

# ON SAMPLING WITH UNEQUAL PROBABILITIES

By P. K. PATHAK

Indian Statistical Institute

**SUMMARY.** This paper deals with the problem of deriving improved estimators in sampling schemes with unequal probabilities of selection. The improved estimator of the population total,  $Y$ , (Basu, 1958), is derived. In addition, two sets of estimators of  $Y$  and  $Y^2$  are given. The first set of estimators is unwieldy to compute, while the second set is simple. The second set of estimators, though less efficient than the first, is more efficient than the usually employed estimators.

It is proved in subfield terminology that if  $\mathcal{S}_1$  and  $\mathcal{S}_2$  are two sufficient subfields and  $K$  is a set common to  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , then  $\mathcal{S}_1K + \mathcal{S}_2K'$  is also a sufficient subfield. Hence the subfield  $\mathcal{S}_1K + \mathcal{S}_2K'$  can be used to derive improved estimators by Rao-Blackwell theorem. Generalisation of this is also given in case of countable number of subfields. Application of this result to sampling with unequal probabilities is given.

## 1. INTRODUCTION

Consider a population containing  $N$  units. Let  $y_j$  be some real-valued characteristic of the  $j$ -th population unit in which we are interested.<sup>1</sup> Suppose that a sample of size  $n$  is drawn from the above population with unequal probabilities of selection (with replacement),  $P_j$  being the probability of selection associated with the  $j$ -th population unit ( $\sum P_j = 1$ ). If for the  $i$ -th sample unit, we record its  $Y$ -characteristic  $y_i$ , probability of selection  $p_i$  and unit-index  $u_i$ , then the sample of observation is

$$S = (x_1, x_2, \dots, x_n),$$

where  $x_i = (y_i, p_i, u_i)$ .<sup>2</sup>

It has been shown by Basu (1958) that the 'order-statistic'

$$T = (x_{(1)}, x_{(2)}, \dots, x_{(r)})$$

(where  $x_{(1)}, x_{(2)}, \dots, x_{(r)}$ , are the distinct units in the sample arranged in ascending order of their unit-indices) is sufficient.

Therefore, if  $g(S)$  is some estimator depending on the sample  $S$ , for any convex (downwards) loss function, an estimator uniformly better than  $g(S)$  is given by  $E[g(S)|T]$ . In the subsequent sections this result is used to derive improved estimators of the population total and its square.

## 2. ESTIMATION OF THE POPULATION TOTAL

The usual estimator of the population total

$$Y = \sum Y_j$$

is given by 
$$\bar{y} = \frac{1}{n} \sum z_i, \quad \dots \quad (2.1)$$

where 
$$z_i = \frac{y_i}{p_i}.$$

<sup>1</sup>  $j$  varies from 1 to  $N$ ,  $i$  from 1 to  $n$  and (i) from (1) to  $(r)$  unless otherwise stated.

<sup>2</sup> Capital letters refer to the population and small letters to the sample.

Theorem 1 : For any convex (downwards) loss function, an estimator uniformly better than  $\bar{z}$  is given by

$$\bar{z}_* = E[\bar{z}|T] = \Sigma C_{(i)} \frac{y_{(i)}}{p_{(i)}} \quad \dots (2.2)$$

$$\text{where } C_{(i)} = \frac{p_{(i)}[(p_{(1)} + \dots + p_{(v)})^{n-1} - \Sigma\{(p_{(1)} + \dots + p_{(r-1)})^{n-1} + \dots + (-)^{r-1} p_{(1)}^{n-1}\}]}{[(p_{(1)} + \dots + p_{(v)})^n - \Sigma\{(p_{(1)} + \dots + p_{(r-1)})^n + \dots + (-)^{r-1} \Sigma_1 p_{(1)}^n\}]} \quad \dots (2.3)$$

the summations  $\Sigma_1$  and  $\Sigma_1^{\#}$  stand for all combinations of  $p$ 's and all combinations of  $p$ 's containing  $p_{(i)}$  (chosen out of  $p_{(1)}, p_{(2)}, \dots, p_{(v)}$ ) respectively.

Proof: Obviously by Rao-Blackwell theorem, an estimator uniformly better than  $\bar{z}$  is given by

$$E(\bar{z}|T) = E\left(\frac{y_1}{p_1} \middle| T\right) = \Sigma \frac{y_{(i)}}{p_{(i)}} P(x_1 = z_{(i)}|T). \quad \dots (2.4)$$

$$\text{But } P(x_1 = z_{(i)}|T) = \frac{p_{(i)} \Sigma' \frac{(n-1)!}{\alpha_{(1)}! \dots \alpha_{(v)}!} p_{(1)}^{\alpha_{(1)}} \dots p_{(v)}^{\alpha_{(v)}}}{\Sigma' \frac{n!}{\alpha_{(1)}! \dots \alpha_{(v)}!} p_{(1)}^{\alpha_{(1)}} \dots p_{(v)}^{\alpha_{(v)}}} \quad \dots (2.5)$$

where  $\Sigma'$  means summation over all integral  $\alpha$ 's such that

$$\alpha_{(i)} > 0 \text{ for } i = 1, \dots, v \text{ and } \alpha_{(1)} + \alpha_{(2)} + \dots + \alpha_{(v)} = n,$$

and  $\Sigma^*$  means summation over all integral  $\alpha$ 's such that

$$\alpha_{(i)} \geq 0, \alpha_{(v)} > 0 \text{ for } i' \neq i = 1, 2, \dots, v$$

and

$$\alpha_{(1)} + \alpha_{(2)} + \dots + \alpha_{(v)} = (n-1).$$

It can be seen by induction over  $v$  that

$$\Sigma' \frac{n!}{\alpha_{(1)}! \dots \alpha_{(v)}!} p_{(1)}^{\alpha_{(1)}} \dots p_{(v)}^{\alpha_{(v)}} = [(p_{(1)} + \dots + p_{(v)})^n - \Sigma_1 (p_{(1)} + \dots + p_{(v-1)})^n + \dots + (-)^{v-1} \Sigma_1 p_{(1)}^n];$$

and

$$\Sigma^* \frac{(n-1)!}{\alpha_{(1)}! \dots \alpha_{(v)}!} p_{(1)}^{\alpha_{(1)}} \dots p_{(v)}^{\alpha_{(v)}} = [(p_{(1)} + \dots + p_{(v)})^{n-1} - \Sigma\{(p_{(1)} + \dots + p_{(v-1)})^{n-1} + \dots + (-)^{v-1} p_{(1)}^{n-1}\}]. \quad \dots (2.6)$$

