Iowa State University

I. Introduction

Let y_j denote a characteristic attached to the jth unit of a finite population of N units with population total $Y = \Sigma y_j$. It is well known that the use of unequal probabilities in selecting a sample may bring about considerable reduction in variance as compared to equal probability sampling. For example, such a procedure may be useful when a 'measure of size' x; is known for all units in the population and it is suspected that these known sizes x_j are correlated with the characteristic y_j . One method lated with the characteristic yj. (though by no means the only method) of utilizing x_j is to draw units with probabilities proportional to sizes x_j (p.p.s.), a technique frequently used in surveys, particularly for selecting primary units in multistage designs. Now the theory of sampling with unequal probabilities is equivalent to multinomial sampling provided the units are drawn with replacement. On the other hand, we know from equal probability sampling that selection with replacement results in estimators which are less precise than those computed from samples selected without replacement, the proportional reduction in the variance being given by the sampling fraction n/N. It is, therefore, natural to investigate the efficiency of unequal probability sampling without replacement as compared to unequal probability sampling with replacement.

A general theory of unequal probability sampling without replacement is given by Horvitz and Thompson (1952). Their estimator of Y is

$$\hat{\mathbf{Y}} = \frac{\mathbf{n}}{\mathbf{1}} \mathbf{y}_{j} / \mathbf{P}_{j}$$
(1)

with variance

$$\mathbf{V}(\hat{\mathbf{Y}}) = \sum_{1}^{N} \frac{\mathbf{y}_{j}}{\mathbf{p}_{j}} + 2 \sum_{1' \ge 1}^{N} \frac{\mathbf{p}_{11'}}{\mathbf{p}_{1'}} \mathbf{y}_{1} \mathbf{y}_{1'} - \mathbf{y}^{2}, \quad (2)$$

where P_j is the probability for the jth unit to be in a sample of size n and P₁₁, is the probability for the units i and i' both to be in the sample.*, Now, when the P_j are proportional to the y_j, Y is constant and hence V(Y) is zero which suggests that, by making P_j proportional to the known sizes x_j, considerable reduction in the variance will result if the x_j are approximately proportional to the y_j. Special devices are needed to satisfy this condition, namely P_j = np_j where p_j = x_j/2x_j, when sampling with p.p.s. and without replacement. Yates and Grundy (1953) suggest an iterative procedure to obtain revised sizes x^{*}_j and draw the sample by selecting the first unit with probabilities propor-

tional to the revised sizes, the second unit with probabilities proportional to the remaining (revised sizes), and so on, such that $P_j = np_j$. Let us call their method procedure 1. One can also use the iterative procedure of obtaining revised sizes for another well known sampling scheme which is as follows: n units are drawn with probabilities proportional to revised sizes with replacement. If any unit is selected more than once in the sample, reject the n selections and make further n selections with replacement and with probabilities proportional to revised sizes, the process being continued until n different units are selected in the sample. Let us call this method procedure 2 (e. g., Durbin (1953)). In order to draw a sample by procedure 1 or 2 it is convenient to use a method suggested by Lahiri (1951).* The application of Lahiri's method leads to the following scheme of drawing the sample for procedure 1: Let us consider, for illustration, n=2. Let x_0^* be a number not smaller than the largest of x_1^* . Select a random integer between 1 and N, say α where $1 \le \alpha \le N$. Select similarly a random number β , subject to the condi-tion $0 < \beta \le x_0^*$. If $\beta \le x_0^*$, then unit α is se-lected for the sample. If $\beta > x_0^*$, unit α is not selected (at this 'draw') and the whole process has to be repeated with all the units until a selection is made. The probability of selecting a unit by this procedure is proportional to the revised sizes. After a unit is selected, repeat the whole process with the remaining x_j^* until another selection is made. For procedure 2, after a unit is selected by Lahiri's method, the whole process is repeated with all the x_j^* for the selection of the second unit. If the second unit selected happens to be the same as the first unit, reject both selections and repeat the whole process with all the units until two different selections are made.

