be modified to yield estimators which are more precise than customary ones in the sense of having smaller expected variance under a super population model. The results established in this paper are extensions of similar results proved by Rao (1966). A result conjectured by Rao (1966) concerning Murthy's estimator (1957) in this connection is shown to be true.

1. INTRODUCTION AND STATEMENT OF THE PROBLEM

Consider a population $\pi = (U_1, U_2, \dots, U_j, \dots, U_N)$ of N elements. Let P_j and Y_j respectively denote the values of some P and Y characteristics of U_j . Unless otherwise stated the subscript j runs from 1 through N . It is assumed that $\sum P_j = 1$. Consider now the problem of estimating $Y = \sum Y_j$ on the basis of a sampling scheme $\{S, P\} = \{(s, p(s)) : s \in S\}$ defined on π where S denotes the set of samples to be selected from π and $p(s)$ denotes the probability of selection associated with the sample $s (s \in S)$. We restrict our attention to the following class of unbiased estimators of Y .

$$t_p(s) = \sum b_{ij}(Y_j/P_j) \quad \dots (1)$$

where $E(b_{ij}) = P_j$ and b_{ij} is the coefficient of (Y_j/P_j) in $t_p(s) (s \in S)$.

It is well-known that when Y_j and P_j are highly correlated, estimators like $t_p(s)$ in sampling schemes with probabilities P_j lead to considerable gain in efficiency as compared with the customary estimators like $n^{-1}\sum y_j$ suitable for equal probability sampling. However, in surveys, where population parameters of several characteristics are to be estimated on the basis of some sampling scheme with probabilities P_j , it is likely that some characteristics will be highly correlated with P_j , while others may be poorly correlated and sometimes even uncorrelated with P_j . For example Rao (1966) considers an example of chicken population wherein the number of chickens in a farm and the farm size are very poorly correlated. In this case estimates of total number of chickens obtained on the basis of farms selected by probabilities proportional to farm size, are likely to be inefficient. For such situations Rao (1966) has suggested alternative estimators of Y under several sampling schemes. These estimators can be obtained by substituting NY_j/P_j for y_j in $t_p(s)$ given in (1). These

* This paper was written by Dr. P. K. Pathak after seeing the paper by Dr. J. N. K. Rao sent for publication in *Sankhyā* and which appears in this number. Dr. Pathak's paper is being published in the same number with the permission of Dr. J. N. K. Rao.—*Editor*.

SANKHYĀ : THE INDIAN JOURNAL OF STATISTICS : SERIES A

estimators are likely to have smaller mean square error than the corresponding unbiased estimators $t'_p(s)$ particularly in small samples. The bias of such estimators remains the same for all sample sizes unless Y_j and P_j are uncorrelated in which case they are unbiased. Under a super population model Rao (1966) has shown that the new estimators have smaller expected variance than the corresponding unbiased estimators of Y for some sampling schemes. For sampling with unequal probabilities without replacement of size two, Rao (1966) leaves open the problem of comparing Murthy's estimator (1957) and the new estimator obtained from it.

In this note we give a general treatment of the problem considered by Rao (1966), and solve the problem concerning Murthy's estimator (1957).

2. A NEW ESTIMATOR

Under the earlier described sampling scheme $(S, P) = \{(s, p(s)) : s \in S\}$ defined on π , consider the following estimator

$$t'_p(s) = N \sum b_j Y_j \quad \dots (2)$$

which is obtained by replacing Y_j by $NY_j P_j$ in $t_p(s)$ given in (1).

Since

$$E\{t_p(s)\} = Y_j, \quad \dots (3)$$

we can easily get $E\{t'_p(s)\}$ by substituting $NY_j P_j$ for Y_j in the right side of (3).

Therefore

$$E\{t'_p(s)\} = N \sum Y_j P_j. \quad \dots (4)$$

Thus the bias in $t'_p(s)$, $B(t'_p)$, is given by

$$\begin{aligned} B(t'_p) &= N \sum Y_j P_j - \Sigma Y_j \\ &= N \Sigma (Y_j - N^{-1} \Sigma Y_j) (P_j - N^{-1}) \\ &= N^2 \text{cov}(Y_j, P_j). \quad \dots (5) \end{aligned}$$

It is evident from (5) that $t'_p(s)$ will be unbiased if Y_j and P_j are uncorrelated. It is also important to note that the bias remains the same for all sample sizes. Especially for small samples when Y_j and P_j are poorly correlated, the bias may be expected to be small relative to the standard error.

Now in order to compare the two estimators, $t'_p(s)$ and $t_p(s)$, we consider the following model (Rao, 1966).

$$Y_j = m + e_j \quad \dots (6)$$

where e_j are mutually independent and satisfy $e(e_j | P_j) = 0$ and $v(e_j | P_j) = a$ ($a > 0$), ϵ and v respectively denoting the expectation and the variance under the given model.

