

On inverse sampling with unequal probabilities†

By P. K. PATHAK‡

University of Illinois

SUMMARY

Sampford (1962) has considered the following sampling scheme to select a sample with n distinct population units. Sampling with unequal probabilities with replacement is carried out until $(n + 1)$ different population units are selected, the last unit is not recorded in the sample to insure some simplicity in the estimation procedure. The method is called inverse sampling with unequal probabilities. In this paper it is shown that this method of sampling is equivalent to sampling with unequal probabilities without replacement in some sense. The estimator of the population total given by Sampford is compared with other existing estimators. It is shown that there exist estimators which are uniformly better than the estimator given by Sampford. Finally, it is demonstrated that Des Raj's estimator (1956) has certain definite advantages over Sampford's estimator suggesting that sampling with unequal probabilities (without replacement) with Des Raj's estimator probably gives a better estimation procedure than inverse sampling with unequal probabilities.

1. INTRODUCTION AND PRELIMINARIES

Although sampling with unequal probabilities with replacement has the merit of simplicity, it has a major drawback that the number of distinct population units selected in a sample is a random variable and varies from sample to sample. Sampling with unequal probabilities without replacement does not have this drawback, but most estimators of the population total in this sampling scheme are either unwieldy or biased. Recently some workers have developed sampling schemes wherein the number of distinct population units selected in a sample is constant from sample to sample, and also the estimator of the population total is easy to compute and has some desirable properties like possessing a non-negative estimator of its variance, etc. Of special interest in this connexion are the works of Des Raj (1956), Stevens (1958), Rao, Hartley & Cochran (1962) and others. In cluster sampling, Sampford (1962) has considered a method of sampling with probabilities proportional to cluster size (with replacement) by which a fixed number of distinct clusters are included in the sample and has given an unbiased estimator of the population total and an estimator of its variance. A slightly different version of this method of sampling is considered here in detail and is referred to as inverse sampling with unequal probabilities.

The following notation and definitions are used in the subsequent sections.

Number of population units:	N .
j th population unit:	U_j ($j = 1, \dots, N$).
Value of a real-valued Y -characteristic of U_j :	Y_j ($j = 1, \dots, N$).
Probability of selection of U_j :	P_j ($\sum_{j=1}^N P_j = 1$).

† Work supported by the National Science Foundation.

‡ On leave from Indian Statistical Institute.

Z-value of U_j :	$Z_j = Y_j/P_j \quad (j = 1, \dots, N).$
Population total of Y -characteristic:	$Y = \sum_{j=1}^N Y_j.$
Population variance:	$\sigma_y^2 = \sum_{j=1}^N P_j(Z_j - Y)^2.$
Number of distinct population units in a sample:	$n.$
Size of a sample:	$r \quad (r = n, n+1, \dots).$
An observed sample:	$s = (u_1, \dots, u_r).$

Small letters will refer to the sample and capital to the population, e.g. u_1, u_2, \dots, u_r are the r sample units selected in order of draw.

DEFINITION 1-1. Let $u_{(1)}, u_{(2)}, \dots, u_{(n)}$ be the n distinct population units observed in a sample and arranged in order of their draw. Then $T_1 = [u_{(1)}, u_{(2)}, \dots, u_{(n)}]$ is called the 'serial-statistic'.

DEFINITION 1-2. Let $u_{(1)}, u_{(2)}, \dots, u_{(n)}$ be the n distinct population units observed in a sample and arranged in increasing order of their unit-indices. Then $T_1 = [u_{(1)}, u_{(2)}, \dots, u_{(n)}]$ is called the 'order-statistic'.

The unit index of an individual whose population representation is U_j is j . The above definition of 'order-statistic' differs from its customary definition where the sample observations are arranged in order of magnitude. This definition has been adopted here to avoid possible ambiguity in defining the 'order-statistic' when a sample may contain two different population units of the same order of magnitude.

DEFINITION 1-3. Let $\alpha_{(1)}, \alpha_{(2)}, \dots, \alpha_{(n)}$ be respectively the number of times $u_{(1)}, u_{(2)}, \dots, u_{(n)}$ occur in the sample. Then $T_2 = \{(u_{(1)}, \alpha_{(1)}), (u_{(2)}, \alpha_{(2)}), \dots, (u_{(n)}, \alpha_{(n)})\}$ is called the 'symmetric-statistic'.

