# ON AVERAGING OVER DISTINCT UNITS IN
# SAMPLING WITH REPLACEMENT

By R. M. Korwar and R. J. Serfling

*Florida State University*

**0. Summary.** Raj and Khamis [2], and Raj [3], have shown that in sampling with replacement the average over distinct units possesses lower variance than the average over the entire sample including repetitions. But no usable results are available regarding the amount of reduction in variance. The relevant expression given by Raj and Khamis is prohibitively cumbersome even for numerical computations. The present note provides the simplest possible exact expression as well as a convenient approximation. In particular, it follows readily that the relative benefit due to averaging over distinct units only is basically an increasing function of the sampling fraction.

**1. Background.** Let $u$ be the number of distinct units obtained in a sample of size $n$ drawn with replacement from a population of size $N$.

To estimate a mean, one may average over all units in the sample including repetitions, or alternatively average over the distinct units only. Each estimate is unbiased but the variances differ. For the first, the variance is proportional to

$$(1.1) \qquad Q - \frac{1}{N}$$

where

$$(1.2) \qquad Q = \frac{1}{n} + \frac{n-1}{nN} \ .$$

In the second case, the quantity $Q$ becomes replaced by $E(1/u)$. Also, in comparison of the mean square errors of certain ratio estimates, the quantities $Q$ and $E(1/u)$ again are involved.

It is thus important to determine $E(1/u)$ and to compare it with $Q$. Raj and Khamis [2] derive an expression for $E(1/u)$ and show that

$$(1.3) \qquad E(1/u) < Q, \qquad\qquad n \geq 3.$$

(For $n = 1$ and 2, $E(1/u) = Q$.) However, their expression is too complicated to yield readily any further theoretical insight, nor is it well suited to numerical computations. It is of interest to make a more refined comparison than that given by (1.3). For further background details, see [2] and [3].

**2. Results.** It will be shown that

$$(2.1) \qquad E(1/u) = N^{-n}\sum_{j=1}^{N} j^{n-1}, \qquad\qquad n \geq 1,$$

and that

$$(2.2) \qquad Q^* - \frac{1}{720N} < E(1/u) \leqq Q^*, \qquad\qquad 3 \leqq n \leqq N,$$

where

$$(2.3) \qquad Q^* = \frac{1}{n} + \frac{1}{2N} + \frac{n-1}{12N^2}.$$

(For $n = 3$, $E(1/u) = Q^*$.) Formula (2.1) affords the simplest possible exact representation for $E(1/u)$ and is amenable to numerical computations. In addition, $Q^*$ gives an excellent approximation, suitable not only for computations but also for theoretical considerations. For example, we have from (1.2) and (2.3) the approximation

$$(2.4) \qquad \frac{Q - 1/N}{Q^* - 1/N} \doteq \frac{1}{1 - \frac{1}{2}f},$$

where $f = n/N$, showing that the relative benefit due to averaging over distinct units only is an increasing function of $f$.

**3. Proofs.** Formula (2.1) may be obtained simply by appealing to the symmetry and independence inherent in the sampling scheme, without invoking the probability distribution of the random variable $u$. To this effect, let us denote the population units by $U_1, \cdots, U_N$. Let $i_1, \cdots, i_n$ be the subscripts of units selected in the sample and let $i_{(1)} < i_{(2)} < \cdots < i_{(u)}$ be their ordered distinct values. Then

$$(3.1) \qquad P[i_1 = i_{(u)}] = EP[i_1 = i_{(u)} \,|\, i_{(1)}, \cdots, i_{(u)}]$$
$$= E(1/u),$$

utilizing the symmetry. On the other hand,

$$P[i_1 = i_{(u)}] = P[i_1 \geqq i_2, i_1 \geqq i_3, \cdots, i_1 \geqq i_n]$$
$$(3.2) \qquad = \sum_{j=1}^{N} P[i_1 = j] P[i_2 \leqq j, \cdots, i_n \leqq j]$$
$$= \sum_{j=1}^{N} \left(\frac{1}{N}\right)\left(\frac{j}{N}\right)^{n-1},$$

again appealing to symmetry and also using the independence of $i_1, \cdots, i_n$. Thus (2.1) is proved. It may be of interest that (2.1) was first obtained by a combinatorial argument, the present argument being discovered after knowing the result.

To obtain (2.2), use will be made of the formulas

$$(3.3) \qquad \sum_{j=1}^{N} j^{n-1} = \frac{N^n}{n} + \frac{N^{n-1}}{2} + \sum_{j=0}^{L} (-1)^j \frac{1}{2(j+1)} \binom{n-1}{2j+1} B_{j+1} N^{n-2-2j},$$

where $L = (n-4)/2$ if $n$ is even and $= (n-3)/2$ if $n$ is odd, and

$$(3.4) \qquad \frac{2^{2r-1} \pi^{2r}}{(2r)!} B_r = \sum_{k=1}^{\infty} k^{-2r},$$

from Adams [1] pages 27 and 140. (The $B_r$ are the Bernoulli numbers.) By (2.1), (2.3) and (3.3), and since $B_1 = \frac{1}{6}$, we have, for $n \geq 3$,

$$(3.5) \qquad E(1/u) = Q^* + \sum_{j=0}^{L} (-1)^j A_j$$

where $A_0 = 0$ and

$$(3.6) \qquad A_j = N^{-2(j+1)} \frac{1}{2(j+1)} \binom{n-1}{2j+1} B_{j+1}, \qquad\qquad 1 \leq j \leq L.$$

Now, for $j \geq 1$, (3.6) and (3.4) yield

$$-A_{2j-1} + A_{2j} = N^{-4j} \frac{1}{4j} \binom{n-1}{4j-1} \left[ -B_{2j} + B_{2j+1} \frac{(2j)(n-4j+1)(n-4j)}{(2j+1)(4j)(4j+1)N^2} \right]$$

$$(3.7) \qquad\qquad < B_{2j} N^{-4j} \frac{1}{4j} \binom{n-1}{4j-1} \left[ -1 + \frac{(n-4j+1)(n-4j)}{4\pi^2 N^2} \right]$$

$$< -B_{2j} N^{-4j} \frac{1}{4j} \binom{n-1}{4j-1} \left( 1 - \frac{1}{4\pi^2} \right).$$

Hence $\sum_0^L (-1)^j A_j \leq 0$, so that $Q^*$ is an upper bound for $E(1/u)$. On the other hand, it may similarly be shown that

$$(3.8) \qquad A_{2j} - A_{2j+1} > B_{2j+1} N^{-4j-2} \frac{1}{4j+2} \binom{n-1}{4j+1} \left( 1 - \frac{1}{4\pi^2} \right),$$

whence $\sum_0^L (-1)^j A_j > -A_1$. Since $B_2 = 1/30$, we have $A_1 = \binom{n-1}{3}/(120N^4) < 1/720N$. Thus (2.2) is proved.

**Acknowledgment.** The authors are grateful to Dr. R. R. Davidson for encouragement in the investigation and to the referee for constructive and helpful suggestions on the presentation.

## REFERENCES

[1] ADAMS, EDWIN P. (1939). *Smithsonian Mathematical Formulae and Tables of Elliptic Functions.* Smithsonian Institution, Washington, D.C.

[2] RAJ, DES and KHAMIS, SALEM H. (1958). Some remarks on sampling with replacement. *Ann. Math. Statist.* **29** 550–557.

[3] RAJ, DES (1968). *Sampling Theory.* McGraw-Hill, New York.