

SOME REMARKS ON SAMPLING WITH REPLACEMENT

BY DES RAJ AND SALEM H. KHAMIS

American University of Beirut

1. Introduction and summary. In order to estimate the mean of a finite population from a random sample, the sample units may be selected in two ways. A pre-determined number m of units may be selected with replacement or sampling with replacement may be continued till a desired number n of distinct units is obtained. The first procedure is called sampling with replacement while the second may be called sampling without replacement. A comparison is generally made between the two procedures with $m = n$. This comparison, however, is not fair since costs are usually proportional to the number of distinct units and in the first procedure this number would be less than or equal to n .

In the first procedure the population mean is usually estimated by the sample mean based on all the units in the sample including repetitions while in the second procedure the estimate is generally made to depend on the distinct units only. The object of this paper is to show that the estimate making use of only the distinct units is superior in either procedure. The following results are proved in this paper:

- (i) In sampling with replacement the estimate of the mean based on distinct units in the sample is superior to the estimate based on the total sample size when (a) the total sample size is fixed in advance, while the number of distinct units in the sample is a random variable, and, (b) when the total sample size is a random variable while the number of distinct units is fixed in advance.
- (ii) The same is true of ratio estimates. It is also shown that the bias is numerically less if the ratio estimate is based on distinct units regardless of whether these are fixed in advance or considered as random variables.
- (iii) Expressions for the estimation of the variances of the various estimates considered in this paper are given.
- (iv) The above results are extended to multistage sampling.

2. Statement of the problem. Let us consider a finite population consisting of N sampling units. Suppose we are interested in estimating the population mean \bar{Y} for a character y , from a sample selected with replacement with equal probabilities. We consider the following two sampling schemes.

Scheme A. We select with replacement a total sample of size m fixed in advance.

We denote the number of distinct sample units selected by the random variable u .

Scheme B. We select with replacement a sample of n distinct units fixed in advance. We denote the total sample size, including repetitions, by the random variable v .

Received February 15, 1957; revised November 18, 1957.

We consider Scheme A first. Let a_1, a_2, \dots, a_u be the u distinct units selected in the sample and let k_r be the number of times the r th unit a_r occurs in the sample, with $\sum k_r = m$. We compare the bias and variances of the usual estimate

$$(1) \quad \bar{y}_m = \frac{1}{m} \sum_{r=1}^u k_r y_r,$$

and the estimate

$$(2) \quad \bar{y}_u = \frac{1}{u} \sum_{r=1}^u y_r$$

based on the distinct units only, where y_r is the value of the character y for the r th distinct unit in the sample.

The estimate (1) is well known to be unbiased and that its variance is given by

$$(3) \quad V(\bar{y}_m) = \frac{1}{m} \left(1 - \frac{1}{N}\right) \sigma^2 = \left(Q - \frac{1}{N}\right) \sigma^2$$

where

$$(4) \quad Q = \frac{1}{m} \left(1 + \frac{m-1}{N}\right)$$

and

$$(5) \quad \sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2.$$

For a given u the expected value of \bar{y}_u is \bar{Y} so that \bar{y}_u is an unbiased estimate of \bar{Y} . With regard to the variance of \bar{y}_u we have

$$(6) \quad V(\bar{y}_u) = \left[E\left(\frac{1}{u}\right) - \frac{1}{N} \right] \sigma^2.$$

The estimate \bar{y}_u is superior to \bar{y}_m if

$$(7) \quad E\left(\frac{1}{u}\right) < Q.$$

The probability distribution of the random variable u is given by (cf. Feller, [1])

$$(8) \quad P(u) = N^{-m} \binom{N}{u} \Delta^u 0^m,$$

where the s th difference of 0^t is defined by

$$\Delta^s 0^t = \sum_{r=0}^s (-1)^{s-r} \binom{s}{r} r^t.$$

Hence

$$(9) \quad E\left(\frac{1}{u}\right) = N^{-m} \sum_{u=1}^m \frac{1}{u} \binom{N}{u} \Delta^u 0^m$$

and the expected sample size is

$$(10) \quad E(u) = N^{-m} \sum_{u=1}^m u \binom{N}{u} \Delta^u \mathbf{0}^m.$$

We consider now Scheme B. Let b_1, b_2, \dots, b_n be the n distinct units selected in the sample and let k_r be the number of times the r th unit b_r occurs in the sample with $\sum k_r = v$. We compare the bias and variances of the estimate

