R. P. Chakrabarty and P. J. Tsai Jackson State University and University of Georgia

Summary

In this paper, we compare the relative merits of four sampling methods. (1) Direct sampling without replacement. (2) Direct sampling with replacement. (3) Inverse sampling without replacement and (4) Inverse sampling with replacement that may be used to estimate population proportions. The stability of the variance estimator in addition to the efficiency of the estimator is taken into consideration in judging the performance of an estimator. Conditions under which one method is preferable to another are obtained. Sample size needed for the use of asymptotic results are indicated by simulation.

1. INTRODUCTION

In socio-economic and ecological surveys, one is frequently interested in estimating either the number NP or proportion P which possesses a specified characteristic in a finite population of size N. The unbiased estimator of the proportions, its variance and the data-base variance estimator for direct sampling with and without replacement are well known (see Cochran 1963). Haldane (1945) first obtained an unbiased estimator of P and its asymptotic variance for the case of inverse sampling with replacement. Recently Best (1974) obtained a closed form for the variance of Haldane's unbiased estimator of the P and compared its efficiency with that of the maximum likelihood estimator given by Feller (1968). Best's numerical evaluation of efficiency showed that the unbiased estimator of P was generally more efficient than the maximum likelihood estimator. Scheaffer (1974) compared the unbiased estimator in direct sampling with replacement and the biased estimator in inverse sampling with replacement, and obtained the asymptotic stabilities of the variance estimators. He has shown that under certain conditions inverse sampling with replacement provides a more stable variance estimator in large samples.

It may be noted that the results of Haldane, Feller, Best and Scheaffer were obtained assuming infinite populations. Some new results for the important case of sampling without replacement from a finite population are obtained in this paper. Only unbiased estimators are considered in view of Bests' results. Special emphasis has been given to obtain results that are exact for any sample size

whenever it is possible. We note that the estimate of the variance of an estimator is often used in drawing statistical inferences (e.g., confidence limits for the parameters). It is, therefore, desirable that the variance estimator should be as stable as possible. The importance of investigation of the stability of the variance estimator was pointed out by Chakrabarty and Rao (1967) and Chakrabarty (1973) in estimation of ratios. We, therefore, consider the stability of the variance estimator associated with an estimator in addition to the efficiency of the estimator itself in judging the performance of an estimator. Without loss of generality, we shall discuss the problem of estimation of the proportion P only in a finite population of N units.

Some notations are introduced here that are used in subsequent sections:

- N : number of units in the population.
- N_1 : number of units in the population that possesses a specified attribute or characteristic or falls into specified class.
- n : total sample size.
- n1 : sample size from the specified class.
- $P : P = N_1 / N$ population proportion to be estimated.
- Q : Q = 1 P.
- P1 : the unbiased estimator of P by direct sampling with replacement; DWR.
- p2: the unbiased estimator of P by direct sampling without replacement; DWTR.
- p3: the unbiased estimator of P by inverse sampling with replacement; IWR.
- p₄ : the unbiased estimator of P by inverse sampling without replacement; IWTR.
- $q_i : q_i = 1 p_i$, i = 1, 2, 3, 4.

: sampling fraction,
$$f = \frac{n}{N}$$

c : finite population correction factor. 1 - f = 1 - n/N. с:

V(p_i): the variance of p_i, i = 1, 2, 3, 4. v(p_i): the estimator of V (p_i), i = 1, 2, 3,4. (v(p_i)): the exact variance of v(p_i), i = 1, 2. (v(pi)): asymptotic variance of v(p;), i = 1, 2, 3, 4.

2. DIRECT SAMPLING

- 2.1 Direct Sampling with Replacement; DWR. The probability density function (p.d.f.) of n_1
- is $f(n_1) = (n_1) P^{n_1 n-n_1}_{Q};$ $n_1 = 0, 1, 2, \dots n$ P+Q=1, $P,Q \ge 0$

It is well known that the unbiased estimator of P, its variance and unbiased variance estimator are given by

$$P_1 = n_1/n$$

 $V(P_1) = PQ/n$
and

 $v(p_1) = p_1 (1-p_1) / (n-1)$ (2.1)

respectively (see Cochran).

