

THEORY AND APPLICATION OF REPLICATE WEIGHTING FOR VARIANCE CALCULATIONS

Robert E. Fay, U.S. Bureau of the Census¹
CDS, U.S. Bureau of the Census, Washington, DC 20233

KEY WORDS: complex samples, half-sample replication, jackknife, sample surveys.

1. Introduction

Most of the initial advances in the theory and application of replication methods to variance estimation have focused on a small number of distinct methods, namely, the jackknife, half-sample replication, the method of random groups, and the bootstrap. A number of specific adaptations and extensions have been proposed. In the original version, each method forms the replicate samples by reusing each of the original observations an integral number of times. For example, the half-sample method constructs subsamples by using half or approximately half of the original observations to form a replicate sample, the random group method uses a relatively small portion of the original sample, and the jackknife uses all but a small portion of the original sample, so that each of these methods omits some observations and forms the replicate estimate based on the remaining observations. The original version of the bootstrap, implemented by resampling with replacement from the original sample, in effect uses each of the original sample observations zero, one, two, or some other integral number of times in composing each replicate sample.

By permitting fractional weighting of observations, however, the class of replication methods becomes considerably broader and more flexible. For example, it is possible to expand the class of methods to include those that give weight 1.5 times the original weight to some of the original sample while weighting the remaining part of the sample by .5 times the original weight. As a precedent, Kreski and Rao (1981) discussed several proposals to extend the jackknife to stratified samples, with the conclusion that the successful extensions would require fractional reweighting of observations in some applications. Several researchers have previously recognized this potential redefinition of replication to include fractional reweighting; for example, Efron (1982) describes this class as "resampling plans" and recounts some of the earlier historical developments.

"Replicate weighting" refers to representing the replication method through associating replicate weights (weights used to form the replicate samples) or replicate factors (multipliers of the original weight to obtain replicate weights) with the characteristics of the sampled cases. This representation has the effect of embedding as data in a file almost all necessary information for the calculation of sampling variances. Dipbo, Fay, and Morganstein (1984) discussed a

number of earlier applications of this idea and its benefits for the analysis of data from samples with highly complex designs. The replicate weighting approach lends itself particularly well to fractional reuse of observations in forming the replicate samples.

This paper extends earlier work (Fay 1984) on a broad class of replication methods that can be easily represented through replicate weights. The general class includes methods constructed from eigenvalues and eigenvectors of a specific matrix. Section 2 extends the earlier theory to cover the replication methodology that has been implemented in the post-censal redesign of the Census Bureau's current demographic surveys after the 1980 census.

Section 3 illustrates an application of this theory to the derivation of replicate weights for the 1985 panel of the Survey of Income and Program Participation (SIPP). The 1985 panel was the first panel of the SIPP to be based on the redesign. Several complex features of the design illustrate the flexibility of the general replication approach; in particular, some of the eigenvalues and eigenvectors required by the general theory of Section 2 had to be computed numerically rather than obtained by simple inspection, as is usually possible in practice.

Section 4 summarizes the possible representation of multiple imputation calculations through a replicate weighting approach and notes other ways in which approximate measures of the uncertainty due to missing data may be obtained with the assistance of replicate weighting.

2. Theory of Generalized Replication

A previous paper (Fay 1984) established through a constructive proof the existence of replication methods to represent a wide class of variance estimators. This section further generalizes the constructive aspects of the previous proof.

Suppose the vector $\mathbf{x}_w = \{w_i x_i\}$ represents the weighted sample observations on a characteristic x_i in a sample s , where the weight w_i for sample case i may depend on s but not x_i . The usual estimate of total, based on the weighted sum of the sample characteristics, may be represented by $\mathbf{1}'\mathbf{x}_w$, where $\mathbf{1}$ denotes a vector of 1's. Virtually all of the familiar variance estimators for the variance of $\mathbf{1}'\mathbf{x}_w$ may be expressed as

$$\text{Var}^*(1'x_w) = x_w' C(s) x_w \quad (2.1)$$

where $C(s)$ is a symmetric matrix that may depend on s but not the values of x . Formula (2.1) may be adapted to estimate the covariance of any two weighted sums, $1'x_w$ and $1'y_w$,

$$\text{Cov}^*(1'x_w, 1'y_w) = x_w' C(s) y_w.$$

By extension of the previous definition (Fay 1984), we may define a *resampling plan of order k* corresponding to $C(s)$ as a set of replicate weights $w_r = \{w_{ir}\}$, $r=1, \dots, k$, which depend on s and are possibly random variables, and a set of coefficients, $b_r(s)$, possibly also depending on s , such that:

$$E\left\{\sum_{r=1}^k b_r(s) (w_r' x - w' x)^2\right\} = x_w' C(s) x_w \quad (2.2)$$

The expectation operator on the left-hand side of (2.2) refers only to any randomization of the replicate weights; consequently, (2.2) is a statement conditional on the realized sample, s . In application, however, interest is in strategies for definitions of resampling plans that can be defined for all or almost all s .