Any statistic $T(s)$, a function of the sample observations, divides the collection of samples into sets such that the given statistic $T(s)$ is constant on each of these sets. This collection of sets is the partition induced by $T(s)$. Two samples are said to be equivalent if they contain the same population units. A sufficient statistic is thus defined as follows (Pathak, 1964).

DEFINITION 1-4. A statistic $T(s)$ is sufficient if the partition induced by $T(s)$ consists of sets of equivalent samples.

It is thus obvious that the above defined statistics, T_1 , T_2 and T_3 , are all sufficient.

For sampling schemes from finite populations, the above definition is a suitable version of the existing notion of sufficiency as defined by Fisher (1921), Halmos & Savage (1949), Bahadur (1954) and others. The following intuitive argument shows that this definition is equivalent to the definition of sufficiency in Fisher sense. Under a given sampling process, a selected sample gives us information about those population units selected in the sample. If two samples consist of the same population units they give the same information about the population units and are equally informative. Given one of these samples one can predict what the other sample is and vice versa. In the same way a sample from a given set of equivalent samples is as informative as any other sample from that set; and it is theoretically possible to work out all samples of this set given any one sample from it. A partition of samples into sets of equivalent samples divides the sample space into sets of

equally informative samples and is, therefore, sufficient. It may be noted that the sufficiency is with regard to the unknown population characteristics Y_1, Y_2, \dots, Y_N .

It is possible to give a rigorous justification of this definition in abstract terms but this has been omitted here for brevity. The author (1964) has dealt with this problem at some length in a paper to be published in *Annals of Mathematical Statistics* soon.

DEFINITION 1-5. Given two estimators $t_1(s)$ and $t_2(s)$ of a population parameter, $t_1(s)$ is said to be uniformly better than $t_2(s)$ if for any convex loss function $l_1(s)$ has smaller expected loss than $l_2(s)$ for all Y_1, \dots, Y_N with strictly less expected loss for at least one Y_1, \dots, Y_N .

It follows as a consequence of the Rao-Blackwell theorem (Pathak, 1964) that if $T(s)$ is a sufficient statistic and if $t_1(s)$ is an estimator of a parameter which does not depend on $T(s)$ then $t_1(s) = E[t_1(s)|T(s)]$ is a uniformly better estimator than $t_2(s)$ unless $t_1(s) = t_2(s)$ with probability one. If the loss function is the squared error then $t_1(s)$ has smaller mean-square error than $t_2(s)$; the decrease in the mean-square error is equal to $E[t_1(s) - t_2(s)]^2$.

DEFINITION 1-6. Random variables u_1, u_2, \dots, u_r are called interchangeable if their joint distribution is invariant under any permutation of u_1, u_2, \dots, u_r .

2. INVERSE SAMPLING WITH UNEQUAL PROBABILITIES

In inverse sampling with unequal probabilities, population units are selected with unequal probabilities (with replacement), P_j being the probability of selection of U_j ($j = 1, \dots, N$), and the selection of units is stopped at the $(r + 1)$ th draw when the sample first contains $(n + 1)$ different population units. The last unit is rejected and the recorded sample consists of sample units corresponding to n different population units. The rejection of the last sample unit introduces simplicity in the distribution of the sample units as then for a given sample size the recorded sample units behave as interchangeable random variables. An observed sample is recorded as

$$s = (u_1, u_2, \dots, u_r) \quad (r = n, n + 1, \dots), \quad (1)$$

where u_1, u_2, \dots, u_r are respectively the first, second and the r th sample units.

The probability of selecting a sample s is given by

$$P(s) = kf(s)p_1 p_2 \dots p_r (1 - p_{(1)} - p_{(2)} - \dots - p_{(n)}), \quad (2)$$

where p_i denotes the probability of selection of the i th sample unit, u_i ($i = 1, \dots, r$),

$$f(s) = \begin{cases} 1 & \text{if } s \text{ contains } n \text{ distinct population units,} \\ 0 & \text{otherwise,} \end{cases}$$

$p_{(1)}, \dots, p_{(n)}$ are the probabilities of selection associated with n different population units selected in the sample, and k is to be determined such that the summation of the right side over all possible samples equals unity.