$$(11) \quad \bar{y}_v = \frac{1}{v} \sum_{r=1}^n k_r y_r$$

and the estimate

$$(12) \quad \bar{y}_n = \frac{1}{n} \sum_{r=1}^n y_r$$

where y_r is the same as in Scheme A with obvious modifications. It is easy to see that the estimates \bar{y}_v and \bar{y}_n are unbiased for estimating \bar{Y} . Also we have

$$(13) \quad V(\bar{y}_v) = E \left(\frac{1}{v} \right) \frac{N-1}{N} \sigma^2$$

and

$$(14) \quad V(\bar{y}_n) = \left(\frac{1}{n} - \frac{1}{N} \right) \sigma^2.$$

The estimate \bar{y}_n is superior if

$$(15) \quad E \left(\frac{1}{v} \right) > \frac{1}{n} \frac{N-n}{N-1}.$$

The probability distribution of the random variable v can be shown to be

$$(16) \quad P(v) = \binom{N-1}{n-1} N^{1-v} \Delta^{n-1} \mathbf{0}^{v-1}.$$

Hence

$$(17) \quad E \left(\frac{1}{v} \right) = \binom{N-1}{n-1} \sum_{v=n}^{\infty} \frac{1}{v} N^{1-v} \Delta^{n-1} \mathbf{0}^{v-1}.$$

We give in Table 1 a numerical table for selected sample sizes which illustrates the numerical magnitudes of the differences discussed above. The theoretical proofs are given in Section 3 below.

3. Proofs of the inequalities. In order to establish inequality (7), we shall first prove the following

LEMMA. *Let*

$$(18) \quad S_{t,N} = \sum_{u=1}^{m-t} \frac{1}{u+t} \binom{N}{u} \Delta^u \mathbf{0}^{m-t}.$$

TABLE 1
 Values of $E(u)$, $E\left(\frac{1}{u}\right)$, Q and $1/E(u)$

N	m	$E(u)$	$E\left(\frac{1}{u}\right)$	Q	$\frac{1}{m}$	$\frac{1}{E(u)}$
$N = 10$	1	1.00	1.000	1.000	1.000	1.000
	2	1.90	.550	.550	.500	.526
	3	2.71	.385	.400	.333	.369
	4	3.44	.303	.325	.250	.291
	5	4.10	.253	.280	.200	.244
	6	4.69	.221	.250	.167	.213
$N = 50$	1	1.00	1.000	1.000	1.000	1.000
	2	1.98	.510	.510	.500	.505
	3	2.94	.343	.347	.333	.340
	4	3.88	.260	.265	.250	.258
	5	4.80	.210	.216	.200	.208
	6	5.71	.177	.183	.167	.175
$N = 100$	1	1.00	1.000	1.000	1.000	1.000
	2	1.99	.505	.505	.500	.503
	3	2.97	.338	.340	.333	.337
	4	3.94	.255	.258	.250	.254
	5	4.90	.205	.208	.200	.204
	6	5.85	.172	.175	.167	.171

Then

$$(19) \quad S_{t,N} \leq (N + 1)S_{t+1,N} \text{ for } t = 0, 1, \dots, m - 1; \quad m > 2,$$

the sign of equality holds only for $t = 0$.

PROOF.

$$\begin{aligned} S_{t,N} &= \sum_{u=1}^{m-t} \frac{u}{u+t} \binom{N}{u} (\Delta^{u-1} 0^{m-t-1} + \Delta^u 0^{m-t-1}) \\ &= \sum_{u=2}^{m-t} \frac{N-u+1}{u+t} \binom{N}{u-1} \Delta^{u-1} 0^{m-t-1} + \sum_{u=1}^{m-t-1} \frac{u}{u+t} \binom{N}{u} \Delta^u 0^{m-t-1} \\ &= \sum_{u=1}^{m-t-1} \left(\frac{N-u}{u+t+1} + \frac{u}{u+t} \right) \binom{N}{u} \Delta^u 0^{m-t-1} \\ &= \sum_{u=1}^{m-t-1} \frac{(N+1)u+tN}{(u+t)(u+t+1)} \binom{N}{u} \Delta^u 0^{m-t-1}. \end{aligned}$$

Thus $S_{0,N} = (N + 1)S_{1,N}$, and $S_{t,N} < (N + 1)S_{t+1,N}$ for $1 \leq t \leq m - 1$. This proves the lemma.

