It may be occasionally necessary or economically convenient to utilize more than one frame in drawing a sample from a population. This paper assumes two frames cover the population and the observational units for the two frames are identical. If the sizes of both frames and the population size are known, the number of elements common to the two frames is also known, as well as the number of elements included only on each individual frame.

Estimators will be developed that have smaller variance than those previously suggested. However, since the number of duplicated elements included in the overall sample is utilized in the proposed estimators, the cost is greater than for previous estimators. Consequently, allocation procedures and situations likely to result in appreciably higher cost will be considered.

The adaptation of a previously suggested method to construct unbiased estimators of the population total with empty domain sample sizes is illustrated.

1.1 Examples

The opportunity to utilize two or more frames of known sizes does occur in practice. Comstock, et. al. (3) describe a study designed to evaluate immunization histories obtained from a sample of the population of Washington County, Maryland. The evaluation was completed by the comparison of historical information from interviews with the results of serologic determinations.

The study was done during the summer of 1968. A 1% systematic sample of the county population was drawn from a list of households obtained in a non-official census conducted in 1963, supplemented by a similar sample of dwelling units added since that time. It seems reasonable to assume that the number of households on both frames were known as well as the total number of households in the county. From this knowledge it is easily determined if the two frames cover the population and if they overlap. If there is an overlap, consideration should be given to the determination of the duplicated households. Since it is not obvious that the costs of doing so can be justified, the improvements in the resulting estimators need to be carefully evaluated.

Other examples of studies that have been conducted based on multiple and possibly overlapping frames include those described by Serfling, Cornell and Sherman (8), Bershad (1) and Cochran (2).

1.2 Notation

It is assumed two frames, A and B, containing N_A and N_B elements respectively are a-vailable. The notation of Hartly (5) is adopted and N denotes the number of elements in-cluded on both frame A and Frame B. N is the number of elements occurring only on frame A and N_b is the number of elements occurring on Frame B. Thus: $N_A = N_0 + N_{ob}$, (1 1)

$$N_{B} = N_{b} + N_{ab}$$
(1.1)
(1.2)

and the total number of elements in the popula-

ion, N, is given by

$$N = N_a + N_b + N_a = N_a + N_B = N_b + N_A.(1.3)$$

The elements contained only on frame A are called domain a, the elements only on frame B domain b and those elements on both frames A and B domain ab. It is assumed that a simple random sample of size n_A is selected from Fr. A & and simple random sample of size n_B is selected from FrB. The number of elements sampled from frame A and contained in domain a is denoted by n_. The number of elements sampled from frame A and contained in domain ab is denoted by n'. The number of sampled elements in domains ab and b drawn from frame B are denoted by n''_{ab} and n_b respectively. Thus

and

 $n_A = n_a + n'_{ab}$

$$n_{B} = n_{ab}'' + n_{b}$$
 (1.5)

(1.4)

This completes the description of the problem. It is only one of several frame problems that the sample designer may face. Kish (6) gives several interesting and informative discussions of these additional problems.

2. ESTIMATORS OF THE POPULATION TOTAL

Assume a sample of size n_A is drawn without replacement from frame A, and a sample of size n_B is drawn without replacement from frame B. Assuming simple random sampling in both frames, the probability of being included in the sample, Π_i , can be calculated for elements in each domain. These probabilities will be utilized in the construction of alternative estimators.

$$Prob\left[\begin{array}{c} i^{\text{th}} \text{ element in domain a is included in the}\\ \text{sample} \right] = \frac{n}{N}, \qquad (2.1)$$

 $Prob\left[i^{th} \text{ element in domain b is included in the sample}\right] = \frac{n_B}{N_B}, \qquad (2.2)$

Prob it and element in domain ab is included in the sample at least once

$$= \frac{n_{A}N_{B} + n_{B}N_{A} - n_{A}n_{B}}{N_{A}N_{P}}$$
(2.3)

Lund (7) proposed the estimator

$$\hat{Y}_{L} = N_{a}\bar{y}_{a} + N_{ab}\bar{y}_{ab} + N_{b}\bar{y}_{b}$$
(2.4)

for the case of known domain sizes. The sample total for domain ab,

$$y_{ab}^{\star} = \sum_{i=1}^{n'ab} y_i + \sum_{i=1}^{n''b} y_i , \qquad (2.5)$$

in \hat{Y}_{L} is based on the $n'_{ab} + n''_{ab}$ elements sampled from frames A and B. If the duplicated elements in domain ab are excluded, the sample mean y_{ab}^{\star} becomes .

$$\bar{y}_{ab} = \sum_{\underline{i=1}}^{n'ab} y_{\underline{i}} + \sum_{\underline{i=1}}^{n'ab} y_{\underline{i}} - \sum_{\underline{i=1}}^{n} y_{\underline{i}}$$
(2.6)

where \mathfrak{n}_d is the number of duplicated items.