From (2) it is evident that for a given r, u_1, u_2, \dots, u_r are interchangeable random variables.

Under this sampling scheme, the order-statistic, the serial statistic and the symmetric-statistic are all sufficient. The probability distribution of the symmetric-statistic

$$T_s = [(u_{(1)}, \alpha_{(1)}), \dots, (u_{(n)}, \alpha_{(n)})]$$

is easily seen to be given by

$$P[T_s] = kf(s) \frac{r!}{\alpha_{(1)}! \dots \alpha_{(n)}!} p_{(1)}^{\alpha_{(1)}} \dots p_{(n)}^{\alpha_{(n)}} (1 - p_{(1)} - \dots - p_{(n)}), \quad (3)$$

where k and $f(s)$ have been defined in (2) and p_{i0} is the probability of selection of u_{i0} ($i = 1, \dots, n$).

A little consideration will show that the probability distribution of the serial-statistic $T_2 = \{u_{i0}, u_{i1}, \dots, u_{in}\}$ is given by

$$P\{T_2\} = p_{10} \frac{p_{11}}{(1-p_{10})} \cdots \frac{p_{n1}}{(1-p_{10}-\dots-p_{(n-1)})}, \quad (4)$$

where p_{i0} is the probability of selection of u_{i0} ($i = 1, \dots, n$).

Lastly the probability distribution of the order-statistic $T_1 = \{u_{(1)}, \dots, u_{(n)}\}$ is given by

$$P\{T_1\} = \Sigma' p_{10} \frac{p_{10}}{(1-p_{10})} \cdots \frac{p_{n0}}{(1-p_{10}-\dots-p_{(n-1)})}, \quad (5)$$

where the summation Σ' extends over all possible orderings of p_{10}, \dots, p_{n0} .

It can be seen from (4) that the distribution of the serial-statistic is the same as that of the serial-statistic in sampling with unequal probabilities (without replacement). Thus in this sense the two sampling schemes are equivalent. The author (1961) has earlier proved a similar result. The equivalence can be made precise as follows. Suppose $g_1(s)$ and $g_2(s)$ are estimators of the population total based on inverse sampling with unequal probabilities and sampling with unequal probabilities (without replacement) respectively. Now if to the statistician the outcome of only one of these sampling schemes is given, he could, if he wished, compute with the help of a random device both $g_1(s)$ and $g_2(s)$ such that their probability distributions will be identical with their original probability distributions under the two given sampling schemes. Thus theoretically, the knowledge of the outcome of one of these sampling schemes enables the statistician to produce an outcome of the second sampling scheme and therefore the two sampling schemes can be said as equivalent. This also shows that in inverse sampling with unequal probabilities if one is using estimators which are based on the serial-statistic then sampling can be stopped as soon as the n th population unit has been selected.

3. ESTIMATION OF THE POPULATION TOTAL

Sampford (1962) has considered the following estimator of the population total

$$\bar{z} = \frac{1}{r} \sum_{i=1}^r z_i, \quad (6)$$

where z_i is the Z -value of u_i , and gave

$$v(\bar{z}) = \frac{1}{r(r-1)} \sum_{i=1}^r (z_i - \bar{z})^2 \quad (7)$$

as an unbiased estimator of $V(\bar{z})$.

Since $E(z_i) = Y$, and for a given r , z_1, \dots, z_r are interchangeable random variables,

$$E(\bar{z}) = E_r E(\bar{z}|r) = E_r E(z_1|r) = E(z_1) = Y. \quad (8)$$

This proves that \bar{z} is an unbiased estimator of the population total Y .