COROLLARY. Applying the lemma $m - 2$ times beginning with $t = 1$, we have, for $m > 2$

$$S_{1,N} < (N + 1)^{m-2} S_{m-1,N} = (N + 1)^{m-2} \frac{N}{m}$$

so that

$$(20) \quad S_{1,N-1} < N^{m-2} \frac{N-1}{m}.$$

Now, inequality (7) is proved in the following

THEOREM.

$$(21) \quad E\left(\frac{1}{u}\right) \leq \frac{1}{N} + \frac{N-1}{N} \frac{1}{m};$$

the sign of equality holds only for $m = 2$.

PROOF.

$$\begin{aligned} E\left(\frac{1}{u}\right) &= N^{-m} \sum_{u=1}^m \binom{N}{u} (\Delta^u 0^{m-1} + \Delta^{u-1} 0^{m-1}) \\ &= \frac{1}{N} + N^{1-m} \sum_2^m \frac{1}{u} \binom{N-1}{u-1} \Delta^{u-1} 0^{m-1} \\ &= \frac{1}{N} + N^{1-m} \sum_1^{m-1} \frac{1}{u+1} \binom{N-1}{u} \Delta^u 0^{m-1} \\ &= \frac{1}{N} + N^{1-m} S_{1,N-1}. \end{aligned}$$

On using (20) we have for $m > 2$,

$$E\left(\frac{1}{u}\right) < \frac{1}{N} + \frac{N-1}{N} \frac{1}{m}.$$

For $m = 2$, $S_{1,N-1} = \frac{1}{2}(N-1)$, so that the last inequality reduces to an equality in this case.

To prove inequality (15), we make use of the following

LEMMA. *If v is a positive random variable, then*

$$(22) \quad E\left(\frac{1}{v}\right) \geq \frac{1}{E(v)}.$$

The proof of this lemma follows from Cauchy's inequality (cf. Hardy and others, [3]),

$$(23) \quad (\sum a^2)(\sum b^2) > (\sum ab)^2,$$

by substituting $a = \sqrt{vP(v)}$, $b = \sqrt{v^{-1}P(v)}$, and noting that the two sides of the inequality are convergent because $E(v)$ and $E(1/v)$ are finite.

Now (cf. Feller, [2])

$$\begin{aligned} E(v) &= N \left(\frac{1}{N} + \frac{1}{N-1} + \dots + \frac{1}{N-n+1} \right), \\ \therefore E\left(\frac{1}{v}\right) &> \frac{1}{N \frac{n}{N-n+1}} = \frac{N-n+1}{nN} > \frac{N-n}{n(N-1)} \quad \text{for } n > 1. \end{aligned}$$

It is easy to see that for $n = 1$ the two procedures are equivalent and the inequality (15) reduces to an equality.

4. Estimation of variances. We shall consider the problem of estimating from the sample the variances of \bar{y}_u and \bar{y}_v . Under Scheme A, it is easy to see that for a given $u \geq 2$, an unbiased estimate of σ^2 is provided by

$$(24) \quad s_u^2 = \frac{1}{u - 1} \sum_{i=1}^u (y_i - \bar{y}_u)^2.$$

Thus considering

$$(25) \quad G_u = \left[\left(\frac{1}{u} - \frac{1}{N} \right) + N^{1-m} \left(1 - \frac{1}{u} \right) \right] s_u^2,$$

we have

$$(26) \quad E[G_u | u \geq 2] = V(\bar{y}_u);$$

so that G_u provides an unbiased estimate of the variance of \bar{y}_u . It is unbiased in the conditional sense, namely when the number of distinct units in the sample exceeds unity. An alternative unbiased estimate is provided by G'_u where

$$(27) \quad G'_u = \left[\left(\frac{1}{u} - \frac{1}{n} \right) + \frac{N - 1}{N^m - N} \right] s^2,$$

$$s^2 = \frac{1}{u - 1} \sum_{i=1}^u (y_i - \bar{y}_u)^2, \quad \text{for } u \geq 2,$$

and

$$s^2 = 0, \quad \text{for } u = 1.$$

Under Scheme B, it is easy to see that

$$(28) \quad E \left[\frac{1}{v} s_v^2 | v \right] = \frac{1}{v} \frac{N - 1}{N} \sigma^2$$

where

$$(29) \quad s_v^2 = \frac{1}{v - 1} \sum_{i=1}^v (y_i - \bar{y}_v)^2.$$

Thus

$$(30) \quad E \left[\frac{1}{v} s_v^2 \right] = E \left(\frac{1}{v} \right) \frac{N - 1}{N} \sigma^2,$$

so that $1/v s_v^2$ is an unbiased estimate of $V(\bar{y}_v)$.