Using (2.4), (2.5) and (2.6), we get

$$\bar{z}_* = E(\bar{z}|T) = \Sigma C_{(i)} \frac{y_{(i)}}{p_{(i)}}$$

Hence the theorem is proved.

ON SAMPLING WITH UNEQUAL PROBABILITIES

The above estimator, though better, is not very useful in large samples because of cumbersome computation of  $C_{(i)}$ 's. In Section 4, we shall derive a simpler estimator of  $Y$  better than  $\bar{y}$ . Table 1 gives exact expressions of  $\bar{z}_p$  for  $n = 3, 4$  and  $5$ . For  $n = 1$  and  $2$ ,  $\bar{z}_p$  and  $\bar{y}$  are identical.

TABLE 1.  $\bar{z}_p$  FOR  $n=3, 4$  AND  $5$

$n \rightarrow$ $p$	3	4	5
1	$\frac{y_{(1)}}{p_{(1)}}$	$\frac{y_{(1)}}{p_{(1)}}$	$\frac{y_{(1)}}{p_{(1)}}$
2	$\frac{\Sigma (2p_{(1)} + p_{(2)}) \frac{y_{(1)}}{p_{(1)}}}{3(p_{(1)} + p_{(2)})}$	$\frac{\Sigma [(p_{(1)} + p_{(2)})^2 - p_{(1)}^2] y_{(1)}}{[(p_{(1)} + p_{(2)})^2 - p_{(1)}^2 - p_{(2)}^2]}$	$\frac{\Sigma [(p_{(1)} + p_{(2)})^2 - p_{(1)}^2] y_{(1)}}{[(p_{(1)} + p_{(2)})^2 - p_{(1)}^2 - p_{(2)}^2]}$
3	$\frac{\Sigma \frac{y_{(1)}}{p_{(1)}}}{p_{(1)}}$	$\frac{\Sigma [2p_{(1)} + p_{(2)} + p_{(3)}] \frac{y_{(1)}}{p_{(1)}}}{4(p_{(1)} + p_{(2)} + p_{(3)})}$	$\frac{\Sigma [12p_{(1)}(p_{(1)} + p_{(2)} + p_{(3)}) + 4(p_{(2)}^2 + p_{(3)}^2) + 6p_{(2)}p_{(3)}] \frac{y_{(1)}}{p_{(1)}}}{5 [4(p_{(2)}^2 + p_{(3)}^2 + p_{(2)}p_{(3)}) + 8(p_{(1)}p_{(2)} + p_{(1)}p_{(3)} + p_{(2)}p_{(3)})]}$
4	—	$\frac{\Sigma \frac{y_{(1)}}{p_{(1)}}}{p_{(1)}}$	$\frac{\Sigma [2p_{(1)} + p_{(2)} + p_{(3)} + p_{(4)}] \frac{y_{(1)}}{p_{(1)}}}{5 [p_{(1)} + p_{(2)} + p_{(3)} + p_{(4)}]}$
5	—	—	$\frac{\Sigma \frac{y_{(1)}}{p_{(1)}}}{p_{(1)}}$

3. ESTIMATION OF  $Y^2$

The problem of finding an unbiased estimator of  $Y^2$  arises in most problems of variance estimation of estimators of  $Y$ . The usual estimator of  $Y^2$  is

$$\bar{z}_p = \frac{1}{n(n-1)} \sum_{i \neq i'}^n z_i z_{i'} \quad \dots (3.1)$$

Theorem 2: For any convex (downward) loss function, an estimator uniformly better than  $\bar{z}_p$  is given by

$$E(\bar{z}_p | T) = \sum_{i=1}^p C_{(i,i)} z_{(i)}^2 + \sum_{i \neq i'}^p C_{(i,i')} z_{(i)} z_{(i')} \quad \dots (3.2)$$

Remark:  $\bar{z}_p$  when  $v=(n-1)$  may be expressed in a simple form as  $z_{(n-1)} = \frac{1}{n} \left\{ \frac{\Sigma y_{(i)}}{p_{(i)}} + \frac{\Sigma y_{(i)}}{\Sigma p_{(i)}} \right\}$

where  $C_{(i,i')} = \frac{p_{(i)}^2 [(p_{(1)} + \dots + p_{(r)})^{n-2} - \Sigma_1^i (p_{(1)} + \dots + p_{(r-1)})^{n-2} + \dots + (-)^{r-1} p_{(i)}^2]}{[(p_{(1)} + \dots + p_{(r)})^n - \Sigma_1 (p_{(1)} + \dots + p_{(r-1)})^n + \dots + (-)^{r-1} \Sigma_1 p_{(1)}^n]}$

and  $C_{(i,i')} =$

$$\frac{p_{(i)} p_{(i')} [(p_{(1)} + \dots + p_{(r)})^{n-2} - \Sigma_1^{i'} (p_{(1)} + \dots + p_{(r-1)})^{n-2} + \dots + (-)^{r-2} (p_{(i)} + p_{(i')})^{n-2}]}{[(p_{(1)} + \dots + p_{(r)})^n - \Sigma_1 (p_{(1)} + \dots + p_{(r-1)})^n + \dots + (-)^{r-1} \Sigma_1 p_{(1)}^n]} \dots (3.3)$$

the summations  $\Sigma_1$  and  $\Sigma_1^i$  have been defined in (2.3) and the summation  $\Sigma_1^{i'}$  stands for all combinations of  $p$ 's containing  $p_{(i)}$  and  $p_{(i')}$ .