Consideration is thus given to $\dot{Y}_d = N_a \bar{y}_a + N_a b \bar{y}_a b + N_b \bar{y}_b$. (2.7) The notation \dot{Y}_d in (2.7) indicates that this es-timator of the total is based on the distinct elements included in the sample.

If it can be assumed that n_a , n_{ab} and n_b are each greater than zero, \hat{Y}_d is unbiased. The conditional expectation becomes $E(\hat{Y}_d \mid n_a > 0, n_a > 0, n_b > 0)$

$$= N_{a}E(\bar{y}_{a}|n_{a}>0) + N_{ab}E(\bar{y}_{ab}|n_{ab}>0) + N_{b}E(\bar{y}_{b}|n_{b}>0)$$
$$= N_{a}\bar{y}_{a} + N_{ab}\bar{y}_{ab} + N_{b}\bar{y}_{b} = Y.$$
(2.8)

It should be noted that Fuller (4) has devised a method of constructing unbiased post-stratified estimators that does not require each domain size to be positive. This method will be considered below.

A comparison of the variances of Y_{L} and Y_{d} reduces to a comparison of the variances of y_{ab} and y_{ab}^{*} since the estimators differ only in the extimated mean of the overlap domain. To facilitate this comparison, let

$$\bar{y}_{ab}^{\star} = \frac{y_{ab} + n_d \bar{y}_d}{n_{ab} + n_d}$$
 (2.9)

and

$$\bar{y}_{ab} = \frac{y_{ab}}{n_{ab}}$$
(2.10)

where $n_{ab} = n'_{ab} + n''_{ab} - n_d$, y_{ab} is the total of the n distinct elements sampled from domain ab and \bar{y}_d is the mean of the n_d duplicated elements. It is assumed that n_{ab} is greater than zero. Conditional on n , and n ,,

$$Var(\bar{y}_{ab}^{\star}) = \frac{Var(y_{ab}) + n_d^2 Var(\bar{y}_d) + 2n_d Cov(y_{ab}, \bar{y}_d)}{(n_{ab} + n_d)^2} (2.11)$$

and

$$\operatorname{Var}(\bar{y}_{ab}) = \frac{\operatorname{Var}(y_{ab})}{\frac{n^2}{n^2}} \quad (2.12)$$

To evaluate (2.11) and (2.12), assume the variance-covariance structure of the y,'s to be

 $\begin{bmatrix} -1/N_{ab} & -1/N_{ab} & -1/N_{ab} \\ \text{Then the variance of } y^{*}_{ab} \text{ is greater than or equal} \end{bmatrix}$

to that of \bar{y}_{ab} if

$$n_{ab}^2$$
 Var $(y_d) + 2n_{ab}^2$ Cov (y_{ab}, y_d)

$$\geq 2n_{d}n_{ab} \quad \forall ar(y_{ab}) + n_{d}^{2} \forall ar(y_{ab}). \quad (2.14)$$

Utilization of the variance-covariance matrix, (2.13), allows us to express each term of (2.14)as follows:

$$n_{ab}^{2} Var(y_{d}) = n_{ab}^{2} n_{d} \left[\frac{1 - (n_{d} + 1)}{N_{ab}}\right] S_{ab}^{2}$$
 (2.15)

$$2n_{ab}^{2}Cov(y_{ab},y_{d}) = 2n_{ab}^{2}n_{d}\left[1-n_{ab}\right]S_{ab}^{2} \quad (2.16)$$

$$2n_{ab}n_{d}Var(y_{ab}) = 2n_{ab}^{2}n_{d} \left[\frac{1-(n_{ab}+1)}{N_{ab}}\right]S_{ab}^{2}$$
(2.17)

a

$$n_{d}^{nd} n_{d}^{2} Var(y_{ab}) = n_{ab} n_{d}^{2} \left[\frac{1 - (n_{ab} + 1)}{N_{ab}} \right] S_{ab}^{2} . \quad (2.18)$$

If n_d equals zero, equality holds in expression (2.14) since the estimators \bar{y}_{ab} and \bar{y}_{ab}^{\star} are identical. If n_{ab} and n_{d} are both greater than zero, (2.14) may be written as

$$\binom{(n_{ab}-n_{d})}{n_{ab}} - \frac{(n_{ab}+n_{d})}{N_{ab}}$$
 (2.19)

Thus \bar{y}_{ab} has smaller variance than does y_{ab}^{\star} for all positive values of n_{ab} and n_{d} .