5. Extension to ratio estimation. We shall now extend the above results to ratio estimates. We make use of the notation and approximate results given by Cochran [1]. The object is to estimate the population ratio $R = Y/X$. Under Scheme A we compare the two estimates

$$(31) \quad \bar{R}_m = \sum_{i=1}^u k_i y_i / \sum_{i=1}^u k_i x_i$$

and

$$(32) \quad \bar{R}_u = \sum_{i=1}^u y_i / \sum_{i=1}^u x_i.$$

We have

$$(33) \quad \begin{aligned} V(\bar{R}_m) &\doteq R^2 \left(Q - \frac{1}{N} \right) [C_{yy} + C_{xx} - 2C_{xy}] \\ B(\bar{R}_m) &\doteq R \left(Q - \frac{1}{N} \right) [C_{xx} - C_{xy}] \\ V(\bar{R}_u) &\doteq R^2 \left[E \left(\frac{1}{u} \right) - \frac{1}{N} \right] [C_{yy} + C_{xx} - 2C_{xy}] \\ B(\bar{R}_u) &\doteq R \left[E \left(\frac{1}{u} \right) - \frac{1}{N} \right] [C_{xx} - C_{xy}] \end{aligned}$$

where $B(\bar{R})$ stands for the absolute value of the bias of the estimate \bar{R} . Using inequality (7), we have

$$(34) \quad V(\bar{R}_u) < V(\bar{R}_m) \text{ and } B(\bar{R}_u) < B(\bar{R}_m).$$

Under Scheme B, the estimates to be compared are

$$(35) \quad \bar{R}_v = \sum_{i=1}^n k_i y_i / \sum_{i=1}^n k_i x_i$$

and

$$(36) \quad \bar{R}_n = \sum_{i=1}^n y_i / \sum_{i=1}^n x_i.$$

It is easy to see that the estimate \bar{R}_n is superior to \bar{R}_v from the point of view of variance and bias.

6. Extension to multistage designs. The result obtained for unistage designs will now be extended to multistage designs. Let a population consist of N first stage sampling units, of which m or n are selected with equal probabilities with replacement according to the Schemes A or B respectively. For the i th first stage unit, let t_i (based on sampling at second and subsequent stages) be an unbiased estimate of y_i , the total value of the character y for the unit. For Schemes A and B the unbiased estimates considered are

$$(37) \quad \bar{y}'_u = \frac{1}{u} \sum_{i=1}^n t_i,$$

$$(38) \quad \bar{y}'_m = \frac{1}{m} \sum_{i=1}^u k_i t_i,$$

and

$$(39) \quad \bar{y}'_n = \frac{1}{n} \sum_{i=1}^n t_i,$$

$$(40) \quad \bar{y}'_v = \frac{1}{v} \sum_{i=1}^n k_i t_i.$$

The variances of the estimates are given by

$$(41) \quad V(\bar{y}'_u) = \left[\left(1 + \frac{\delta^2}{\sigma^2} \right) E \left(\frac{1}{u} \right) - \frac{1}{N} \right] \sigma^2$$

$$(42) \quad V(\bar{y}'_m) = \left[\frac{\delta^2}{\sigma^2} Q + \frac{N-1}{mN} \right] \sigma^2$$

$$(43) \quad V(\bar{y}'_n) = \frac{\delta^2}{n} + \left(\frac{1}{n} - \frac{1}{N} \right) \sigma^2,$$

$$(44) \quad V(\bar{y}'_v) = \frac{N-1}{N} E \left(\frac{1}{v} \right) (\delta^2 + \sigma^2) + \frac{\delta^2}{N},$$

where

$$\delta^2 = \frac{1}{N} \sum_1^N V(t_i).$$

Using inequalities (7) and (15) it is found that \bar{y}'_u is superior to \bar{y}'_m while \bar{y}'_n is superior to \bar{y}'_v .

REFERENCES

- [1] W. G. COCHRAN, "*Sampling Techniques*", John Wiley and Sons, New York, 1953, pp. 117, 118.
- [2] W. FELLER, "*An introduction to probability theory and its applications*", John Wiley and Sons, New York, 1950, pp. 69, 175.
- [3] G. H. HARDY, J. E. LITTLEWOOD, AND G. POLYA, "*Inequalities*", Cambridge University Press, 1934, pp. 16, 124.