Proof : Obviously

$$E(z_p | T) = E \left[ \frac{1}{n(n-1)} \sum_{i \neq i'=1}^n z_i z_{i'} | T \right] = E(z_1 z_2 | T) \\ = \Sigma z_{(i)} z_{(i')} P[x_1 = x_{(i)}, x_2 = x_{(i')} | T]. \dots (3.4)$$

It is easy to see that

$$P[x_1 = x_{(i)}, x_2 = x_{(i')} | T] = \frac{p_{(i)}^2 \Sigma^{i'} \frac{(n-2)!}{\alpha_{(1)}! \dots \alpha_{(r)}!} p_{(i)}^{\alpha_{(1)}} \dots p_{(r)}^{\alpha_{(r)}}}{\Sigma^{i'} \frac{n!}{\alpha_{(1)}! \dots \alpha_{(r)}!} p_{(i)}^{\alpha_{(1)}} \dots p_{(r)}^{\alpha_{(r)}}} \dots (3.5)$$

and  $P[x_1 = x_{(i)}, x_2 = x_{(i')} | T] = \frac{p_{(i)} p_{(i')} \Sigma^{i''} \frac{(n-2)!}{\alpha_{(1)}! \dots \alpha_{(r)}!} p_{(i)}^{\alpha_{(1)}} \dots p_{(r)}^{\alpha_{(r)}}}{\Sigma^{i''} \frac{n!}{p_{(1)}! \dots \alpha_{(r)}!} p_{(i)}^{\alpha_{(1)}} \dots p_{(r)}^{\alpha_{(r)}}}, \dots (3.5.1)$

where  $\Sigma^i, \Sigma^{i'}$  have been defined in (2.5) and  $\Sigma^{i''}$  means summation over all integral  $\alpha$ 's such that

$$\alpha_{(1)} + \dots + \alpha_{(r)} = (n-2), \alpha_{(i)} \geq 0, \alpha_{(i')} \geq 0 \text{ and } \alpha_{(k)} > 0 \text{ for } k \neq i \neq i' = 1, \dots, v.$$

It can be proved on lines similar to (2.8) (by induction over  $v$ ) that

$$\left. \begin{aligned} P[x_{(1)} = x_{(i)}, x_{(2)} = x_{(i)} | T] &= c_{(i,i)} \\ \text{and } P[x_{(1)} = x_{(i)}, x_{(2)} = x_{(i')} | T] &= c_{(i,i')} \end{aligned} \right\} \dots (3.6)$$

Using (3.4) and (3.6), we get

$$E[z_p | T] = \sum_{i=1}^r C_{(i,i)} z_{(i)}^2 + \sum_{i \neq i'=1}^r C_{(i,i)} z_{(i)} z_{(i')} \dots (3.7)$$

which proves the theorem.

Improved estimator of  $\sigma_i^2$  : The usual estimator of  $\sigma_i^2 = \Sigma P_j \left( \frac{Y_j}{P_j} - Y \right)^2$  is

given by  $e_i^2 = \frac{1}{(n-1)} \Sigma (z_i - \bar{z})^2 = \frac{1}{2n(n-1)} \sum_{i \neq i'=1}^n (z_i - z_{i'})^2.$

## ON SAMPLING WITH UNEQUAL PROBABILITIES

Corollary 1 : Thus an estimator uniformly better than  $s_y^2$  is given by

$$E[s_y^2 | T] = E \left[ \frac{(z_1 - z_n)^2}{2} | T \right] = \sum_{i \neq l=1}^n O_{(i,l)} (z_{(i)} - z_{(l)})^2 \quad \dots (3.8)$$

Corollary 2 : An unbiased estimator of  $V(z_r)$  is given by

$$v(z_r) = s_r^2 - \sum_{i=1}^n O_{(i,i)} z_{(i)}^2 - \sum_{i \neq l=1}^n O_{(i,l)} z_{(i)} z_{(l)}. \quad \dots (3.9)$$

Since this estimator is quite complicated for use in large samples, Basu (1958) has suggested the use of

$$\frac{1}{n(n-1)} \sum (z_i - \bar{z})^2 \quad \dots (3.10)$$

as an estimator of  $V(z_r)$ . As it over-estimates  $V(z_r)$ , we are always on the safe side to use (3.10) as our estimate.

The estimators derived in this and preceding sections, though superior to the usually employed estimators, are not of much use for large scale sample surveys owing to their cumbersome nature. In the next section, we give simpler estimators of  $Y$  and  $Y^2$ . These estimators, though less efficient than the above derived estimators, are superior to the usually employed estimators.

