MULTIPLE FRAME SAMPLE SURVEYS

Robert S. Cochran, University of Wyoming

I. INTRODUCTION

In the theory and practice of sampling finite populations two concepts are carefully distinguished in the literature. They are

(1) the population of units, and

(2) the frame for sampling the units.

Concept (1) does not involve much more than a clear definition of the units deemed to belong to the population. However, concept (2) goes beyond a mere definition of an aggregate of units.

In many situations it is not possible to designate a unique reference frame for some reason. It then becomes necessary to supplement the original frame with an additional frame or frames in order to obtain full coverage of the population, and the investigator must design a survey based upon a multiplicity of sampling frames. In other situations it is possible to designate one frame that will give complete coverage but it may be possible to use another frame to cover a subset of the original frame. Here, again, it may be advantageous to the investigator to consider his problem as being one of multiple frame sampling.

Historically, most cases of multiple frames have been concerned with a "master" frame with 100 coverage and a "cheap" frame not possessing 100 per cent coverage. Usually a sample design had been designated for both frames, but whenever a unit sampled from the "master" frame was encountered which was also in the "cheap" frame it was discarded.

At the 1962 meetings of the American Statistical Association, H. O.Hartley presented a paper entitled "Multiple Frame Surveys." (1962) In his paper he suggested a weighting system whereby it would not be necessary to discard any information obtained in order to arrive at an estimate of the population total or the population mean. In this paper the basis for Hartley's work will be reviewed and a comparison will be presented between using the multiple frame estimator and the screening estimator outlined above.

II. MULTIPLE FRAME TERMINOLOGY AND NOTATION IN THE TWO FRAME CASE

Consider the general situation of two interlocking frames, A and B, where their union gives 100 per cent coverage to the population of interest. In this situation there are three subsets consisting of those units that are only in A, (domain a), those units that are only in B (domain b), and those units in both frames (domain ab).

Also assume a simple random sample is to be drawn from each frame and that the number of population units in the three domains are known. The notation introduced by Hartley will be used for this discussion. It is presented in Table I. III. ESTIMATION OF THE POPULATION TOTAL

Hartley's method of weight coefficients actually creates non-overlapping strata out of the overlapping frames. Letting y be the value of the Y characteristic: attached to the in sampling unit define

$$u_{Ai} = y_i$$
, if the ith unit is in a
= py_i , if the ith unit is in ab
when sampling from the A frame. Also define

 $u_{p_i} = qy_i$, if the ith unit is in <u>ab</u>

Bi
$$= y_i$$
, if the ith unit is in b

when sampling from the B frame. Thus, the number of units in the intersection, ab, have been artificially doubled, creating two mutually exclusive strata with U characteristic attached to the units of each strata. Clearly the population total of the original Y characteristic is equal to the population total of the newly constructed U characteristic whenever p + q = 1.

The estimate of the population total of the Y characteristic is

$$\hat{Y} = N_a \bar{y}_a + N_{ab} p \bar{y}_{ab} + N_{ab} q \bar{y}_{ab} + N_b \bar{y}_b \quad (1)$$

This estimator is in the form of **a** post-stratified sampling estimator for the U characteristic. Because of this the variance of \hat{Y} when the usual finite population correction can be ignored and the sample is sufficiently large from each frame is approximately

$$V(\hat{Y}) = \frac{N_{A}}{n_{A}} \left\{ \sigma_{a}^{e}N_{a} + \sigma_{ab}^{e}N_{ab}p^{2} \right\} + \frac{N_{B}}{n_{B}} \left\{ \sigma_{b}^{e}N_{b} + \sigma_{ab}^{e}N_{ab} q^{2} \right\}$$
(2)

where σ_a^2 , σ_b^2 , and σ_{ab}^2 are the wihin post stratum variances.

After some algebra this becomes

$$V(\hat{Y}) = \frac{N_A^2}{n_A} \left\{ \sigma_a^2(1-\alpha) + \alpha p^2 \sigma_{ab}^2 \right\}$$
$$\frac{N_B^2}{n_B} \left\{ \sigma_b^2(1-\beta) + \beta q^2 \sigma_{ab}^2 \right\}$$

where $\alpha = N_{ab}/N_A$, and $\beta = N_{ab}/N_B$. (3)

Assuming a simple cost function

$$C = c_A n_A + c_B n_B$$
(4)