To prove the unbiasedness of $v(\bar{z})$ it suffices to observe that

$$\begin{aligned} E[v(\bar{z})] &= E\left[\bar{z}^2 - \frac{1}{r(r-1)} \sum_{i \neq j=1}^r z_i z_j\right] \\ &= E(\bar{z}^2) - E_r E\left[\frac{1}{r(r-1)} \sum_{i \neq j=1}^r z_i z_j | r\right] \\ &= E(\bar{z}^2) - E_r E(z_1 z_2 | r) \\ &= E(\bar{z}^2) - E(z_1 z_2) = E(\bar{z}^2) - Y^2 \\ &= V(\bar{z}). \end{aligned} \quad (9)$$

Further if c_1, c_2, \dots, c_r are r variables depending on r only and such that $\sum_{i=1}^r c_i = 1$, then it is easily verified in a similar manner that $\sum_{i=1}^r c_i z_i$ is also an unbiased estimator of Y . It is proved below that \bar{z} is the best estimator among this class of estimators.

THEOREM 1. \bar{z} is uniformly better than any estimator of the above type.

Proof. Since the symmetric statistic

$$T_3 = \{(u_{(1)}, \alpha_{(1)}), \dots, (u_{(n)}, \alpha_{(n)})\}$$

is sufficient, it follows that an estimator uniformly better than $\sum_{i=1}^r c_i z_i$ is given by

$$E \left[\sum_{i=1}^r c_i z_i | T_3 \right].$$

Further, when T_3 is given (r is fixed), c_1, c_2, \dots, c_r are fixed and also z_1, z_2, \dots, z_r are interchangeable. Therefore

$$\begin{aligned} E \left[\sum_{i=1}^r c_i z_i | T_3 \right] &= E \left[\left(\sum_{i=1}^r c_i \right) z_1 | T_3 \right] = E[z_1 | T_3] \\ &= \sum_{i=1}^n z_{(i)} P[u_1 = u_{(i)} | T_3]. \end{aligned} \quad (10)$$

It can be easily verified that

$$P[u_1 = u_{(i)} | T_3] = \frac{P[u_1 = u_{(i)} \cap T_3]}{P[T_3]} = \frac{\alpha_{(i)}}{r}, \quad (11)$$

so that

$$E[\sum_{i=1}^r c_i z_i | T_3] = \frac{1}{r} \sum_{i=1}^n \alpha_{(i)} z_{(i)} = \bar{z}. \quad (12)$$

This completes the proof of the theorem.

COROLLARY 1.1. On taking $c_1 = c_2 = \dots = c_n = 1/n$ and $c_{n+1} = \dots = c_r = 0$, it is seen that \bar{z} is uniformly better than the corresponding estimator $\bar{z}_n = \frac{1}{n} \sum_{i=1}^n z_i$ of sampling with unequal probabilities (with replacement).

4. VARIANCE OF \bar{z}

From (9), it is clear that

$$\begin{aligned} V(\bar{z}) &= E \left[\frac{1}{r(r-1)} \sum_{i=1}^r (z_i - \bar{z})^2 \right] \\ &= E \left[\frac{1}{2r(r-1)} \sum_{i \neq i'=1}^r (z_i - z_{i'})^2 \right] \\ &= E_r E \left[\frac{1}{2r^2(r-1)} \sum_{i \neq i'=1}^r (z_i - z_{i'})^2 | r \right] \\ &= E_r E \left[\frac{1}{2r} (z_1 - z_2)^2 | r \right], \end{aligned} \quad (13)$$

since for a given r , z_1, \dots, z_r are interchangeable. Thus

$$\begin{aligned} V(\bar{z}) &= E \left[\frac{(z_1 - z_2)^2}{2r} \right] \\ &= \sum_{j \neq j'=1}^N \frac{1}{r} P_j P_{j'} (Z_j - Z_{j'})^2 C_{jj'}, \end{aligned} \quad (14)$$

where

$$C_{jj'} = E[1/r | u_1 = U_j, u_2 = U_{j'}] \quad (j \neq j' = 1, \dots, N).$$

In the particular case when $n = 2$

$$C_{jF} = -\frac{(1 - P_j - P_j)}{(P_j + P_j)^2} [\log(1 - P_j - P_j) + P_j + P_j] \quad (j \neq j' = 1, \dots, N). \quad (15)$$