### 4. SIMPLE IMPROVED ESTIMATORS OF $Y$ AND $Y^2$

Let us suppose that the observed samples are segregated into groups of equal  $p_i$ 's. For instance, consider the problem of estimating the total yield of a crop from a sample of farms. Every sample-farm is to be selected with probability proportional to its area. Here, if some crude approximation (say correct to an acre) is used to measure their areas, we expect to get number of farms with same  $p_i$  in the sample. In the sequel, by the  $p$ -value of a unit, we mean the probability of selection associated with that unit. Let  $p_{(1)}, p_{(2)}, \dots, p_{(k)}$  be the distinct  $p$ -values of the sample units arranged in an increasing order of their magnitude. Let  $n_{(i)}$  be the number of sample units having  $p_{(i)}$  as their  $p$ -value. However, not all these  $n_{(i)}$  units will be distinct, let  $v_{(i)}$  be the number of distinct units among them. Now, if we arrange these  $v_{(i)}$  distinct units in an increasing order of their unit-indices and call them  $x_{(1)}, x_{(2)}, \dots, x_{(v_{(1)})}$ , then it is not difficult to see that the statistic

$$T^* = \{[x_{(11)}, \dots, x_{(1v_{(1)})}; n_{(1)}], \dots, [x_{(k1)}, \dots, x_{(kv_{(k)})}; n_{(k)}]\} \quad \dots (4.1)$$

is sufficient.

It should be noted that if we take away the ancillary statistics  $n_{(1)}, \dots, n_{(k)}$  from the sufficient statistic  $T^*$ , then it reduces to the 'order-statistic'  $T$  defined in the earlier section. The 'unnecessarily wide' sufficient statistic  $T^*$  is used here for the purpose of deriving estimators of  $Y$  and  $Y^2$  that are much simpler (though somewhat less efficient) than those considered in the previous section. Theorems 3 and 4 below give simple improved estimators of  $Y$  and  $Y^2$  respectively.

Theorem 3 : For any convex (downwards) loss function, an estimator uniformly better than  $\bar{x}$  is given by

$$x_r^* = \frac{1}{n} \sum_{i=1}^k \frac{n_{(i)}}{p_{(i)}} y_{r_{(i)}}, \quad \dots (4.2)$$

where

$$\bar{y}_{r_{(i)}} = \frac{1}{v_{(i)}} \sum_{r=1}^{v_{(i)}} y_{(r)}.$$

*Proof:* Evidently an estimator uniformly better than  $\bar{x}$  is given by

$$E(\bar{x} | T^*) = E \left( \frac{y_1}{p_1} | T^* \right) \quad \dots (4.3)$$

Further, the probability of getting a sample with a given  $T^*$  is

$$P(T^*) = \frac{n!}{n_{(1)}! \dots n_{(k)}!} p_{(1)}^{n_{(1)}} \dots p_{(k)}^{n_{(k)}} C_{v_{(1)}}(n_{(1)}) \dots C_{v_{(k)}}(n_{(k)}), \quad \dots (4.4)$$

where  $C_{v_{(i)}}(n_{(i)}) = v_{(i)}^{n_{(i)}} - \binom{v_{(i)}}{1} (v_{(i)}-1)^{n_{(i)}} + \dots + (-1)^{v_{(i)}-1} \binom{v_{(i)}}{v_{(i)}-1} 1^{n_{(i)}}$ ;  
( $i = 1, \dots, k$ )

and

$$\begin{aligned} P(x_1 = x_{(r)} | T^*) &= \frac{p_{(1)}^{n_{(1)}} \dots p_{(i-1)}^{n_{(i-1)}} p_{(i)}^{n_{(i)}-1} \dots p_{(k)}^{n_{(k)}}}{\frac{n!}{n_{(1)}! \dots n_{(i-1)}! \dots n_{(i)}! \dots n_{(k)}!} p_{(1)}^{n_{(1)}} \dots p_{(i)}^{n_{(i)}} \dots p_{(k)}^{n_{(k)}}} \\ &\times \frac{C_{r_{(1)}}(n_{(1)}) \dots C_{r_{(i)}}(n_{(i)}) \dots C_{r_{(k)}}(n_{(k)})}{C_{v_{(1)}}(n_{(1)}) \dots C_{v_{(i)}}(n_{(i)}) \dots C_{v_{(k)}}(n_{(k)})} \\ &= \frac{n_{(i)}}{n} \cdot \frac{1}{v_{(i)}}. \quad \dots (4.5) \end{aligned}$$

From (4.3) and (4.5), it follows that

$$E(\bar{x} | T^*) = \frac{1}{n} \sum_{i=1}^k \frac{n_{(i)}}{p_{(i)}} \bar{y}_{r_{(i)}}$$

which completes the proof of the theorem.

A simple comparison of  $x_r^*$  and  $\bar{x}$  will show that  $x_r^*$  will be superior to  $\bar{x}$  if and only if the sample size is greater than two and at least three population units have the same  $p$ -value, otherwise  $x_r^*$  and  $\bar{x}$  will be identical. It is not difficult to give a direct proof of the fact that  $V(x_r^*) \leq V(\bar{x})$ . The strict sign of inequality holds only when the above condition is satisfied.

Theorem 4: For any convex (downwards) loss function, an estimator uniformly better than

$$z_p = \frac{1}{n(n-1)} \sum_{i \neq i'-1} \frac{y_i}{p_i} \cdot \frac{y_{i'}}{p_{i'}}$$

is given by

$$z_p^* = \frac{1}{n(n-1)} \left[ \left\{ \left( \sum_{i=1}^k n_{(i)} \frac{g_{v(i)}}{p_{(i)}} \right)^2 - \sum_{i=1}^k n_{(i)} \frac{g_{v(i)}^2}{p_{(i)}^2} \right\} - \sum n_{(i)}(n_{(i)}-1) \frac{O_{v(i)-1}(n_{(i)}-1)}{c_{v(i)}(n_{(i)})} \cdot \frac{g_{v(i)}^2}{p_{(i)}^2} \right], \dots (4.8)$$

where  $e_{v(i)}^2 = \frac{1}{(v_{(i)}-1)} \sum_{r=1}^{v_{(i)}} (y_{v(i)} - g_{v(i)})^2$ , and  $O_{v(i)}(n_{(i)})$

and  $C_{v(i)-1}(n_{(i)}-1)$  have meaning similar to those defined in (4.4).