AN ESTIMATOR IN PPS SAMPLING

It is easy to see that under the model (6), we have

$$\begin{aligned} \epsilon\{V(t_p(s))\} &= \epsilon\{E[\Sigma(b_{ij} - Eb_{ij})Y_j/P_j]^2\} \\ &= E\{\epsilon[\Sigma(b_{ij} - Eb_{ij})Y_j/P_j]\} + E\{\{E[\Sigma(b_{ij} - Eb_{ij})Y_j/P_j]\}^2\} \\ &= a\Sigma V(b_{ij})/P_j^2 + m^2V[\Sigma b_{ij}/P_j] \end{aligned} \quad \dots (7)$$

where E and V respectively denote the expectation and the variance under the given sampling scheme.

It is also straightforward to verify that

$$\epsilon\{V(t_p^*(s))\} = aN^2\Sigma V(b_{ij}) + m^2N^2V[\Sigma b_{ij}] \quad \dots (8)$$

The following theorem asserts that under some regularity conditions, satisfied by many estimators in practice, $\epsilon\{V(t_p^*(s))\} \leq \epsilon\{V(t_p(s))\}$.

Theorem 1 : *Let the estimator $t_p(s)$ given in (1) satisfy $\Sigma b_{ij} = 1$ for all samples and $\Sigma V(b_{ij})/P_j^2 \geq N^2\Sigma V(b_{ij})$. Then under the model (6)*

$$\epsilon\{V(t_p^*(s))\} \leq \epsilon\{V(t_p(s))\} \quad \dots (9)$$

where $t_p^*(s)$ is the estimator obtained by replacing Y_j by NY_jP_j in $t_p(s)$.

Proof : Evidently

$$\begin{aligned} &\epsilon\{V(t_p(s))\} - \epsilon\{V(t_p^*(s))\} \\ &= a\Sigma V(b_{ij})/P_j^2 - aN^2\Sigma V(b_{ij}) + m^2V(\Sigma b_{ij}/P_j) - m^2N^2V(\Sigma b_{ij}) \\ &\geq m^2V(\Sigma b_{ij}/P_j) \geq 0 \end{aligned} \quad \dots (10)$$

since by assumption $V(\Sigma b_{ij}) = 0$ and

$$\Sigma V(b_{ij})/P_j^2 \geq N^2\Sigma V(b_{ij}).$$

This completes the proof.

Corollary 1 : $\epsilon\{V(t_p^*(s))\} \leq \epsilon\{V(t_p(s))\}$ if $V(b_{ij}) = c_1P_j - c_2P_j^2$, where c_1 and c_2 are non-negative scalars.

Remark : In practice the verification of the inequality $\Sigma V(b_{ij})/P_j^2 \geq N^2\Sigma V(b_{ij})$ will be greatly facilitated if we know $V(b_{ij})$. An easy way of computing $V(b_{ij})$ when we know $V(t_p(s))$, is to let $Y_j = P_j$ and $Y_{j'} = 0$ for all $j' \neq j$ in $V(t_p(s))$.

We mention without going into details that

$$(i) \quad V(b_{ij}) = P_j(1 - P_j)/n$$

in sampling with unequal probabilities with replacement,

$$(ii) \quad V(b_{ij}) = (1 - (n - 1)/n(N - 1) + k(n - k)/(N - 1))P_j(1 - P_j)$$

in the Rao-Hartley-Cochran sampling scheme (1962) and

$$(iii) \quad V(b_{ij}) = P_j/n - P_j^2$$

for the Horvitz-Thompson estimator (1952) in pps sampling without replacement of size n . As a consequence of Corollary 1, we have

$$\{e\{V(t_{\mu}^*(\theta))\}\} \leq e\{V(t_{\mu}(\theta))\}$$

in all these cases.

In the next section we prove the conjecture by Rao (1966) concerning Murthy's estimator (1957).

3. MURTHY'S ESTIMATOR

In pps sampling without replacement of size two, Murthy's estimator (1957) of the population total is given by

$$t_{\mu}(\theta) = \frac{1}{(2-p_1-p_2)} ((1-p_2)y_1/p_1 + (1-p_1)y_2/p_2) \quad \dots (11)$$

where the symbols have their usual meanings.

The modification of $t_{\mu}(\theta)$ leads to the following estimator

$$t_{\mu}^*(\theta) = \frac{N}{(2-p_1-p_2)} ((1-p_2)y_1 + (1-p_1)y_2). \quad \dots (12)$$

Rao (1966) has conjectured that $t_{\mu}^*(\theta)$ has smaller expected variance than $t_{\mu}(\theta)$ under the model (6). We proceed now to prove this conjecture.