An exact expression for C_{jF} is derived in the Appendix given at the end of the paper. When $P_j = 1/N$ ($j = 1, \dots, N$), $V(\bar{z})$ can be expressed as (Sampford)

$$V(\bar{z}) = \sigma_z^2 \binom{N-2}{n-1} \sum_{k=1}^n (-)^{n-k} \binom{n-1}{k-1} \left(\frac{1}{k} - \frac{N(k-1)}{k^2} \log \left(1 - \frac{k}{N} \right) \right). \quad (16)$$

An asymptotic expression for $V(\bar{z})$ in this case in ascending powers of $1/N$ is given by

$$V(\bar{z}) = \frac{1}{n} \sigma_z^2 \left\{ 1 - \frac{1}{2N} \left[n - \frac{2}{n+1} \right] - \frac{1}{12N^2} \left[n^3 + \frac{15n^2 + 2n - 24}{(n+1)(n+2)} \right] \dots \right\}. \quad (17)$$

5. COMPARISON OF SAMPFORD'S ESTIMATOR WITH DES RAJ'S ESTIMATOR

If the serial-statistic is recorded then it is seen that Des Raj's estimator (1956) is obtained from

$$\left. \begin{aligned} t_1 &= \frac{y_{11}}{p_{11}}, \\ t_2 &= y_{11} + \frac{y_{12}}{p_{12}} (1 - p_{11}), \\ &\dots \dots \dots \\ t_i &= y_{11} + \dots + y_{(i-1)} + \frac{y_{i1}}{p_{i1}} (1 - p_{11} - \dots - p_{(i-1)}), \\ &\dots \dots \dots \\ t_n &= y_{11} + \dots + y_{(n-1)} + \frac{y_{n1}}{p_{n1}} (1 - p_{11} - \dots - p_{(n-1)}). \end{aligned} \right\} \quad (18)$$

Using the fact that t_1, \dots, t_n are uncorrelated, Des Raj considered

$$\bar{l} = \frac{1}{n} \sum_{i=1}^n t_i \quad (19)$$

as an unbiased estimator of the population total and gave

$$v(\bar{l}) = \frac{1}{n(n-1)} \sum_{i=1}^n (t_i - \bar{l})^2 \quad (20)$$

as an unbiased estimator of $V(\bar{l})$.

The above estimator of the population total has an advantage over \bar{z} that it does not take into account the number of times a particular unit u_{11} is included in the sample whereas the Sampford's estimator does. Des Raj's estimator, however, takes into account the order in which the distinct population units are included in the sample and Sampford's estimator does not. This is an advantage of Sampford's estimator over Des Raj's. It is felt that perhaps the advantage of Des Raj's estimator is more desirable than that of Sampford's. A still more desirable estimator than Des Raj's and Sampford's would be one which disregards both the order and the number of times different population units are included in the sample (Theorem 2).

It seems rather difficult to give a direct comparison of the relative efficiencies of \bar{l} and

\bar{z} in the general case. However, when $n = 2$, and if the loss function is the squared error, $V(\bar{l})$ is given by

$$\begin{aligned} V(\bar{l}) &= E\{v(\bar{l})\} = E\left\{\frac{1}{2}(t_1 - t_2)^2\right\} \\ &= \frac{1}{2}E\left\{\left(\frac{y_{(1)}}{p_{(1)}} - \frac{y_{(2)}}{p_{(2)}}\right)^2 (1 - p_{(1)})^2\right\} \\ &= \frac{1}{8} \sum_{j \neq j'=1}^N (Z_j - Z_{j'})^2 P_j P_{j'} (2 - P_j - P_{j'}). \end{aligned} \quad (21)$$

From this, it is obvious that $V(\bar{l}) < V(\bar{z})$ provided

$$2 - P_j - P_{j'} < -4 \frac{(1 - P_j - P_{j'})}{(P_j + P_{j'})^2} \left\{ \log \{ (1 - P_j - P_{j'}) + P_j + P_{j'} \} \right\} \quad (22)$$

for all $j \neq j' = 1, \dots, N$.

It can be seen that this will be so if $(P_j + P_{j'}) < \frac{1}{2}$ for all $j \neq j' = 1, \dots, N$. Thus while sampling from a reasonably big population where this condition is automatically satisfied Des Raj's estimator will have smaller variance.