Proof: Obviously, an estimator uniformly better than  $z_p$  is given by

$$E[z_p | T^*] = E \left[ \frac{y_i}{p_i} \cdot \frac{y_{i'}}{p_{i'}} | T^* \right]. \dots (4.7)$$

Further, it can be shown that

$$P[x_1 = x_{(i)}, x_2 = x_{(i')} | T^*] = \frac{n_{(i)}(n_{(i)}-1)}{n(n-1)} \cdot \frac{O_{v(i)}(n_{(i)}-1)}{v_{(i)} C_{v(i)}(n_{(i)})},$$

$$P[x_1 = x_{(i)}, x_2 = x_{(i'')} | T^*] = \frac{n_{(i)}(n_{(i)}-1)}{n(n-1)} \times \left[ \frac{C_{v(i)}(n_{(i)}) - C_{v(i)}(n_{(i)}-1)}{v_{(i)}(v_{(i)}-1) C_{v(i)}(n_{(i)})} \right], (v \neq v')$$

and  $P[x_1 = x_{(i)}, x_2 = x_{(i'')} | T^*] = \frac{n_{(i)}n_{(i'')}}{n(n-1)} \cdot \frac{1}{v_{(i)}} \cdot \frac{1}{v_{(i')}} (i \neq i''). \dots (4.8)$

Therefore,

$$\begin{aligned} E[z_p | T^*] &= \sum_{i=1}^k \frac{n_{(i)}(n_{(i)}-1)}{n(n-1)} \sum_{r=1}^{v_{(i)}} \frac{y_{v(i)}^2}{p_{(i)}^2} \frac{O_{v(i)}(n_{(i)}-1)}{v_{(i)} C_{v(i)}(n_{(i)})} \\ &+ \sum_{i=1}^k \frac{n_{(i)}(n_{(i)}-1)}{n(n-1)} \sum_{r \neq r'-1}^{v_{(i)}} \frac{y_{v(i)} y_{v(i')}}{p_{(i)}^2} \frac{[O_{v(i)}(n_{(i)}) - C_{v(i)}(n_{(i)}-1)]}{v_{(i)}(v_{(i)}-1) C_{v(i)}(n_{(i)})} \\ &+ \sum_{i \neq i'-1}^k \frac{n_{(i)}n_{(i')}}{n(n-1)} \sum_{r=1}^{v_{(i)}} \sum_{r'=1}^{v_{(i')}} \frac{y_{v(i)} y_{v(i')}}{p_{(i)} p_{(i')}} \cdot \frac{1}{v_{(i)}} \cdot \frac{1}{v_{(i')}}. \dots (4.9) \end{aligned}$$

Using the equality

$$\begin{aligned} \frac{O_{v(i)-1}(n_{(i)}-1)}{C_{v(i)}(n_{(i)})} e_{v(i)}^2 &= g_{v(i)}^2 - \frac{C_{v(i)}(n_{(i)}-1)}{v_{(i)} C_{v(i)}(n_{(i)})} \sum_{r=1}^{v_{(i)}} y_{v(i)}^2 \\ &- \frac{[C_{v(i)}(n_{(i)}) - C_{v(i)}(n_{(i)}-1)]}{v_{(i)}(v_{(i)}-1) C_{v(i)}(n_{(i)})} \sum_{r \neq r'-1}^{v_{(i)}} y_{v(i)} y_{v(i')} \end{aligned}$$

and simplifying (4.9) we get

$$z_p^* = E[z_p | T^*] = \frac{1}{n(n-1)} \left[ \left\{ \left( \sum_{i=1}^k n_{(i)} \frac{\bar{y}_{v_{(i)}}}{p_{(i)}} \right)^2 - \sum_{i=1}^k n_{(i)} \frac{\bar{y}_{v_{(i)}}^2}{p_{(i)}} \right\} - \sum_{i=1}^k n_{(i)}(n_{(i)}-1) \cdot \frac{C_{v_{(i)}-1}(n_{(i)}-1)}{C_{v_{(i)}}(n_{(i)})} \cdot \frac{s_{v_{(i)}}^2}{p_{(i)}} \right].$$

This completes the proof.\*

Corollary 3 : It is easy to see that

$$\begin{aligned} E[s_1^2 | T^*] &= E \left[ \frac{1}{(n-1)} \sum (z_i - \bar{z})^2 | T^* \right] = E \left[ \frac{y_1^2}{p_1^2} | T^* \right] - E \left[ \frac{y_1}{p_1} \cdot \frac{y_2}{p_2} | T^* \right] \\ &= \sum_{i=1}^k \frac{n_{(i)}}{n} \cdot \frac{1}{v_{(i)}} \sum_{r=1}^{v_{(i)}} \frac{y_{(i)r}^2}{p_{(i)}} + \sum_{i=1}^k \frac{n_{(i)}(n_{(i)}-1)}{n(n-1)} \cdot \frac{C_{v_{(i)}-1}(n_{(i)}-1)}{C_{v_{(i)}}(n_{(i)})} \cdot \frac{s_{v_{(i)}}^2}{p_{(i)}} \\ &\quad - \frac{1}{n(n-1)} \left[ \left( \sum_{i=1}^k n_{(i)} \frac{\bar{y}_{v_{(i)}}}{p_{(i)}} \right)^2 - \sum_{i=1}^k n_{(i)} \frac{\bar{y}_{v_{(i)}}^2}{p_{(i)}} \right] \quad \dots \quad (4.10) \end{aligned}$$

is a simple improved estimator of  $s_1^2$ .