There is, however, a well known procedure of drawing a sample with p.p.s. and without replacement which insures that $P_j = np_j$ with the <u>original</u> sizes x_j . This procedure, which we denote as procedure 3, is as follows: The N units in the population are listed in a <u>random</u> order and their x_j are cumulated and a systematic selection of n elements from a random start is then made on the cumulation (e.g., Goodman and Kish (1950)). It may be noted that for all the three procedures all x_j 's should satisfy the necessary condition $P_j = np_j \leq 1$. Therefore, if there are some units for which np > 1, either select them automatically in the sample or use some other device like stratification or splitting a unit into smaller units, etc.

Recently (Hartley and Rao (1959), (1962)), the mathematical difficulties involved in evaluating the probabilities P_{11} , for procedure 3, are resolved with the help of an asymptotic

[&]quot;It may be pointed out that in unequal probability sampling there are several different classes of linear estimators. The estimator (1) belongs to one of these classes and we concentrate here only on the estimator (1). For a detailed discussion of different classes of linear estimators the reader is referred to Koop (1957).

^{*}Professor Leslie Kish suggested that I point out how Lehiri's method can be used to draw a sample by procedure 1 or 2.

theory, and compact expressions for the variance of the estimate of the population total together with variance estimates have been obtained. The following approach is used in developing the asymptotic theory: In sampling with p.p.s. and with replacement we have

$$\hat{\mathbf{Y}}^{\mathbf{l}} = \sum_{i}^{n} \mathbf{y}_{i} / n \mathbf{p}_{i}, \qquad (3)$$

where $\tilde{\Sigma}$ denotes the summation over the n units selected with replacement, as the estimator of Y. Also the variance of \tilde{Y}^1 is

$$V(\hat{\mathbf{Y}}^{1}) = \sum_{1}^{N} P_{j} \left(\frac{y_{j}}{P_{j}} - \frac{y}{n}\right)^{2}$$
(4)

where $P_j = np_j$. Therefore, by assuming that each P_j is of the order of N⁻¹ for large N, it is seen that $V(\hat{Y}^1)$ is of order N². In sampling without replacement the term of order N² will be the leading term in $V(\hat{Y})$ and hence the next lower order terms, namely terms of order N¹, will represent the gain in precision due to sampling without replacement. Therefore, Hartley and Rao evaluated (for their procedure 3) $V(\hat{Y})$ to order N¹. This is equivalent to evaluating P_{11} to order N⁻³ and substituting it in the variance formula (2). Also, for the benefit of smaller size populations, they evaluated P_{11} to order N⁻⁴ and hence $V(\hat{Y})$ to order N⁻⁴.

The purpose of the present paper is to present compact expressions for the variance and estimated variance together with simplified formulas for revised sizes in terms of the original sizes for procedures 1 and 2, obtained by using the asymptotic theory. The mathematical derivations are not given here and will be published elsewhere. It is shown that the three procedures have exactly the same value of P_{11} ; to order N⁻³ and hence identical V(\hat{Y}) to order N¹. Since the terms of order N^1 are the important terms in $V(\hat{Y})$, which contribute to the gain in precision of sampling without replacement over sampling with replacement for large N, it follows that the three procedures have practically the same efficiency. However, with procedure 3 there is no need to compute the revised sizes, and this procedure, therefore, circumvenes an operation which may be cumbersome as N becomes large. It is also shown for the case n = 2 that, to order N⁰, the estimator from procedure 1 has smaller variance than that from procedure 2 and that from procedure 2 has smaller variance than that from procedure 3. However, as N increases, the contribution from terms of order N^O becomes negligible.

II. The case
$$n = 2$$

Let $p_j^* = x_j^*/\Sigma x_j$. Then, for procedure 1,
we find

$$p_{j}^{*} = p_{j} \left[1 + \frac{1}{4} (2p_{j} - 2\sum_{1}^{N} p_{t}^{2} + 4p_{j}^{2} - 3p_{j}\sum_{1}^{N} p_{t}^{2} + 3(\sum_{1}^{N} p_{t}^{2})^{2} - 4\sum_{1}^{N} p_{t}^{3} \right]$$
(5)

to order N^{-3} . If the sample is selected by procedure 1 using the revised sizes from (5), we obtain