One of the reasons for calculating the probability of selection, Π_i , for each element in the population now is evident. The Π ,'s clearly demonstrate that every element in domain ab has the same probability of being included in the sample. This is required if the simple mean of the distinct elements sampled from domain ab is to be used as an unbiased estimator of the domain mean.

3. UNBIASED ESTIMATOR OF Y WITH EMPTY DOMAINS

It was noted above that Fuller (4) has devise ed a scheme to construct unbiased post-stratified estimators. This scheme does not require the usual assumption of non-empty strata. Fuller's general construction will be reviewed and then his approach will be applied to this problem.

Assume a random sample of size n has been drawn from a population. After the sample has been taken, the sampled elements are classified as members of two strata. Assume that the population is such that the population proportion of elements contained in stratum one, P1, and the proportion contained in stratum two, $P_2 = 1 - P_1$, are known. Fuller then considers the general

estimator

$$A_{i}\bar{y}_{1} + (1-A_{i})\bar{y}_{2}$$
, (3.1)

where

and

- \bar{y}_1 = the sample mean of the characteristic y for stratum one,
- y_2 = the sample mean of the characteristic y for stratum two

 A_i = the weight applied to the mean of strat tum one for samples with i (i = 0, 1,. . . , n) sample elements in stratum one.

It is further assumed that $0 \stackrel{<}{-} A_{i} \stackrel{<}{-} 1$ for all i, $A_0 = 0$ and $A_n = 1$.

Minimization of the conditional mean square error of (3.1) yields

$$A_{i} = \frac{if_{2i}s_{2}^{2} + i(n-i)P_{1}(\bar{Y}_{1}-\bar{Y}_{2})^{2}}{(n-i)f_{1i}s_{1}^{2} + if_{2i}s_{2}^{2} + i(n-i)(\bar{Y}_{1}-\bar{Y}_{2})^{2}}, \quad (3.2)$$

 $t_{1i} = \frac{N_1 - 1}{N_1}$

where

and

$$f_{2i} = \frac{N_2 - n}{N_2}$$

The estimator (3.1) employing the weight (3.2) is in general biased. However, Fuller shows it is possible to derive weights A, such

that the estimator (3.1) is unbiased.

Fuller also extends the above development for two post-strata to the general case of more than two strata. First the strata must be arranged in a natural order. Then the strata are repeatedly divided into groups of two. Beginning at the finest subdivision, an unbiased estimator constructed for each pair of strata.

In this problem, the situation is as shown in Table 1.

Table 1 A DESCRIPTION OF THE POST-STRATA RESULT-ING FROM TWO OVERLAPPING SAMPLING FRAMES

Stratum	Pro portion of pop. in stratum	Stratum ID	Sample number	Sample mean
a	N _a /N	1 1	na	۶ _а
ab	N _{ab} /N	2 1	nab	y _{ab}
Ъ	N _b /N	22	пъ	y _b

Table 1 indicates that the first division of the strata is into stratum a and strata ab and b. Thus an unbiased estimator will be constructed first for strata ab and b. It is to be remembered that n_a , n_{ab} or n_b may be zero in the development of this unbiased estimator.

In the development of the scheme, it is convenient to let $n = n_a + n_{ab} + n_b n_1 = n_a n_2 = n_{ab}$ + n_b , $P_{21} = N_{ab}/N_B$, $P_{22} = N_b/N_B$ and $P_1 = N_a/N$. The estimator for this specific problem is $N_{\overline{y}} = N \left[A_{1} \overline{y}_{a} + A_{2} (A_{21} \overline{y}_{ab} + A_{22} \overline{y}_{b}) \right]$, (3.3)

~... ..

ΓN.