Corollary 4 : An unbiased estimator of  $V(z_p^*)$  is given by

$$\begin{aligned} v(z_p^*) &= z_p^{*2} + \sum_{i=1}^k \frac{n_{(i)}(n_{(i)}-1)}{n(n-1)} \cdot \frac{C_{v_{(i)}-1}(n_{(i)}-1)}{C_{v_{(i)}}(n_{(i)})} \cdot \frac{s_{v_{(i)}}^2}{p_{(i)}} \\ &\quad - \frac{1}{n(n-1)} \left[ \left( \sum_{i=1}^k n_{(i)} \frac{\bar{y}_{v_{(i)}}}{p_{(i)}} \right)^2 - \sum_{i=1}^k n_{(i)} \frac{\bar{y}_{v_{(i)}}^2}{p_{(i)}} \right]. \quad \dots \quad (4.11) \end{aligned}$$

However, in practice it seems reasonable to use

$$\frac{1}{n(n-1)} \sum \left( \frac{z_i}{p_i} - \bar{z} \right)^2 \quad \dots \quad (4.12)$$

as an estimator of  $V(z_p^*)$ . First, because it is simple to compute, secondly, because it is always non-negative. Besides this, we are always on the safe side as it over-estimates the variance of  $z_p^*$ .

\* The estimator  $Z_p$  requires the computation of the ratio  $\frac{C_{r-1}(n-1)}{C_r(n)}$ . Values of  $\frac{C_{r-1}(n-1)}{C_r(n)}$  can be obtained from the relation  $\frac{C_{r-1}(n-1)}{C_r(n)} = \frac{1}{r} - \frac{C_r(n-1)}{C_r(n)}$ .

Values of  $\frac{C_r(n-1)}{C_r(n)}$  have been tabulated for all  $r$  and  $n=1$  to 50 in a paper published elsewhere in the same issue (Pathak, 1962), pp. 287-303.



ON SAMPLING WITH UNEQUAL PROBABILITIES

5. A REMARK ON RAO-BLACKWELL THEOREM

Let  $X$  be the sample space of all possible outcomes  $x$ . Let  $\mathfrak{G}$  be the field of subsets of  $X$  on which a set  $P = \{p\}$  of probability measures is defined. Any statistic  $T(x)$  [an  $\mathfrak{G}$ -measurable function defined from  $X$  onto another space  $Y$ ] generates a subfield  $\mathfrak{G}_T \subseteq \mathfrak{G}$ . The statistic  $T$  or the subfield  $\mathfrak{G}_T$  is called sufficient for  $P$  (Bahadur, 1954) if corresponding to each  $\mathfrak{G}$ -measurable set  $A$ , there exists an  $\mathfrak{G}_T$ - $P$ -integrable function  $\phi_A$  such that

$$P\{x \in A_0 \cap A\} = \int_{A_0 \cap A} dp = \int_{A_0} \phi_A(x) dp \text{ for } A_0 \in \mathfrak{G}_T, p \in P. \quad \dots (5.1)$$

It is known that any estimator  $g(x)$  based on the sample  $x$  can be uniformly improved by taking the conditional expectation of  $g(x)$  given a sufficient subfield. If several sufficient subfields are available, the minimum condensation of  $g(x)$  is mostly obtained by employing the minimal sufficient subfield ( $\mathfrak{G}_3$ , say). There are situations, e.g., in sample surveys where this minimal condensation is unwieldy and it is not possible for practical reasons to use this condensation. But it is sometimes possible to divide  $X$  into subsets  $K(\in \mathfrak{S}_1)$  and  $K'$  such that the condensation is simple on  $K$  and unwieldy on  $K'$ . In such cases simpler condensation can be achieved with the help of some other subfield,  $\mathfrak{S}_2$ , which contains the minimal one. It follows as a consequence of Theorem 5 (stated below) that condensation smaller than that of  $\mathfrak{S}_3$  (but larger than that of  $\mathfrak{S}_1$ ) can be obtained by employing  $\mathfrak{S}_1 K + \mathfrak{S}_2 K'$  as the sufficient subfield, and this condensation will still have the merit of simplicity.

Theorem 5: Let  $\mathfrak{S}_1$  and  $\mathfrak{S}_2$  be two sufficient subfields of  $(X, \mathfrak{G}, P)$ , and  $K$  a set common to  $\mathfrak{S}_1$  and  $\mathfrak{S}_2$ . Then the subfield\*

$$\mathfrak{S}_3 = \mathfrak{S}_1 K + \mathfrak{S}_2 K' \quad \dots (5.2)$$

is also sufficient.

Proof: Since  $\mathfrak{S}_1$  and  $\mathfrak{S}_2$  are two sufficient subfields, there exist for each  $\mathfrak{G}$ -measurable set  $A$ , an  $\mathfrak{S}_1$ - $P$ -integrable function  $\phi_{1A}$  and an  $\mathfrak{S}_2$ - $P$ -integrable function  $\phi_{2A}$  such that

$$\int_{A_i \cap A} dp = \int_{A_i} \phi_{iA}(x) dp \text{ for } A_i \in \mathfrak{S}_i \text{ (} i = 1, 2 \text{), } p \in P. \quad \dots (5.3)$$

Now for any  $A_3 = A_1 K + A_2 K' \in \mathfrak{S}_3$ , and for each  $A \in \mathfrak{G}$

$$\int_{A_3 \cap A} dp = \int_{A_1} \phi_{3A} dp. \quad \dots (5.4)$$

where  $\phi_{3A} = \phi_{1A} \chi_K + \phi_{2A} \chi_{K'}$ , and  $\chi_K = 1 - \chi_{K'}$  is the characteristic function of the set  $K$ .

It is easily seen that  $\phi_{1A} \chi_K$  and  $\phi_{2A} \chi_{K'}$  are both  $\mathfrak{S}_3$ - $P$ -integrable and thus  $\phi_{3A}$  is  $\mathfrak{S}_3$ - $P$ -integrable. This with (5.4) implies that  $\mathfrak{S}_3$  is a sufficient subfield.