and

$$\mathbf{v}^{(1)}(\hat{\mathbf{Y}}) = \left[1 - (\mathbf{P}_{i} + \mathbf{P}_{i}) + \frac{1}{2}\sum_{1}^{N} \mathbf{P}_{t}^{2} - \frac{1}{2}(\mathbf{P}_{i}^{2} + \mathbf{P}_{i}^{2}) + \frac{5}{8}\mathbf{P}_{i}\mathbf{P}_{i}, -\frac{3}{32}(\sum_{1}^{N} \mathbf{P}_{t}^{2})^{2} - \frac{1}{16}(\mathbf{P}_{i} + \mathbf{P}_{i}) \sum_{1}^{N} \mathbf{P}_{t}^{2} + \frac{1}{2}\sum_{1}^{N} \mathbf{P}_{t}^{3} - \frac{3}{12}(\sum_{1}^{N} \mathbf{P}_{t}^{2})^{2} - \frac{1}{16}(\mathbf{P}_{i} - \frac{\mathbf{P}_{i}}{\mathbf{P}_{i}}) \sum_{1}^{N} \mathbf{P}_{t}^{2} + \frac{1}{2}\sum_{1}^{N} \mathbf{P}_{t}^{3} - \frac{3}{12}(\sum_{1}^{N} \mathbf{P}_{t}^{2})^{2} - \frac{1}{16}(\mathbf{P}_{i} - \frac{\mathbf{P}_{i}}{\mathbf{P}_{i}})^{2} - \frac{1}{\mathbf{P}_{i}} \mathbf{P}_{t}^{2} - \frac{1}{2}\sum_{1}^{N} \mathbf{P}_{t}^{3} \right]$$

$$(7)$$

to order N^0 , where $P_j = 2p_j$. Here $V^{(1)}(\hat{Y})$ and $v(1)(\hat{Y})$ denote the variance and the estimated variance respectively for procedure 1, and i and i' are the two units included in the sample. For procedure 2 we have

$$p_{j}^{*} = p_{j} \left[1 + (p_{j} - \sum_{1}^{N} p_{t}^{2} + 2p_{j}^{2} - 2p_{j} \sum_{1}^{N} p_{t}^{2} + 2(\sum_{1}^{N} p_{t}^{2})^{2} - 2\sum_{1}^{N} p_{t}^{3}) \right]$$
(8)

to order N^{-3} . If the sample is selected by procedure 2 using the revised sizes from (8), we find

$$\mathbf{v}^{(2)} (\hat{\mathbf{Y}}) = \sum_{1}^{N} \mathbf{P}_{j} (\mathbf{1} - \frac{\mathbf{P}_{j}}{2}) (\frac{\mathbf{y}_{j}}{\mathbf{P}_{j}} - \frac{\mathbf{Y}}{2})^{2} - \frac{1}{2} \sum_{1}^{N} (\mathbf{P}_{j}^{3} - \frac{\mathbf{P}_{j}^{2} \sum_{1}^{N} \mathbf{P}_{j}^{2}}{2}) (\frac{\mathbf{y}_{j}}{\mathbf{P}_{j}} - \frac{\mathbf{Y}}{2})^{2} + \frac{1}{8} (\sum_{1}^{N} \mathbf{P}_{t} \mathbf{y}_{t} - \frac{\mathbf{Y} \sum_{1}^{N} \mathbf{P}_{t}^{2}}{2})$$
(9)

and

$$v^{(2)}(\hat{\mathbf{Y}}) = \left[1 - (P_{1} + P_{1}) + \frac{1}{2} \prod_{i=1}^{N} P_{t}^{2} - \frac{1}{2} (P_{1}^{2} + P_{1}^{2}) + \frac{1}{2} P_{1}P_{1}, -\frac{1}{8} (\prod_{i=1}^{N} P_{t}^{2})^{2} + \frac{1}{2} \prod_{i=1}^{N} P_{t}^{3}\right] (\frac{y_{1}}{P_{1}} - \frac{y_{1}}{P_{1}})^{2}$$
(10)

to order N⁰, where V⁽²⁾ (\hat{Y}) and v⁽²⁾ (\hat{Y}) denote the variance and the estimated variance respectively for procedure 2. Finally, for procedure 3, denoting the variance and the estimated variance by V(3) (\hat{Y}) and v(3) (\hat{Y}) respectively, we find