For Murthy's estimator (1957) the condition $\sum b_{ij} = 1$ is satisfied. Further we have

$$V(b_{ij}) = \sum_{j=1 \neq i}^N P_j P_j (1-P_j - P_j) / (2-P_j - P_j). \quad \dots (13)$$

Thus

$$e\{V(t_{\mu}^*(\theta))\} \leq e\{V(t_{\mu}(\theta))\}$$

$$\text{if } \sum_j \sum_{j=1 \neq i}^N P_j P_j (1-P_j - P_j) (2-P_j - P_j)^{-1} (P_j^2 - N^2) \geq 0. \quad \dots (14)$$

To prove (14), we make use of the following lemma.

Lemma 1: Let $(U_1, V_1), \dots, (U_p, V_p), \dots, (U_N, V_N)$ be N pairs of real numbers.

$$\text{Then } \sum U_i V_i \geq N^{-1} (\sum U_i) (\sum V_i) \quad \dots (15)$$

if $U_i - U_k \geq 0$ if and only if (iff) $V_i - V_k \geq 0$ for all i, k ($i, k = 1, 2, \dots, N$).

AN ESTIMATOR IN PPS SAMPLING

The proof of the lemma is omitted.

To prove (14), it suffices to prove that

$$\sum_j \sum_{j' \neq j} P_j (1 - P_j - P_{j'}) (2 - P_j - P_{j'})^{-1} (P_j^{-1} - N^2 P_j) \geq 0. \quad \dots (10)$$

Letting

$$U_j = P_j^{-1} - N^2 P_j$$

and
$$V_j = \sum_{j' \neq j} P_j (1 - P_j - P_{j'}) (2 - P_j - P_{j'})^{-1},$$

we find that
$$U_i - U_k \geq 0 \quad \text{iff} \quad P_k - P_i \geq 0$$

and
$$V_i - V_k = \left[\sum_{j' \neq i, j' \neq k} P_j (P_k - P_j) (2 - P_i - P_j)^{-1} (2 - P_k - P_j)^{-1} \right]$$

$$+ (P_k - P_i) (1 - P_i - P_k) (2 - P_i - P_k)^{-1} \geq 0$$

iff
$$P_k - P_i \geq 0.$$

Therefore,
$$U_i - U_k \geq 0 \quad \text{iff} \quad V_i - V_k \geq 0.$$

Hence, from Lemma 1, we have

$$\sum_j \sum_{j' \neq j} \{P_j (1 - P_j - P_{j'}) (2 - P_j - P_{j'})^{-1} (P_j^{-1} - N^2 P_j)\}$$

$$\geq N^{-1} \left(\sum_j \sum_{j' \neq j} P_j (1 - P_j - P_{j'}) (2 - P_j - P_{j'})^{-1} \right) \sum_j (P_j^{-1} - N^2 P_j) \geq 0 \quad \dots (17)$$

since $\sum_j P_j^{-1} \geq N^2$.

Thus Rao's conjecture that

$$c\{V(\theta^*(\theta))\} \leq c\{V(\theta(\theta))\}$$

is indeed true.

Incidentally for estimating $V(\theta^*(\theta))$, it is trivial to verify that if $c(\theta_j(\theta))$ is an unbiased estimator of $V(\theta_j(\theta))$, then an unbiased estimator of $V(\theta^*(\theta))$ can be got by replacing Y_j by $NY_j P_j$ in $c(\theta_j(\theta))$.

For interesting applications of Theorem 1 in a variety of situations, the reader may refer to the paper by Rao (1966).

ACKNOWLEDGEMENT

The author is greatly indebted to Professor C. R. Rao for suggesting him to work on the problem discussed in this paper.

SANKHYĀ : THE INDIAN JOURNAL OF STATISTICS : SERIES A

REFERENCES

- HORVITZ, D. G. and THOMPSON, D. J. (1952): A generalization of sampling without replacement from a finite universe. *J. Amer. Stat. Assoc.*, **47**, 663-685.
- MURTHY, M. N. (1957): Ordered and unordered estimators in sampling without replacement. *Sankhyā, Series A*, **19**, 379-390.
- RAO, J. N. K. (1958): Alternative estimators in p.p.s. sampling for multiple characteristics. *Sankhyā, Series A*, **20**, 47-60.
- RAO, J. N. K., HARTLEY, H. O. and COCHRAN, W. G. (1962): On a simple procedure of unequal probability sampling without replacement. *J. Roy. Stat. Soc., B*, **24**, 492-491.

Paper received : September, 1965.