In another special case when $P_j = 1/N$ ($j = 1, \dots, N$), $V(\bar{l})$ is given (Murthy, 1957) by

$$V(\bar{l}) = \frac{\sigma^2}{n} \left\{ 1 - \frac{(n-1)}{N} + \frac{(n-1)(n-2)}{3N(N-1)} \right\}. \quad (23)$$

A comparison of (23) with the asymptotic expression for the variance of \bar{z} shows that \bar{l} has smaller variance than \bar{z} .

The above considerations lead the author to believe that Des Raj's estimator will probably be better than Sampford's estimator in practice.

6. AN ESTIMATOR UNIFORMLY BETTER THAN SAMPFORD'S ESTIMATOR

The theorem below gives an estimator uniformly better than Sampford's estimator.

THEOREM 2. An estimator uniformly better than \bar{z} is given by

$$\bar{z}^* = \sum_{i=1}^n c_{i0} y_{(i)} \quad (24)$$

where $c_{i0} = \{P(T_1 | u_1 = u_{(i)})\} / P(T_1)$, $P(T_1)$ is the probability of getting the order-statistic $T_1 = [u_{(1)}, \dots, u_{(n)}]$ and $P(T_1 | u_1 = u_{(i)})$ is the conditional probability of getting the order-statistic, T_1 , given that $u_{(i)}$ has been selected as the first draw.

Proof. Since T_1 is sufficient, an estimator uniformly better than \bar{z} is given by

$$\bar{z}^* = E \left[\frac{1}{r} \sum_{i=1}^r z_i | T_1 \right] = E_r E \left[\frac{1}{r} \sum_{i=1}^r z_i | T_1, r \right]. \quad (25)$$

Evidently z_1, \dots, z_r are interchangeable for a given r and T_1 , and therefore

$$\begin{aligned} \bar{z}^* &= E_r E[z_1 | T_1, r] = E[z_1 | T_1] \\ &= \sum_{i=1}^n \frac{y_{(i)}}{p_{(i)}} P(u_1 = u_{(i)} | T_1) \\ &= \sum_{i=1}^n c_{i0} y_{(i)}. \end{aligned} \quad (26)$$

This proves the theorem.

The author (1961) had earlier proved that \bar{z}^* is uniformly better than Des Raj's estimator. Murthy (1967) proved that \bar{z}^* has smaller variance than Des Raj's estimator. This shows that \bar{z}^* is uniformly better than Des Raj's as well as Sampford's estimator. Unfortunately \bar{z}^* cannot be of much use in practice because of the cumbersome computation of the coefficients α_{ij} . It can be seen that $\alpha_{ij} = N/n$ when $P_j = 1/N$ ($j = 1, \dots, N$) with the help of (5) showing that in simple random sampling (without replacement) the estimator based on the sample mean is uniformly better than both Des Raj's and Sampford's estimators.

Other estimators uniformly better than Sampford's estimator can be suggested, as for example, $E[\bar{z}|T_1, r]$, $E[\bar{z}|T_2, r]$, etc. The reader may find it an illustrative exercise to work out these estimators. The author's (1962) paper may be referred to for reference. It will be interesting to see how the above considered methods of estimation compare in two-state sampling.

The author wishes to thank a referee for some valuable suggestions.

APPENDIX

The probability distribution of r . The probability of getting the order-statistic $T_1 = \{u_{(1)}, \dots, u_{(n)}\}$ for a given r , is given by (Pathak, 1962)

$$P(T_1, r) = \{p_{(1)} + \dots + p_{(n)}\}^r - \Sigma_i \{p_{(1)} + \dots + p_{(n-i)}\}^r + \dots + (-1)^{n-1} \Sigma_i p_{(1)}^r [1 - p_{(1)} - \dots - p_{(n)}], \quad (27)$$

where the summation Σ_i is taken over all combinations of p 's chosen out of $p_{(1)}, \dots, p_{(n)}$. Therefore, the probability distribution of r is given by

$$P(r) = \Sigma_i \{ \{p_{(1)} + \dots + p_{(n)}\}^r - \Sigma_i \{p_{(1)} + \dots + p_{(n-i)}\}^r + \dots + (-1)^{n-1} \Sigma_i p_{(1)}^r [1 - p_{(1)} - \dots - p_{(n)}] \}, \quad (28)$$

where Σ_i is taken over all $\binom{N}{n}$ combinations of p 's chosen out of P_1, P_2, \dots, P_N .