TABLE I HARTLEY'S NOTATION FOR TWO FRAME DESIGNS AND ESTIMATES

Category	Fram	.e		Domain		
	A	В	а	Ъ	ab	
Population number	NA	NB	Na	Nb	Nab	
Sample number*	ⁿ A	n B	na	n b	¹¹ ab'	n " ab
Population total	YA	Y _B	Ya	Чb	Y _{ab}	
Population mean	Ϋ́ _A	Ϋ́ _B	Ŷ	Ϋ́ _b	Ϋ́ _{ab}	
Sample total*	УА	У _В	y a	У _Ъ	y '	y"
Sample mean*	y _A	y _B	y _a	y _b	y _{ab}	y _{ab}
Cost of sampling unit*	с _А	с _в				

*Applies to case of drawing random samples from both frames

where C is the total cost of sampling, $\mathbf{c}_{\scriptscriptstyle \Delta}$ is the cost of obtaining an observation from the A frame and c_B is the cost of obtaining an observation

from the B frame, it is possible to contemplate finding the values of n_A , n_B , and p that will

give a minimum value for the variance whenever the cost is fixed or vice versa.

After some labor, the optimum value of p is

found to be one of the solutions of $p^2 \rho \quad [\phi_B(1-\beta) + \beta \quad q^2] = q^2 \quad [\phi_A(1-\alpha) + \alpha p^2]$

where $\rho = c_A/c_B$, $\phi_B = \sigma_b^2/\sigma_{ab}^2$, and $\phi_A = \sigma_a^2/\sigma_{ab}^2$. (5) Once the value of p has been determined from the above the values of n_A and n_B can be found from

$$\frac{{}^{n}\underline{A}}{{}^{N}\underline{A}} = \Theta \left\{ \left(\sigma_{a}^{2} (1-\alpha) + \alpha p^{2} \sigma_{ab}^{2} \right) f c_{\underline{A}} \right\}^{\frac{1}{2}}$$

$$\frac{{}^{n}\underline{B}}{{}^{N}\underline{B}} = \Theta \left\{ \frac{\left(\sigma_{b}^{2} (1-\beta) + \beta q^{2} \sigma_{ab}^{2} \right)}{c_{\underline{B}}} \right\}^{\frac{1}{2}}$$

$$(6)$$

where θ would be determined by the budget available.

IV. ONE FRAME IS 100 PER CENT

An important special case of the above exists whenever one of the frames, say A, gives 100 per cent coverage to the population of interest and the other frame covers only a subset of the population. An example of such a situation is a survey of households in a city which would be sampled by the block sampling plan while it may be possible to sample most of the population by using the telephone directory. However, there are households that do not have a telephone and can only be sampled by more expensive personal interviews.

When using simple random sampling from both frames we have a special case $\beta = 1$ and $\sigma_B^2 = \sigma_{ab}^2$. With this simplification the estimator becomes

 $Y = N_a \overline{y}_a + p N_{ab} y_{ab}' + q N_{ab} y_{ab}''$ (7)

and the variance becomes

Using the same cost equation as above the optimum p now becomes

$$p = \sqrt{\frac{\mathcal{I}_{A}(1-\alpha)}{-\alpha}}$$
(9)

An alternate sampling plan for such a situation is the one mentioned in the introduction as the plan that has been usually used when such situations have developed historically. This alternate is actually a special case of the above with p = 0 and q = 1. Therefore, the estimator is

$$\hat{Y}_{0} = \frac{N_{a}}{n_{a}} y_{a} + \frac{N_{B}}{n_{B}} y_{B}$$
(10)

and the approximate variance is $V(\dot{Y}_{0}) = \frac{N_{A}^{2}}{A} (1-\alpha) \sigma_{0}^{2} + \frac{N_{B}^{2}}{B} \sigma_{0}^{2}$

$$\begin{aligned} \chi_{0} &= \frac{\Lambda_{A}}{n_{A}} \begin{pmatrix} 1-\alpha \end{pmatrix} \sigma_{a}^{2} + \frac{\Lambda_{B}}{n_{B}} \sigma_{ab}^{2} \\ &= \frac{\Lambda_{A}^{2}}{n_{A}} \sigma_{ab}^{2} \left[(1-\alpha) \phi_{A} + \frac{n_{A}}{n_{B}} \alpha^{2} \right] \end{aligned} (11)$$

For this procedure the cost equation would be

$$C = n_a c_a + n_{ab}' c_A' + n_B c_B$$
(12)

because n and n ' are random we have

$$E(C) = (1 - \alpha) c_A + \alpha c_A n_A + c_B n_B$$
$$= c_A^* n_A + c_B n_B$$

and

- $c_A =$ the cost of interviewing sampled units from the 100 per cent frame
- $c_{B}^{}$ = the cost of interviewing sampled units from the list
- c_{A}^{\prime} = the cost of determining that units sampled from the 100 per cent frame are also on the list (screening cost)