It can be shown that the exact variance of $\mathbf{v}(\mathbf{p}_1)$ is given by

$$V_{e}(v(p_{1})) = (A_{1} + A_{2} + A_{3} + A_{4}) / (n^{4}(n-1)^{2})$$
where
$$A_{1} = m_{1}^{(4)} - 2(n-3)m_{1}^{(3)}$$

$$A_{2} = (n^{2} - 6n + 7)m_{1}^{(2)}$$

$$A_{3} = (n-1)^{2}m_{1}^{(1)}$$

$$- A_{4} = \{(n-1)m_{1}^{(1)} - m_{1}^{(2)}\} 2 .$$

$$m_{1}^{(k)} = E(n_{1}^{(k)}) = E[n_{1}(n_{1}-1)...(n_{1}-k+1)];$$

$$k = l_{2} 2 . 3 . 4$$
(2)

$$= 1, 2, 3, 4$$
 (2.2)

The asymptotic variance of $v(p_1)$ by using the well known formula (Kendall and Sturat, 1968)

$$V(g(x)) \stackrel{:}{=} g'(E(x)) V(x) ,$$

is obtained as
$$V_{a}(v(p_{1})) = \{v'(p)\}^{2} V(p_{1})$$
(2.3)
$$= (1-2P)^{2} PQ/(n(n-1)^{2})$$

2.2 Direct Sampling without Replacement; DWTR

The p.d.f. of n₁ is

$$f(n_{1}) = \frac{\begin{pmatrix} N_{1} \\ n_{1} \end{pmatrix}}{\begin{pmatrix} N_{1} \\ n_{1} \end{pmatrix}} \begin{pmatrix} N-N_{1} \\ n-n_{1} \end{pmatrix}}, \\ \begin{pmatrix} N \\ n \end{pmatrix}}{\begin{pmatrix} N \\ n \end{pmatrix}}, \\ n_{1} = 0, 1, \dots, \min(n, N_{1})$$
(2.4)

The unbiased estimator of P, its variance and unbiased variance estimator are given by:

$$p_2 = n_1/n$$

 $V(p_2) = PQ(1-(n-1)/N-1))/n$
and $v(p_2) = p_2q_2(1-n/N)/(n-1)$ (2.5)

respectively. Using the relations between moments about origin and descending factorial moments of n_1 for the $p \cdot d \cdot f \cdot (2 \cdot 4)$ it can be shown that the exact variance of $v(p_2)$ is given by:

$$V_e(v(p_2) = C^2 (B_1 + B_2 + B_3 + B_4)/(n^4 (n-1)^2)$$

where

$$B_{1} = m_{2}^{(4)} - 2 (n-3)m_{2}^{(3)}$$

$$B_{2} = (n^{2} - 6n + 7)m_{2}^{(2)}$$

$$B_{3} = (n-1)^{2}m_{2}^{(1)}$$

$$B_{4} = -((n-1)m_{2}^{(1)} - m_{2}^{(2)}).$$

$$m_{2}^{(k)} = E (n_{1}^{(k)}) = n^{(k)}N_{1}^{(k)} / N^{(k)};$$

$$k = 1, 2, 3, 4$$
(2.6)

The asymptotic variance of $v(p_2)$ is given by

$$V_{a} (v(p_{2}) = (\frac{c}{(n-1)} (1-2P))^{2} V(p_{2})$$
$$= (1-2P)^{2} PQ (1-n/N)^{2} (1-(n-1)/(1-(n-1))) / (n(n-1)^{2}) (2.7)$$

3. INVERSE SAMPLING

Inverse Sampling is a procedure where sampling is continued until n1 units possessing the specified characteristic have been obtained. The number n1 is pre-determined and the Sample Size $n(\ge n_i)$ is thus a random variable.