Theorem 1 of Fay (1984) established the existence of plans satisfying (2.2) for any $C(s)$. Furthermore, if $C(s)$ is positive definite or semi-definite, plans of orders as low as $k=1$ are possible with an appropriate choice of b_1 . Of course, a single replicate-weighting of the data, even if (2.2) holds, would not give a satisfactory variance estimate, but the significance of the theoretical result was to show how replicates could be formed so that the b_r were equal.

Other general methods to construct these generalized replicate weightings are important in application. If the rank of $C(s)$ is k , let $\lambda_1, \dots, \lambda_k$ be an enumeration of the non-zero eigenvalues of $C(s)$, including multiplicities, and $v_{(1)}, \dots, v_{(k)}$ a corresponding set of orthonormal eigenvectors. One possible resampling plan of order k is given for any arbitrary set of $c_r > 0$, $r=1, \dots, k$, by

$$f_r = 1 + c_r v_{(r)} \quad (2.3)$$

with $b_r = \lambda_r c_r^{-2}$. The vector f_r represents replicate factors by which to multiply each respective weight in w , that is, $w_{ir} = f_{ir} w_i$. Each eigenvector is involved in exactly one replicate in this approach. With this choice, the expectation operator in (2.2) is no longer required, that is,

$$\sum_{r=1}^k b_r (w_r' x - w' x)^2 = x_w' C(s) x_w \quad (2.4)$$

To show (2.4), note that

$$C(s) = \sum_{r=1}^k \lambda_r v_{(r)} v_{(r)}'$$

Substituting (2.3) into the left-hand side of (2.4), the question reduces to

$$\sum_{r=1}^k \lambda_r c_r^{-2} (c_r v_{(r)}' x_w)^2 = \sum_{r=1}^k \lambda_r x_w' v_{(r)} v_{(r)}' x_w,$$

which in fact is an equality. This is a stronger result than the more general (2.2). On the other hand, because the relative practical importance of the respective eigenvalues in their contribution to the overall variance is often unclear, this method should not be used unless all k replicates are included or a random selection of replicates is carried out. In other words, a variance estimate based on an arbitrary subset of the k replicates may prove unsatisfactory.

Note that the fact that c_r in (2.3) denoted an arbitrary set of positive constants provides great flexibility in choosing the replicates to meet other desired properties, for example, that all the b_r be equal.

Another replication option resembles half-sampling replication, in the sense that each eigenvector participates in each replicate, just as half-sample replication selects one of two halves from each of the strata in forming each replicate sample. Suppose all non-zero eigenvalues are positive. Let $H = \{H_{mn}\}$ be a Hadamard matrix of order $k' \geq k$ with

$$\sum_{n=1}^{k'} H_{mn} H_m' n = 0$$

for all $m \neq m'$. Then

$$f_r = 1 + c \sum_{m=1}^k H_{mr} \lambda^{\frac{1}{2}} v(m)$$

provides k' replicate factors with $b_r = 1/(k' c^2)$. Again, this choice satisfies (2.4). Unlike (2.3), however, this approach could be expected generally to yield reasonable estimates if it became necessary to employ only a subset of the initially constructed replicates, since each eigenvector appears in each replicate.

3. Application to the 1985 Panel of SIPP

The Census Bureau's program of current demographic surveys includes the Current Population Survey, the American Housing Survey, the National Crime Survey, the National Health Interview Survey, the Consumer Expenditure Survey, the Survey of Income and Program Participation, and others. As part of the redesign of these surveys following the 1980 census, a replicate weighting approach was introduced. In most instances, the replication approach represented a relatively simple adaptation of half-sample replication, substituting replicate factors of 1.5 and .5 in place of factors 2 and 0 implicit in half-sample replication. Variances may be computed from the formula:

$$\text{Var}_r^*(Y(0)) = (4/k) \sum_{r=1}^k (Y(r) - Y(0))^2 \quad (3.1)$$

where $Y(0)$ represents the estimate of a characteristic based on the original weights and $Y(r)$ represents the estimate using the replicate weights for replicate r . Note that $Y(0)$ is used in (3.1) to denote not only weighted totals, as in the previous section, but also smooth nonlinear functions of weighted totals, such as means, ratios, and many descriptive statistics defined for the weighted sample. Formula (3.1) resembles one of the usual variance estimators from half-sample replication, with the adjustment by the factor of 4 arising from the use of .5 and 1.5 in lieu of 0 and 2.0. For most current surveys, $k=48$ replicates were created.