Using this cost equation the optimum values of n, and n_R yield

$$\frac{n_{A}}{n_{B}} = \frac{1}{\alpha} \sqrt{\frac{\phi_{A}}{\rho^{*}}}$$
(13)

where $\rho^* = \frac{c_A^*}{c_B}$ Using this expression for n_A/n_B in $V(\hat{Y}_0)$ yields $V(\hat{Y}_0) = \frac{N_A^2}{n_A} \frac{\sigma_{ab}^2}{\sigma_{ab}} \begin{bmatrix} (1-\alpha) & \emptyset_A \\ & \alpha \sqrt{\frac{\emptyset_A (1-\alpha)}{\rho *}} \end{bmatrix}$ (14)

Likewise in C it yields

$$C = n_{A} c_{A}^{*} \left(\begin{array}{c} 1 & + \frac{\alpha}{\sqrt{\rho * \sqrt{\theta_{A}} (1-\alpha)}} \end{array} \right)$$
(15)

Thus the variance of the screening estimator is

$$V_{0} = \frac{N_{A}^{2} \sigma^{2}}{C} \qquad c_{A}^{*} \left(1 + \frac{\alpha}{\sqrt{\rho^{2} \sqrt{\beta}} \sqrt{\beta} (1-\alpha)}\right)^{2} \phi_{A} (1-\alpha)$$
(16)

On page 4 equation (3), the corresponding multiple frame estimator is given to be

$$V(\hat{Y}) = \frac{N_A^2}{n_A} [(1-\alpha) \sigma_a^2 + \alpha p^2 \sigma_{ab}^2] + \frac{N_B^2}{n_B} q^2 \sigma_{ab}^2.$$
(17)

The optimum values for p and q given values of $n_{\underline{A}}$ and $n_{\underline{R}}$ lead to

$$\frac{\mathbf{p}}{\mathbf{q}} = \frac{\mathbf{n}_{\mathbf{A}}}{\mathbf{n}_{\mathbf{B}}} \quad \frac{\mathbf{N}_{\mathbf{B}}}{\mathbf{N}_{\mathbf{A}}} = \frac{\mathbf{n}_{\mathbf{A}}}{\mathbf{n}_{\mathbf{B}}} \quad \alpha.$$
(18)

Substitution of the above into $V(\hat{Y})$ leads to

$$\mathbf{V}(\mathbf{\hat{Y}}) = \frac{\mathbf{N}_{\mathbf{A}}^{2}}{\mathbf{n}_{\mathbf{A}}} \quad \mathbf{d}_{\mathbf{ab}}^{2} \quad [(1-\alpha) \quad \mathbf{\mathbf{A}}_{\mathbf{A}} + \alpha \mathbf{p}]$$
(19)

Also substituting p/q above into the cost equation on page 4. equation (4), leads to

$$C = c_A^n n_A \left(1 + \frac{q}{\rho} \frac{\alpha}{p} \right)^2$$
(20)

Therefore. the variance of the multiple frame estimator is

$$V(\hat{Y}) = \frac{N_A^2}{C} \frac{\sigma_A^2}{ab} - c_A p^2 \rho \left(1 + \frac{q}{p} \frac{\alpha}{\rho}\right)^2 \quad (21)$$

When the total budget C is the same for the two types of investigations the ratio of V(Y) to V(\dot{Y}_0) will give an indication of the relative efficiency of the screening estimator as compared

After some algebra this ratio becomes

to the multiple frame estimator.

$$\frac{V(\hat{Y})}{V(\hat{Y}_{0})} = \frac{V_{P}}{V_{0}} = \frac{c_{A}}{c_{A}^{*}} \frac{\rho}{\rho - \alpha} \frac{(1 + \frac{q \alpha}{p \rho})^{2}}{(1 + \frac{\alpha}{p \sqrt{\rho^{*} \sqrt{\rho - \alpha}}})^{2}}$$

Letting

$$w = \frac{c_A}{c_A^*}$$

$$\rho^* \text{ becomes}$$

$$\rho^* = \frac{\rho}{w}$$

and the variance becomes $\frac{q \alpha_{1}}{1 + q \alpha_{2}}$

$$\frac{\mathbf{V}_{\mathbf{p}}}{\mathbf{V}_{\mathbf{0}}} = \mathbf{w} \frac{\rho}{\rho \cdot \alpha} \qquad \frac{(1 + \frac{p}{p} \rho)}{(1 + \frac{\alpha \sqrt{w}}{p \sqrt{\rho} \sqrt{\rho - \alpha}})^2} \qquad (24)$$