۱ ٦

where the general expressions for the weights in (3.3) are

$$A_{1} = (1-A_{2}) = \frac{n_{a} + n_{a}(n_{ab} + n_{b})}{n + n_{a}(n_{ab} + n_{b})M_{1}}$$
(3.4)

and

$$A_{21} = (1 - A_{22}) = \frac{n_{ab} + n_{ab} n_{b}}{(n_{b} + n_{ab}) + n_{ab} n_{b} M_{2}} + \frac{\Lambda_{M}}{n_{b}} 2 \int_{A_{21}} (3.5)$$

where



$$s_1^2 = s_2^2 = s_w^2 > 0$$
,

$$\lambda_{M_{1}} = \frac{F_{1} - \sum_{i=1}^{n-1} \frac{i+i(n-i)P_{1}M_{1}}{n+i(n-i)M_{1}}}{\sum_{i=1}^{n-1} \frac{P_{1}i(n-i)}{n+i(n-i)M_{1}}}$$

and

$$\lambda_{M_{2}} = \frac{F_{21} - \sum_{i=1}^{n_{2}-1} \frac{i+i(n_{2}-i)P_{21}M_{2}}{n_{2}+i(n_{2}-i)M_{2}}}{\sum_{i=1}^{n_{2}-1} \frac{P_{1}i(n_{2}-i)}{n_{2}+i(n_{2}-i)M_{2}}} \cdot (3.6)$$

 F_{21} in (3.6) is P_{21} -Prob $[n_{ab} = (n_{ab}+n_b)$ given strata ab and b contain $n_{ab} + n_b$ sample elements] . That is, F_{21} is N_{ab}/N_B minus the probability that all n_R sampled elements fall in domain ab. P, in expression (3.6) is the probability domain

ab contains i units given that n_2 units have been selected from domains ab and b. 4. COST CONSIDERATIONS

Cost considerations are now introduced into the efficiency comparisons. The employment of the proposed estimator, \hat{Y}_{i} , necessitates the identification of duplicated elements in domain ab. Lund's estimator (2.4) is a special case of Hart ley's(5) procedure which utilizes a weighted average of y'_{ab} and y''_{ab} and does not require this identification. Hartley gave expressions for the sampling fractions n_A/N_A and n_B/N_B that minimize the variance of his estimator subject to the cost restraint

$$C = c_A n_A + c_B n_B . \qquad (4.1)$$

Lund's estimator results when these optimum sampling fractions are used to solve the bi-quadratic equation given by Hartley for the value of the weight p. Another possibility is the retention of elements from domain ab from one frame only. This procedure was employed by the Bureau of the Census in a 1949 study (1). This procedure is a special case of Hartley's procedure with the weight p = 1.

A third procedure is to merge the two frames before sampling and remove the duplicated elements Once the merging has been completed, any number of sampling schemes could be utilized. For example, one could employ stratified random sampling where the strata are the three domains a, b and ab.

The variance of Y_d (2.7) can be minimized subject to the cost constraint (4.1). A system of three equations in three unknowns, $\boldsymbol{f}_{A}^{},\;\boldsymbol{f}_{B}^{}$ and $\boldsymbol{\lambda}$

results. Utilization of the ratio of two of these equations reduces the system to the following system of two equations in two unknowns, f_A

and f_B:

$$\frac{c_{A}N_{A}}{c_{B}N_{B}} = \frac{f_{B}^{2} \left[N_{a}\sigma_{a}^{2} (f_{A} + f_{B} - f_{A}f_{B})^{2} + N_{ab}\sigma_{ab}^{2} f_{A}^{2} (1 - f_{B}) \right]}{f_{A}^{2} \left[N_{b}\sigma_{b}^{2} (f_{A} + f_{B} - f_{A}f_{B})^{2} + N_{ab}\sigma_{ab}^{2} f_{B}^{2} (1 - f_{A}) \right]}{C = c_{A}f_{A}N_{A} + c_{B}f_{B}n_{B}}$$
(4.2)
(4.2)
(4.2)

Solution of these two equations requires the solution of a sixth degree equation.