3.1. Inverse Sampling with Replacement:
IWR. The p^d f of n is

$$f(n) = {n-1 \choose n_q-1} P^{n_q} Q^{n-n_1}$$
; $n=n_1, n_1+1, ...$
(3.1)

Haldane (1945) first gave the unbiased estimator of P and its variance as

$$P_{3} = (n_{1}-1)/(n-1)$$

and
$$V(p_3) = \frac{p^2 Q}{n_1} \left(1 + \frac{2! Q}{n_1 + 1} + \frac{3! Q^2}{(n_1 + 1)(n_1 + 2)} + \dots \right)$$

(3.2)

respectively. Note that $V(n) = n_1 Q/P^2$ (3.3) \mathbf{E} (n) = n₁ / P andconsequeltly, the asymptotic variance of p₃ is

$$V(p_{3}) \stackrel{!}{=} V(n) \left[(n_{1}-1)^{2} / (E(n) - 1)^{4} \right]$$

$$= n_{1} (n_{1}-1)^{2} P^{2} Q / (n_{1} - P)^{4} (3.4)$$
Now, $E \left[\frac{(n_{1}-1) (n_{1}-2)}{(n-1) (n-2)} \right] = \sum_{n > n_{1}}^{\infty} \frac{(n_{1}-1) (n_{1}-2)}{(n-1) (n-2)} x$

$$\left(\frac{n-1}{n_{1}-1} \right) P^{n_{1}} Q^{n-n_{1}} = P^{2}$$

$$E \left[\left(\frac{n_{4}-1}{n-1} \right)^{2} - \frac{(n_{4}-1)(n_{5}-2)}{(n-1) (n-2)} \right]$$

$$= E (p_{3}^{2}) - E^{2} (p_{3}) = V(p_{3})$$
Thus
$$v(p_{3}) = \left(\frac{n_{1}-1}{n-1} \right)^{2} - \frac{(n_{1}-1) (n_{1}-2)}{(n-1) (n-2)}$$

$$= \frac{(n_{1}-1) (n-n_{1})}{(n-1) (n-2)} = \frac{p_{3}(1-p_{3})}{(n-1) (n-2)}$$

$$\frac{1}{(n-1)^2(n-2)} = \frac{1}{(n-2)} (3.5)$$

is an unbiased estimator of V(p3). The asymptotic variance of $v(p_2)$ is obtained by using (2.3) as

$$Va(v(p_3)) = \left[n_1(n_1-1)^2 \left\{ (3P-2)n_1^2 + (2-5P)n_1P+2P^2 \right\}^2 \right]^2$$

$$P^4Q \int \left[(n_1-P)^6 (n_1-2P)^4 \right] \quad (3.6)$$

3.2 Inverse Sampling without Replacement; IWTR

The p.d.f. of sample size n is

$$f(n) = \frac{\binom{N_1}{n_1 - 1} \binom{N - N_1}{n - n_1}}{\binom{N}{n - 1}} \cdot \frac{N_1 - n_1 + 1}{N - n + 1};$$

$$n = n_1, n_1 + 1, \dots, N \qquad (3.7)$$

The kth ascending factorial moment of n is

$$m_{4}^{[k]} = E(n^{[k]}) = n_{1}^{[k]} (N_{1} + 1)^{[k]} / (N_{1} + 1)^{[k]}$$
Where $n^{[k]} = n(n+1) (n+2) \dots (n+k-1)$ (3.8)

consequently, E (n) =
$$n_1$$
 (N+1) / (N₁+1) (3.9)
and V(n) = n_1 (N+1) (N-N₁) (N₁-n₁+1)

 $(N_1+1)^2(N_1+2)$

or

$$V(n) = \frac{n_1(1+\frac{1}{N}) (P - \frac{n_1 - 1}{N}) Q}{(P + \frac{1}{N})^2 (P + \frac{2}{N})} (3.10)$$

Since

$$E \left(\frac{n_{1}-1}{n-1}\right) = \sum_{n=1}^{N} \left(\frac{n_{1}-1}{(n-1)} \left(\frac{n_{1}-1}{n-1}\right)\right)$$

Ν

$$\frac{N_{1} - n_{1} + 1}{N - n + 1} = \frac{N_{1}}{N} = P$$

$$P_{4} = \frac{n_{1} - 1}{n - 1}$$
(3.11)

is an unbiased estimator of P,

The asymptotic variance of p_A is

$$V(p_4) \stackrel{:}{=} \frac{{(n_1 - 1)}^2}{\left\{ E(n) - 1 \right\}^4} V(n)$$
(3.12)

where E(n) and V(n) are given in (3.9) and (3.10) respectively.