The screening estimator will have the lower variance whenever this ratio is greater than 1. In order to determine parametric conditions that will result in such a situation let

$$A = \frac{\rho}{\rho - \alpha}$$

$$B = (1 + \frac{q \alpha}{p \rho})^{2}$$

$$C = \frac{-1}{p \sqrt{\rho \sqrt{\rho - \alpha}}}$$
(25)

With these definitions $V/V_0>1$ becomes

$$\frac{\mathbf{w} \mathbf{AB}}{(1 \div \sqrt{\mathbf{w}})^2} > 1$$
(26)

From this it can be shown that $\sqrt{w} > (\sqrt{AB} - C^{-1})^{-1}$

and because ρ is usually greater than αit can be shown that

$$\sqrt{AB} > c^{-1}$$
. (28)

(27)

-2

Therefore,

$$w > (\sqrt{AB - c^{-2}})$$
 (29)

which yields

$$w \neq \left[\left(1 + \frac{q \alpha}{p}\right) \sqrt{\frac{\rho}{\rho - \alpha}} - \frac{\alpha}{p \sqrt{\rho} \sqrt{\rho - \alpha}} \right] = \frac{\rho}{\rho - \alpha} \neq 1$$
(30)

Thus w will be greater than ρ whenever $\rho - \alpha$ V/V₀ is greater than 1. It can also be shown that w greater than ρ implies that V/V₀ will $\rho - \alpha$ (31) be greater than 1.

In terms of the cost conditions in its definition

$$w = \frac{c_A}{c_A^*} = \frac{c_A}{(1-\alpha)c_A + c_A^* \alpha} = \frac{1}{(1-\alpha) + w^* \alpha}$$

where

$$w^{\dagger} = \frac{c_{A}^{\dagger}}{c_{A}} .$$

Since the advantage is to the screening estimator whenever $w_{7-\alpha}$ using the definition of w $\rho - \alpha$ in (31) we find

$$\frac{1}{(1-\alpha) + w' \alpha} > \frac{\rho}{\rho - \alpha}$$
(32)

or

 $\frac{\rho^{-\alpha}}{\rho} > 1 - \alpha + \alpha w'$

and

$$1 - \frac{1}{\rho} > w'.$$
 (33)

Using the definitions of ρ and w' this can be written as

$$1 - \frac{c_B}{c_A} - \frac{c_A'}{c_A} \quad \text{or} \quad c_B < c_A - c_A' \quad (34)$$

This indicates that on the average the screening estimator will have the lower variance whenever the cost of sampling from the supplementary frame is less than the difference between sampling from the 100 per cent frame and screening members of the 100 per cent frame in the supplementary frame. For example, if it does cost less to ask informative questions of a person on the telephone than it does to ask them face to face, the screening estimator will have the lower variance.

Some illustrations of the relationship between the variance ratio V/V₀ and the screening cost ratio w' under various parametric conditions are given in the following tables and graphs. These are set up for a low (.25) and a high (.75) values of $\phi_A = o_a^2/o_{ab}^2$, for a low (.20), a medium (.50), and a high (.90) values of $\alpha = N_{ab}/N_A = N_B/N_A$, and for a low (1), a medium (2), and high (10) values of $\rho = c_A/c_B$. The range presented for w' = c_A^2/c_A is from 0.0 to 1.0 in increments of 1.

TABLE II $\phi_{A} = .75$

TABLE III $\phi = .25$

w' p	<u>1</u> •5	2 .33	20 10 .14	$\frac{\alpha}{1} = \frac{1}{2}$	<u>50</u> 10 .11	$\frac{\alpha}{1} = .9$ $\frac{1}{.50} \cdot .15$	<u>10</u> .05
.0	1.00	1.09	1.19	1.00 1.20	1.52	$\begin{array}{c} 1.00 & 1.20 \\ .96 & 1.14 \\ .93 & 1.09 \\ .91 & 1.06 \\ .89 & 1.02 \end{array}$	1.67
.1	.99	1.07	1.17	.96 1.14	1.44		1.50
.2	.97	1.05	1.14	.94 1.11	1.36		1.38
.3	.96	1.03	1.12	.91 1.07	1.29		1.30
.4	.94	1.02	1.10	.88 1.03	1.23		1.23
.5	•93	1.00	1.07	.86 1.00	1.18	.87 1.00	$1.17 \\ 1.12 \\ 1.08 \\ 1.04 \\ 1.00$
.6	•92	.99	1.06	.83 .96	1.13	.85 .97	
.7	•90	.96	1.04	.81 .94	1.08	.84 .95	
.8	•89	.95	1.02	.81 .92	1.05	.83 .94	
.9	•88	.94	1.00	.79 .91	1.00	.81 .91	
1.0	.87	.92	.98	.77 .86	•97	.80 .90	•97