From (27) and (28), it can be seen that

$$P(T_1, r | u_1 = u_{(1)}, u_2 = u_{(2)}) = \{ \{p_{(1)} + \dots + p_{(n)}\}^{r-2} - \Sigma_i \{p_{(1)} + \dots + p_{(n-i)}\}^{r-2} + \dots + (-1)^{n-2} \{p_{(1)} + p_{(2)}\}^{r-2} [1 - p_{(1)} - \dots - p_{(n)}] \}, \quad (29)$$

where Σ_i stands for all combinations of p 's containing $p_{(1)}$ and $p_{(2)}$ chosen out of $p_{(1)}, \dots, p_{(n)}$.

From (29), it follows that

$$\begin{aligned} E[1/r | u_1 = u_{(1)}, u_2 = u_{(2)}, T_1] &= \left[\frac{1}{\{p_{(1)} + \dots + p_{(n)}\}^2} \left(\log \frac{1}{1 - p_{(1)} - \dots - p_{(n)}} - p_{(1)} - \dots - p_{(n)} \right) \right. \\ &\quad - \Sigma_i \left(\frac{1}{\{p_{(1)} + \dots + p_{(n-i)}\}^2} \left(\log \frac{1}{1 - p_{(1)} - \dots - p_{(n-i)}} - p_{(1)} - \dots - p_{(n-i)} \right) \right. \\ &\quad \left. \left. + \dots + (-1)^{n-2} \frac{1}{\{p_{(1)} + p_{(2)}\}^2} \left(\log \frac{1}{1 - p_{(1)} - p_{(2)}} - p_{(1)} - p_{(2)} \right) \right] \right. \\ &\quad \left. \times [1 - p_{(1)} - \dots - p_{(n)}] = C_{1/r}(T_1). \end{aligned} \quad (30)$$

Finally, $C_{1/r}$, as required in (14), is given by

$$C_{1/r} = \Sigma_i' C_{1/r}(T_1), \quad (31)$$

where Σ_i' is taken over all $\binom{N-2}{n-2}$ combinations of n P 's chosen out of P_1, P_2, \dots, P_N and containing P_1 and P_2 .

REFERENCES

- BARADUR, R. R. (1954). Sufficiency and statistical decision function. *Ann. Math. Statist.* **25**, 423-62.
 DES RAJ (1956). Some estimators in sampling with varying probabilities without replacement. *J. Amer. Statist. Ass.* **51**, 260-84.
 FURBER, R. A. (1921). On the mathematical foundations of theoretical statistics. *Phil. Trans. A*, **222**, 309-88.
 HALMOS, P. R. & SAVAGE, L. J. (1949). Applications of the Radon-Nikodym theorem to the theory of sufficient statistics. *Ann. Math. Statist.* **20**, 226-41.

- MURPHY, M. N. (1957). Ordered and unordered estimators in sampling without replacement. *Sankhyā*, 18, 379-90.
- PATNAK, P. K. (1961). Use of 'order-statistic' in sampling without replacement. *Sankhyā*, A, 23, 409-14.
- PATNAK, P. K. (1962). On sampling with unequal probabilities. *Sankhyā*, A, 24, 316-28.
- PATNAK, P. K. (1964). Sufficiency in sampling theory (to be published in *Ann. Math. Statist.*).
- RAO, J. N. K., HARTLEY, H. O. & COCHRAN, W. G. (1962). On a simple procedure of unequal probability sampling without replacement. *J. R. Statist. Soc. B*, 24, 482-91.
- SAMFORD, M. R. (1962). Methods of cluster sampling with and without replacement for clusters of unequal sizes. *Biometrika*, 49, 27-40.
- STEVENS, W. L. (1963). Sampling without replacement with probability proportional to size. *J. R. Statist. Soc. B*, 20, 393-7.