Corollary 5: If  $\mathfrak{S}_1 \subseteq \mathfrak{S}_2$ , then  $\mathfrak{S}_1 \subseteq \mathfrak{S}_3 \subseteq \mathfrak{S}_2$ .

\*  $\mathfrak{S}_3$  is the field of sets of the form  $\mathfrak{S}_3 = \{A_1 K + A_2 K'\}$ ,  $A_i \in \mathfrak{S}_i$  ( $i=1, 2$ ).

Corollary 6 : Let  $\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_k$  be  $k$  sufficient subfields. If  $A_1, A_2, \dots, A_k$  [ $\bigcup_{i=1}^k A_i = X, A_i \cap A_j = \phi$  ( $i \neq j = 1, \dots, k$ )] are  $k$  sets such that  $A_i \in \mathcal{S}_i$  ( $i = 1, \dots, k$ ), then

$$\mathcal{S}^* = A_1\mathcal{S}_1 + A_2\mathcal{S}_2 + \dots + A_k\mathcal{S}_k \quad \dots (5.5)$$

is also a sufficient subfield.

Corollary 7 : Let  $g(x)$  be an estimator based on  $x$ . Let  $\mathcal{S}_1$  and  $\mathcal{S}_2$  be two sufficient subfields and  $K$  a set common to  $\mathcal{S}_1$  and  $\mathcal{S}_2$ . Then for any convex (downward) loss function, an estimator uniformly better than  $g(x)$  is given by

$$E[g(x)|\mathcal{S}_2] = \begin{cases} E[g(x)|\mathcal{S}_1] & \text{if } x \in K \\ E[g(x)|\mathcal{S}_2] & \text{otherwise,} \end{cases} \quad \dots (5.6)$$

where  $\mathcal{S}_3 = \mathcal{S}_1K + {}_xK'$ .

This result is useful in estimation problems when the improved estimator  $E[g(x)|\mathcal{S}_1]$  is difficult to compute for all  $x \in X$ . In such cases, this result may be utilized by employing subfield  $\mathcal{S}_1$  and another subfield  $\mathcal{S}_2$  such that  $E[g(x)|\mathcal{S}_1]$  is simple to compute when  $x \in K$  and  $E[g(x)|\mathcal{S}_2]$  is simple to compute when  $x \in K'$  ( $K \in \mathcal{S}_1, i = 1, 2$ ) and using (5.6) as the improved estimator. The resulting estimator will still be better than  $g(x)$  and in addition will be simple to compute. Further if  $\mathcal{S}_1 \subseteq \mathcal{S}_2$ , this estimator will also be better than  $E[g(x)|\mathcal{S}_2]$ . For completeness a straightforward generalisation of Theorem 5 to the case of countable number of sufficient subfields is given below.

Theorem 6 : Let  $\{\mathcal{S}_i\}$  be a countable number of sufficient subfields. Let  $\{K_i\}$  be a sequence of mutually exclusive and exhaustive sets such that  $K_i \in \mathcal{S}_i$  ( $i = 1, 2, \dots$ ). Then

$$\mathcal{S}^* = \sum_{i=1}^{\infty} K_i\mathcal{S}_i$$

is a sufficient subfield.

*Proof:* Since  $\mathcal{S}_i$  is a sufficient subfield, there exists for each  $A \in \mathcal{S}$  an  $\mathcal{S}_i$ - $P$ -integrable function  $\phi_{iA}$  such that

$$P[x \in A_i \cap A] = \int_{A_i \cap A} dp = \int_{A_i} \phi_{iA}(x) dp \quad \text{for } A_i \in \mathcal{S}_i \quad (i = 1, \dots, \infty), p \in P. \quad \dots (5.7)$$

Now for any  $A^* = \sum A_i K_i \in \mathcal{S}^*$  and for each  $A \in \mathcal{S}$ , we have from (5.7)

$$\begin{aligned} P[x \in A^* \cap A] &= \lim_{n \rightarrow \infty} P[x \in (\sum_{i=1}^n A_i K_i) \cap A] \\ &= \lim_{n \rightarrow \infty} \int_{(\sum_{i=1}^n A_i K_i)} \left[ \sum_{i=1}^n \phi_{iA}(x) \chi_{K_i}(x) \right] dp \\ &= \lim_{n \rightarrow \infty} \int_{A^*} \left[ \sum_{i=1}^n \phi_{iA}(x) \chi_{K_i}(x) \right] dp. \quad \dots (5.8) \end{aligned}$$

## ON SAMPLING WITH UNEQUAL PROBABILITIES

Since  $0 < f_n = \left( \sum_{i=1}^n \phi_{iA}(x) \chi_{K_i}(x) \right) < 1$  a.e. (P),  $\lim_{n \rightarrow \infty} f_n = f_A^*$  exists a.e. (P). We, thus, have by Lebesgue's monotone convergence theorem

$$P[x \in A^* \cap A] = \int_A f_A^*(x) dP. \quad \dots (5.9)$$

Since for each  $n$ ,  $f_n$  is  $S^*$ - $P$ -integrable,  $f_A^*(x)$  is  $S^*$ - $P$ -integrable. Hence the theorem.

In the following section Theorem 5 is used in sampling with unequal probabilities to derive some simple improved estimators of  $Y$ .