The sample design for the 1985 Panel of the Survey of Income and Education includes several challenging or complex features that affect design-based estimation of variance. Briefly, these are:

1. The use of a Durbin-Sanford method to select two primary sampling units (PSU's) within most of the non-self-representing strata. This scheme to select two units with probability proportional to size, without replacement, gives each possible pair of units a

known positive probability of selection, thus enabling the use of the Yates-Grundy variance estimator. The more usual forms of replication do not properly represent this variance estimator.

2. After the sample had been selected, a sample reduction was implemented. The reduction in non-self-representing areas was accomplished through a random selection of whole strata for removal. Again, joint inclusion probabilities were known, but incorporation of the corresponding Yates-Grundy variance estimator to measure the resulting between-strata variance represented another layer of complexity.
3. Ernst, Huggins, and Grill (1986) developed a new weighting of between and within components of variance to be used in conjunction with the Yates-Grundy estimator. Special adaptations of the replication approach were required to integrate the results of this research.

Each of these issues poses separate problems in developing a suitable replication approach. For clarity, each issue will be addressed separately first, followed by a discussion of the integrated solution to the overall problem. Although the creation of replicate weights for the SIPP was a complex process, the resulting weights were consistent with (3.1). A total of $k=100$ replicate factors were derived, although use of a smaller number of factors may be adequate for many applications to the SIPP.

3.1 Durbin Selection of PSU's Suppose that a scheme of sampling without replacement draws samples of fixed size $n \geq 2$, Y_1, \dots, Y_n , such that each Y_i 's unconditional probability of selection, π_i , is known, as is each unconditional joint inclusion probability, π_{ij} , with $\pi_{ij} > 0$ for all $i \neq j$. For simplicity, assume that each Y_i is measured without sampling error; Section 3.3 discusses the effect of sampling error in the measurement of the Y_i 's from further stages of sampling. Then,

$$Y^* = \sum_i Y_i / \pi_i \quad (3.2)$$

is an unbiased estimator of the population total. An unbiased estimator of its variance, $\text{Var}^*(Y^*)$, is given by the Yates-Grundy variance estimator,

$$\text{Var}^*(Y^*) = \sum_{i < j} (\pi_i \pi_j / \pi_{ij} - 1) (Y_i / \pi_i - Y_j / \pi_j)^2 \quad (3.3)$$

Durbin (1967) proposed a sampling scheme satisfying these conditions, with the advantage that the leading coefficient

$(\pi_i \pi_j / \pi_{ij} - 1)$ in (3.3) is nonnegative. Negative values for this coefficient can contribute to instability in (3.3) as a variance estimator.

In the redesign of the SIPP following the 1980 census, most of the non-self-representing areas of the country were divided into strata from which $n=2$ primary sampling units (PSU's) were selected according to Durbin's scheme. Application of (3.3) to this problem therefore takes the form:

$$\text{Var}^*(Y^*) = \sum_s (\pi_{is} \pi_{js} / \pi_{ijs} - 1) (Y_{is} / \pi_{is} - Y_{js} / \pi_{js})^2, \quad (3.4)$$

where Y_{is} and Y_{js} represent the totals for the two distinct PSU's in stratum s . Similarly, π_{is} represents the probability of selecting PSU i within stratum s , etc. (Again, issues of within-PSU variance are deferred to Section 3.3.) The realized sample for the SIPP satisfied the condition that the factor $(\pi_{is} \pi_{js} / \pi_{ijs} - 1)$ was less than 4 for each s .

When (3.4) is given the general representation (2.1), then each non-zero eigenvalue of $C(s)$ corresponds to a stratum, s , and takes the value $2(\pi_{is} \pi_{js} / \pi_{ijs} - 1)$. The corresponding eigenvector \mathbf{v} whose components are 0's outside s and $2^{-\frac{1}{2}}$, and $-2^{-\frac{1}{2}}$ within s corresponding to Y_{is} and Y_{js} . Application of (2.5) with $bk' = 4$ gives replicate factors within each stratum weighting one PSU by $1 + (\pi_{is} \pi_{js} / \pi_{ijs} - 1)^{\frac{1}{2}} / 2$ and the other by $1 - (\pi_{is} \pi_{js} / \pi_{ijs} - 1)^{\frac{1}{2}} / 2$. Further comments on the exact form of the replicate factors appear in Section 3.4.