We show that an unbiased estimator of the exact variance of p_4 can be obtained even though an explicit expression for $V(p_4)$ could not be obtained. By definition, we have

Since

$$V(p_{4}) = E(p_{4}^{2}) - P^{2}.$$

$$\left(\frac{n_{1}-1}{n-1}\right) - \frac{(n_{1}-1)(n_{1}-2)}{(n-1)(n-2)} = \frac{p_{4}q_{4}}{n-2}$$

and $E((n_1-1)(n_1-2)/(n-1)(n-2)) = P(N_1-1)/(N-1)$ We get $E(p_4q_4/(n-2)) = V(p_4) + PQ/(N-1)$ Further, $E(p_4q_4) = PQ - V(p_4)$

Consequently, $v(p_4) = (p_4q_4)(1 - (n-1)/N)/(n-2)$

(3.13)

is an unbiased estimator of $V(p_4)$.

Now, the derivative of $v(p_4)$ with respect to n is

$$\frac{dv(p_4)}{dn} = (n_1 - 1)(\frac{-2n^2 + (2 + 3n_1)n + 2 - 5n_1}{(n-1)^3(n-2)^2}$$

$$\frac{(1-\frac{n-1}{N})-\frac{n-n_1}{(n-1)^2(n-2)}\cdot\frac{1}{N}}{(say)}$$

Accordingly, the asymptotic variance of $v(p_4)$ is 2

$$V_a(v(p_4)) \triangleq (g(E(n)))^2 V(n)$$

where E(n) and V(n) are given in (3.9) and (3.10) respectively. The above expression is too complicated for analytical comparison but can be easily evaluated numerically.

4. COMPARISON OF DIRECT AND INVERSE SAMPLING

Any comparison of direct sampling with a fixed sample size n to inverse sampling with variable sample size n would be made on the basis of expected sample size, E(n). Note that in

IWR: $E(n) = n_1 / \dot{P} \text{ or } n_1 = P.E(n)$

IWTR:
$$E(n) = n_1(1+1/N)/(P+1/N)$$

= n_1/P for large N (4.1)

Therefore, n is replaced by P.E(n) in variance formulas in inverse sampling for comparison with variance formulas in direct sampling. The results are given in the sequal.

4.1. Efficiencies of Unbiased Estimators pi It can be shown that for large n or E(n)

$$V(p_1) = PQ/n$$
 , $V(p_2) \neq PQ(1-n/N)/n$
 $V(p_3) \neq PQ/E(n)$, $V(p_4) \neq PQ(1-E(n)/N)/E(n)$
(4.2)

Thus p_2 in DWTR is more efficient than p_2 in

DWR and p_{4} in IWTR is more efficienct than

 p_3 in IWR. But p_1 and p_3 as well as p_2 and

 p_4 are equally efficient in large samples.

Sampling without replacement is definitely preferable but one cannot choose between DWTR and IWTR on the basis of the efficiency of unbiased estimators of P.

4.2. Stabilities of Variance Estimators v(pi)

Following Rao and Chakrabarty (1967) and Chakrabarty (1973) the stability of the variance estimator $v(p_i)$ relative to $v(p_j)$ is given by

 $\frac{(CV(v(p_i)))^2/(CV(v(p_i)))^2}{(V(v(p_i)))^2 = V(v(p_i))/(V(p_i))^2}$ $\frac{(4.3)}{(V(p_i))^2} = \frac{V(v(p_i))}{(V(p_i))^2}$ $\frac{(4.3)}{(V(p_i))^2} = \frac{V(v(p_i))}{(V(p_i))^2}$ $\frac{(4.3)}{(V(p_i))^2}$ $\frac{(4.3)}{(V(p_i))^2}$

Now from (2.2), (2.4), (2.6), and (2.7) it can be shown that $Ve(v(p_1)) \ge Ve(v(p_2))$

and $Va(v(p_1)) > Va(v(p_2))$ for all P.