$$v^{(3)} (\hat{Y}) = \frac{N}{1} P_{j} (1 - \frac{P_{j}}{2}) (\frac{y_{j}}{P} - \frac{Y}{2})^{2} - \frac{1}{2} \frac{N}{1} (P_{j}^{3} - \frac{P_{j}^{2} \frac{N}{2} P_{t}^{2}}{2}) (\frac{y_{j}}{P_{j}} - \frac{Y}{2})^{2} + \frac{1}{4} (\frac{N}{1} P_{t} y_{t} - \frac{Y \frac{N}{2} P_{t}^{2}}{2})^{2}$$
(11)

and

$$v^{(3)}(\hat{Y}) = \left[1 - (P_{1} + P_{1'}) + \frac{1}{2}\sum_{i=1}^{N} P_{t}^{2} - \frac{1}{2}(P_{1}^{2} + P_{1'}^{2}) - \frac{1}{4}(\sum_{i=1}^{N} P_{t}^{2})^{2} + \frac{1}{4}(P_{1} + P_{1'})\sum_{i=1}^{N} P_{t}^{2} + \frac{1}{2}\sum_{i=1}^{N} P_{t}^{3}\right] - (\frac{y_{1}}{P_{1}} - \frac{y_{1'}}{P_{1'}})^{2}$$
(12)

to order \mathbb{N}^{0} . From (6), (9) and (11) it follows that V(1) (\hat{Y}) < V(2) (\hat{Y}) < V(3) (\hat{Y}), to order \mathbb{N}^{0} . However, the three procedures have exactly the same variance to order N1, namely

$$V(\hat{Y}) = \sum_{1}^{N} P_{j}(1 - \frac{P_{j}}{2})(\frac{y_{j}}{P_{j}} - \frac{Y}{2})^{2}.$$
 (13)

Also, the three procedures have exactly the same estimated variance to order N¹, namely

$$w(\hat{\mathbf{Y}}) = \left[1 - (\mathbf{P}_{1} + \mathbf{P}_{1}) + \frac{1}{2}\sum_{1}^{N} \mathbf{P}_{t}^{2}\right] \\ \cdot (\frac{\mathbf{y}_{1}}{\mathbf{P}_{1}} - \frac{\mathbf{y}_{1}}{\mathbf{P}_{1}})^{2}.$$
(14)

Equation (13) when compared with the variance in sampling with replacement, namely (4), clearly demonstrates the reduction in variance achieved by sampling without replacement through the 'finite population corrections' $(1 - P_1/2) < 1$.

Since the three procedures have exactly the same variance to order N^{\perp} and since the terms of order N^{\perp} are the important terms that contribute to the reduction in the variance achieved by sampling without replacement for large N. one can conclude that in most of the practical situations there is little to choose between the three procedures on the basis of efficiency alone. However, since there is no need to compute revised sizes with procedure 3, this may be preferred to other procedures. If one is using procedure 1 and is satisfied with the variance to order N^{1} , then the revised sizes x] are obtained from the simplified formula

$$\mathbf{p}_{j}^{*} = \mathbf{p}_{j} \left[\mathbf{1} + \frac{1}{2} (\mathbf{p}_{j} - \sum_{1}^{N} \mathbf{p}_{t}^{2}) \right]$$
(15)

to order N⁻². Similarly, for procedure 2, the revised sizes x are obtained from

$$p_{j}^{*} = p_{j} \left[1 + (p_{j} - \frac{N}{2} p_{t}^{2}) \right]$$
 (16)

to order N⁻².

An example for the evaluation of the variance formulas

Horvitz and Thompson (1952, Table 2) give the data of a population of N = 20 blocks in Ames, Iowa. Here y_1 and x_1 denote the number of house-holds and the 'eye-estimated' number of households respectively in the ith block. Using their data the following values are obtained:

$$V(\hat{Y}^{L}) = 3,241 \text{ using } (4)$$

$$V(\hat{Y}) = 3,025$$
 to order N^{\perp} using (13).