### 6. APPLICATION TO SAMPLING WITH UNEQUAL PROBABILITIES

We have seen that  $\bar{z}_v$  is uniformly better than  $\bar{z}$ . For  $n = 3$ ,  $\bar{z}_v$  can be expressed as

$$\bar{z}_v = \begin{cases} z_{(1)} & \text{if } v = 1, \\ \frac{1}{3} \left[ z_{(1)} + z_{(2)} + \frac{y_{(1)} + y_{(2)}}{p_{(1)} + p_{(2)}} \right] & \text{if } v = 2 \\ \frac{1}{3} [z_{(1)} + z_{(2)} + z_{(3)}] & \text{if } v = 3. \end{cases} \quad \dots (6.1)$$

It is not simple to compute  $\bar{z}_v$  when  $n > 3$  owing to cumbersome computation of  $C_{(i)}$ 's. However, if in a sample of size  $n$ ,  $v = (n-1)$ ,  $\bar{z}_v$  is expressible in the simple form

$$\bar{z}_{(n-1)} = \frac{1}{n} \left[ \sum_{i=1}^{(n-1)} \frac{y_{(i)}}{p_{(i)}} + \frac{\sum_{i=1}^{(n-1)} y_{(i)}}{\sum_{i=1}^{(n-1)} p_{(i)}} \right]. \quad \dots (6.2)$$

As a direct consequence of Corollary 7, it follows that a simple estimator uniformly better than  $\bar{z}_v^*$  (and hence better than  $\bar{z}$ ) is given by

$$\bar{z}'_v = \begin{cases} \bar{z}_v & \text{if } v = (n-1) \\ \bar{z}_v^* & \text{otherwise.} \end{cases} \quad \dots (6.3)$$

Two points in favour of utilising  $\bar{z}'_v$  are: (i) it is as simple as  $\bar{z}$  or  $\bar{z}_v^*$  and (ii) it is more efficient than  $\bar{z}_v^*$ .

*Another simple improved estimator of  $Y$ .* Another simple improved estimator of  $Y$  can be derived by using the following sufficient statistic

$$T_2 = [(x_{(1)}, \alpha_{(1)}), \dots, (x_{(v)}, \alpha_{(v)})] \quad \dots (6.4)$$

where

$$\alpha_{(i)} = \begin{cases} \lambda_{(i)} & \text{if } \lambda_{(i)} > 2 \\ 1 & \text{otherwise } (i = 1, \dots, v) \end{cases}$$

and  $\lambda_{(i)}$  is the number of times  $x_{(i)}$  is included in the sample.

Assuming without any loss of generality that  $\alpha_{(i)} = 1$  if  $i = 1, \dots, k$  and  $\alpha_{(i)} > 1$  if  $i = k+1, \dots, v$ , and  $\sum_{i=1}^k \lambda_{(i)} = m$ , it can be shown that an estimator better than  $\bar{z}$  is given by

$$\bar{z}_{(k)} = E[\bar{z} | T_2] = \frac{1}{n} \left[ \sum_{i=1}^k \alpha_{(i)} z_{(i)} + \sum_{i=1}^k \hat{\alpha}_{(i)} z_{(i)} \right], \quad \dots (6.5)$$

where  $\bar{d}_{(i)} = \frac{\sum_1^i p_{(1)} \cdots p_{(m-k)}}{\sum_1^i p_{(1)} \cdots p_{(m-k)}}$  the summations  $\Sigma_1$  and  $\Sigma_i$  have been defined in (2.3)

and are taken over  $p_{(1)}, \dots, p_{(k)}$ .

In practical situations, it is much simpler to compute this estimator than to compute  $\bar{d}_r$ . For  $m-k=1$  and  $m-k=2$ ,  $\bar{d}_{r_{(k)}}$  is given by

$$\bar{d}_{r_{(k)}} = \begin{cases} \frac{1}{n} \left[ \sum_{i=1}^n \alpha_{(i)} z_{(i)} + \frac{\sum_{i=1}^k y_{(i)}}{\sum_{i=1}^k p_{(i)}} \right] & \text{if } m-k=1 \\ \frac{1}{n} \left[ \sum_{i=1}^n \alpha_{(i)} z_{(i)} + \frac{2k(y, p)}{k(p, p)} \right] & \text{if } m-k=2 \end{cases} \quad \dots (6.6)$$

where  $k(y, p) = \sum_{i=1}^k y_{(i)} p_{(i)} - (\sum y_{(i)}) (\sum p_{(i)})$  and  $k(p, p)$  is defined similarly.

In general when  $(m-k)$  is large, this estimator may also involve some extra computation. If the statistician is not even in favour of this extra computation, the author, as a consequence of Corollary 7, recommends the following improved procedure of estimation

- (i) use  $\bar{d}$  if  $(m-k) > 2$  ... (6.7)  
 (ii) use  $E[\bar{z} | T_k]$  if  $(m-k) \leq 2$ .

For estimating the variance of these estimators, author suggests (4.12) as an estimator.

#### 7. CONCLUDING REMARK

In case of large samples if one is interested in altogether dispensing with the extra computation, the observed sample of size  $n$  may be divided into sub-samples of sizes  $n_1, n_2, \dots, n_k$  etc., ( $\sum n_i = n$  and  $n_i = 3, 4$  or  $5$  etc.). This division should, however, be independent of sample observations. Each sub-sample may then be treated as a sample in itself and simple improved estimators may be obtained for each sub-sample by using the estimators given in the preceding section. The over-all improved estimator can now be obtained by averaging the estimators obtained from each sub-sample with weights proportional to the sub-sample sizes.

#### ACKNOWLEDGMENT

The author wishes to express his gratefulness to Dr. D. Basu for his guidance and encouragement in writing the paper.

#### REFERENCES

- BARADUJ, R. R. (1954): Sufficiency and statistical decision function. *Ann. Math. Stat.*, 26, 423-462.  
 BASU, D. (1958): On sampling with and without replacement. *Sankhyā*, 20, 287-294.  
 FRASER, D. A. S. (1957): *Nonparametric Methods in Statistics*, John Wiley and Sons, New York.  
 PATHAK, P. K. (1962): On simple random sampling with replacement. *Sankhyā*, 24, 287-302.

Paper received: August, 1961.