Rather than solving a sixth degree equation, the following iterative procedure may be used. Let f_{\star}

 $= \frac{f_A}{f_B}$

An initial value for r is

$$r_1 = \sqrt{\frac{c_B}{c_A}}$$

- **г** -

Note that (4.2) may be expressed as

(4.4)

$$\frac{n_{A}^{2}}{n_{B}^{2}} = \frac{c_{B}N_{A} \left[N_{a}\sigma_{a}^{2} \left[f_{B}+f_{A}(1-f_{B})\right]^{2} + N_{ab}\sigma_{ab}^{2}f_{A}^{2}(1-f_{B})\right]}{c_{A}N_{B} \left[N_{b}\sigma_{b}^{2} \left[f_{B}+f_{A}(1-f_{B})\right]^{2} + N_{ab}\sigma_{ab}^{2}f_{B}^{2}(1-f_{A})\right]}$$

Set 1-f and 1-f equal to 1 and divide each term of the right hand side of (4.4) by f_B^2 to obtain

$$r_{i}^{2}+1 = \frac{C_{B}N_{B}\left[N_{a}\sigma_{a}^{2}(1+r_{i})^{2} + N_{ab}\sigma_{ab}^{2}r_{i}^{2}\right]}{C_{A}N_{A}\left[N_{b}\sigma_{b}^{2}(1+r_{i})^{2} + N_{ab}\sigma_{ab}^{2}\right]}$$
(4.5)

Thus this procedure is repeated until r_i shows an arbitrarily small change from one iteration to the next.

The c_A term in (4.1) includes all the costs involved in taking a sample of size n_A from frame A. Therefore, define c_A as

$$c_{A} = c_{sA} + c_{cA} + c_{oA} \qquad (4.6)$$

That is, c_A includes the cost of selection, the

cost of classification into the proper domain and the cost of observation. The scheme used in constructing Y_d includes the costs of selection and

classification from both frames. It may be possible to select elements to be included in the sample and remove the duplicated elements before the actual observations are made. In this situation, the costs of observation are diminished since n_d fewer observations are made; however,

the elements drawn from domain ab must now be checked for duplication. The new cost equation then becomes

=
$$(c_{sA}+c_{cA}+c_{oA})n_{A}+(c_{sB}+c_{cB}+c_{oB})n_{B}$$

- $n_{d}c_{o}'+c_{d}n_{ab}'n_{ab}''$,

where

C'

 $c'_{OA} = \max(c_{OA}, c_{OB})$ (4.7)

and \boldsymbol{c}_d is the cost of determining duplications.

In (4.7) it is assumed that each sampled element is classified into its proper domain upon selection. Thus c_{1} is multiplied by the product of n'_{ab} and n''_{ab} .

We conclude that if (4.1) is the correct cost equation, \hat{Y}_d is preferred to \hat{Y}_L . The superiority of \hat{Y}_d over \hat{Y}_L is not so clear if (4.7) is the appropriate cost equation, but it should be noted that the increase in cost when employing , may be small. It is quite likely the cost of observation, which would include travel costs and expenses of an enumerator, would be larger than the cost of checking elements for duplication in the office before the fieldwork is started. That is, if it is reasonable to assume that $n_d c'_o$ is of the same magnitude as $n'_a n''_{ab} c_d$, \hat{Y}_{d} would be superior to \hat{Y}_{L} or any of the other special forms of Hartley's original estimator. However, it may be impossible to determine the duplicated elements before the fieldwork is done. In this case, \hat{Y}_d would require more expense than YL.

REFERENCES

- Bershad, M.A., "A Sample Survey of Retail Stores," <u>Sample Survey Methods and Theory</u>, <u>Vol. I</u>, New York: John Wiley and Sons, Inc., 1953, 516-58.
- Cochran, R. S., "The Estimation of Domain Sizes when Sampling Frames are Interlocking," Proceedings of the Social Science Section of the American Statistical Association Meeting, Washington, D.C., 1967.
- Comstock, G. W., et. al., "Validity of Interview Information in Estimating Community Immunization Levels," <u>Health Services Reports</u> <u>88</u> (October, 1973), 750-7.
- Fuller, W. A., "Estimation Employing Post Strata," <u>Journal of the American Statistical</u> <u>Association</u>, <u>61</u> (December, 1966), 1172-83.
- 5. Hartley, H. O., "Multiple Frame Surveys," Proceedings of the Social Science Section of the American Statistical Association meeting, Minneapolis, Minnesota, 1962.
- Kish, L., <u>Survey Sampling</u>, New York: John Wiley and Sons, Inc., 1965.
- J. Lund, R. E., "Estimators in Multiple Frame Surveys," Proceedings of the Social Science Section of the American Statistical Association meeting, Pittsburg, Pennsylvania, 1968.
- Serfling, R.E., Cornell, R.G. and Sherman, I. L., "The CDC Quota Sampling Technique with Results of 1959 Poliomyelitis Vaccination Surveys," <u>American Journal of</u> <u>Public Health</u>, 50 (December, 1960), 1847-57.