Incidentally, for this example, $\Sigma P_{ty_{t}} = 49.63$ and $Y \Sigma P_{t}^{2/2} = 49.95$. Therefore, we find that, to the nearest integer and to order N⁰, V(1) (\hat{Y})= V(2) (\hat{Y}) = V(3) (\hat{Y}) = 3,007. It may be of inter-est to exhibit the nature of convergence of the various approximations to the variance by regarding the variance formula for sampling with replacement as an approximation to order N^2 as set out in Table 1 below.

Table 1. Terms in the approximations to the variance formula

Order of approximation	Sampling procedure	Variance	Difference
0(N ²)	with replacement	3,241	216
0(N ¹)	procedures 1, 2 and 3	3,025	18
0(N ⁰)	procedures 1, 2 and 3	3,007	

The convergence in this example appears to be satisfactory although the population size (N = 20)is much smaller than those usually encountered in survey work. The variance in equal probability sampling without replacement for this example is 16,219. Therefore, all these procedures of p.p.s. sampling are vastly superior compared to equal probability sampling. It must not be forgotten, however, that there are other devices of decreasing the variance in the latter case with the help of the known x1 values (e.g., ratio estimation). The gain in precision through unequal probability sampling without replacement as compared to sampling with replacement is about 7% (235/3241).

III. The general case $n \ge 2$ Most of the published literature on unequal probability sampling deals only with the case n=2 and does not have anything to offer for n > 2 due to difficulties in evaluating Piir, and hence $V(\hat{Y})$ and $v(\hat{Y})$. Though the case $\hat{n} = 2$ is important, particularly in stratified designs, situations often arise when n is greater than 2. A striking feature of our asymptotic approach is that it permits an easy evaluation of $V(\tilde{Y})$ and $v(\ddot{Y})$ for n > 2. It can be shown that the three procedures have exactly the same variance and estimated variance to order N1, namely

234

$$V(\hat{Y}) = \sum_{1}^{N} P_{j} \left[1 - \frac{(n-1)}{n} P_{j} \right] \left(\frac{y_{j}}{P_{j}} - \frac{y_{n}}{n} \right)^{2} (17)$$

and

$$r(\hat{\mathbf{Y}}) = (n-1)^{-1} \sum_{\substack{i>1\\1>1}}^{n} \left[1 - (\mathbf{P_i} + \mathbf{P_{i'}}) + \frac{1}{n} \sum_{\substack{1\\1}}^{N} \mathbf{P_t^2}\right] (\frac{\mathbf{y_i}}{\mathbf{P_i}} - \frac{\mathbf{y_{i'}}}{\mathbf{P_{i'}}})^2, \qquad (18)$$

where $P_{j} = np_{j}$. The assumption that P_{j} is of cr-der N⁻¹ implies that n is relatively small com-pared to N. For procedure 1 the revised sizes x_{j}^{*} are obtained from

$$p_{j}^{*} = p_{j} \left[1 + \frac{(n-1)}{2} p_{j} - \frac{(n-1)}{2} \sum_{l}^{N} p_{t}^{2} \right],$$
(19)

and for procedure 2 the x_j^* are found from

$$p_{j}^{*} = p_{j} \left[1 + (n - 1)_{p_{j}} - (n - 1) \sum_{l}^{N} p_{l}^{2} \right]$$
 (20)

to order N^{-2} .

The extension of the theory, presented here

for simple sampling, to multi-stage sampling, etc., is fairly straightforward and will not be discussed here.

References

- Durbin, J. (1953). Journal of the Royal Sta-
- tistical Society, Series B, p. 253. Goodman, R. and Kish, L. (1950). Journal of the American Statistical Association, p. 350.
- Hartley, H. O. and Rao, J. N. K. (1959). Proceedings of the Social Statistics Section, American Statistical Association, 1960, p. 56.
- Hartley, H. O. and Rao, J. N. K. (1962). Annals of Mathematical Statistics (in press). Horvitz, D. G. and Thompson, D. J. (1952).
- Journal of the American Statistical
- Association, p. 663. Koop, J. C. (1957). Institute of Statistics, Mimeo Series No. 296, North Carolina. Lahiri, D. B. (1951). Bulletin of the Inter-
- national Statistical Institute, Part II, p. 133.
- Yates, F. and Grundy, P. M. (1953). Journal of the Royal Statistical Society, Series B, p. 235.