Substituting n by P.E(n) in variance form-

ulas in inverse sampling and equating E(n) in inverse sampling to n in direct sampling we obtain after some algebraic manipulations the conditions under which one variance estimator is more stable than another. The details which are given in a technical report by the authors (1976) are omitted here and only the results are summarized below:

- (1) v(p₂) in DWTR is more stable than v(p₁) in DWR for all P.
- (2) v(p₃) in IWR is more stable than v(p₁) in DWR iff P>0.6.
- (3) $v(p_4)$ in IWTR is more stable than $v(p_1)$ in DWR iff

P>
$$\frac{2-d-(1-d)^{\frac{1}{2}}}{3-2d+2(1-d)^{\frac{1}{2}}}$$

(4) $v(p_3)$ in IWR is more stable than $v(p_2)$ in in DWTR iff $2(1-d)^{\frac{1}{2}}$ $2(1-d)^{\frac{1}{2}}$

$$\frac{2 - (1 - d)^{\frac{2}{2}}}{3 + 2(1 - d)^{\frac{1}{2}}} \zeta P \zeta \frac{2 + (1 - d)^{2}}{3 - 2(1 - d)^{\frac{1}{2}}}$$

(5) v(p₄) in IWTR is more stable than v(p₂) in DWTR iff

P>
$$\frac{3-2d}{5-4d}$$

(6) v(p₄) in IWTR is more stable than v(p₃) in IWR iff

P>
$$\frac{2-d+2(1-d)}{3-2d+3(1-d)^{\frac{1}{2}}}$$
.

where d=n/N

These results are shown graphically indicating the regions of stable variance estimator. The implication of th**es**e results are that if we have prior knowledge about the range of the proportion parameter P, then it is possible to choose the optimum sampling techniques from the four techniques we have discussed, particularly between DWTR and IWTR.

We may mention that we have also evaluated numerically the stabilities of the variance estimators by generating 1000 samples, each of size 100 from a population with N=1000 and P=.8 on cyber 74. DWTR and IWTR was equally efficient but IWTR provided the most stable variance estimator.

We have also investigated the sample size needed for the use of asymptotic resluts by simulation for different values of P. The asymptotic variances were found to be almost equal to exact and/or simulated variances for moderate sample sizes. These results are given in a technical report by the authors (1976).

REFERENCES

- Best, D. J., (1974). "The Variance of the Inverse Binomial Estimator," <u>Biometrika</u>, 385-386.
- Chakrabarty, R. P. and Rao, J. N. K., (1967). "The Bias and Stability of Jack-Knife Variance Estimator in Ratio Estimation," <u>Proc. Amer. Statist. Assoc.</u>, (Social Statistics section), 326-331.
- Chakrabarty, R. P., (1973). "A Note on the Small Sample Theory of the Ratio Estimator in Certain Specified Populations."

Jour. Indian Soci. Agri. Statist., Vol. 25, No. 2, 49-57.

Cochran, W. G., (1963). <u>Sampling Techniq</u>ues, Wiley, New York.

- Feller, W., (1968). <u>An Introduction to Proba-</u> bility Theory and Its Applications, Vol 1, 3rd Ed., Wiley, New York.
- Haldane, J. B. S., (1945). "On a Method of Estimating Frequencies," <u>Biometrika</u>, 33:222-225.
- Kendall, M. G. and Stuart, A., (1968). <u>The</u> <u>Advanced Theory of Statistics</u>, Vol. 1, London, Griffin.
- Scheaffer, R. L., (1974). "On Direct Versus Inverse Sampling form Mixed Populations", <u>Biometrics</u>, 30:187-198.
- Tsai, P. J. and Chakrabarty, R. P., (1976). "Comparison of Direct and Inverse Sampling for Estimation of Population Proportions," Technical Report, Dept. of Statistics and Computer Science, Univ. of Georgia.







Figure 2: The Region of the Stable Variance Estimators in DWR and IWR