3.2 Sample Reduction in Non-Self-Representing Areas After selection of the PSU's in non-self-representing areas, reductions in the budget forced cuts in the sample sizes. The reductions in non-self-representing areas were implemented primarily through random reductions in strata. Most of the strata were grouped into sets of three or four strata, and a single strata was selected randomly from each group for elimination. Both of the sampled PSU's within each eliminated stratum were dropped from the sample. The selection was based on the Hartley-Rao-

Cochran method (Cochran 1977: pp. 266-267). For example, in the instance of groupings of four strata, one of the six possible pairs of strata was selected at random with equal probability. From the selected pair, one stratum was retained with probability proportional to size and the other dropped. If P_s represents the measure of size for stratum s , then the probability of retention in sample, π_s , is given by:

$$\pi_s = 1/2 + 1/6 P_s \sum_{t \neq s} (P_s + P_t)^{-1}. \quad (3.5)$$

Although the original Hartley-Rao-Cochran method employed a conditional weighting of the sampled data given the pairing, weights for the SIPP were based on the unconditional probability given by (3.5). Hence, the Yates-Grundy variance estimator, (3.3), is again appropriate for this problem. In this instance, however, the summation in (3.3) was over three pairs of retained strata for each grouping of four from which one stratum had been dropped. Each such grouping yields two algebraic degrees of freedom in (3.3), since (3.3) has a zero eigenvalue corresponding to the vector with three identically weighted stratum totals, Y_s / π_s .

The joint inclusion probabilities required in (3.3) are given by

$$\pi_{st} = 1/6 \{1 + P_s((P_s + P_u)^{-1} + (P_s + P_v)^{-1}) + P_t((P_t + P_u)^{-1} + (P_t + P_v)^{-1})\} \quad (3.6)$$

where t , u , and v denote the remaining strata in the grouping. Letting $g_{st} = (\pi_s \pi_t / \pi_{st} - 1)$, the matrix $C(s)$ corresponding to retained strata s , t , and v , takes the form:

$$C(s) = \begin{pmatrix} g_{st} + g_{su} & -g_{st} & -g_{su} \\ -g_{st} & g_{st} + g_{tu} & -g_{tu} \\ -g_{su} & -g_{tu} & g_{su} + g_{tu} \end{pmatrix}$$

In the instance of four strata of the same size, $g_{st} = 1/8$ for each s and t . In turn, $C(s)$ would consist of a matrix with $1/4$ on the main diagonal and $-1/8$ elsewhere. One eigenvector $(1, 1, 1)'$, has eigenvalue 0, but the eigenvalue $3/8$ has multiplicity two. Possible corresponding

eigenvectors include $(1, -1, 0)'$ and $(1, 1, -2)'$.

The actual eigenvalues and eigenvectors of $C(s)$ for each reduction of four to three strata were computed with a public-domain version of the EISPACK routines. The problem of reduction of three strata to two was handled in a similar manner, but the single eigenvector of interest was a scalar multiple of $(1, -1)'$, so that numerical methods were not required.

3.3 Incorporation of Estimates of Within Variance in the Yates-Grundy Variance Estimator The estimator (3.3) is appropriate when primary units are measured without sampling error. In the application to the SIPP, however, the Y_i/π_i represent weighted totals for the primary sampling units based on further stages of sampling within the PSU. To extend (3.4)

$$\begin{aligned} \text{Var}^*(Y^*) = & \\ & \sum_s \{ (\pi_{is}\pi_{js}/\pi_{ijs}-1)(Y_{is}/\pi_{is}-Y_{js}/\pi_{js})^2 \\ & + f_{is}\text{Var}^*(Y_{is}/\pi_{is}) + f_{js}\text{Var}^*(Y_{js}/\pi_{js}) \} \end{aligned} \quad (3.7)$$

where $\text{Var}^*(Y_{is}/\pi_{is})$ represents an unbiased estimate of the within-PSU variance of the weighted estimate Y_{is}/π_{is} . A possible choice for f_{is} , $1 - (\pi_{is}\pi_{js}/\pi_{ijs}-1)$, gives an unbiased variance estimate, but f_{is} can take negative values in some circumstances. Ernst, Huggins, and Grill (1986), proposed a modification to (3.7) that is also unbiased but which avoids negative f_{is} . These researchers furnished the actual values of these factors to be used in this application to SIPP.

3.4 Replication Strategy The previous sections discuss contributions to total variance from the sampling of strata, the sampling of PSU's within strata, and further stages of sampling within PSU's, that occurred in the majority of non-self-representing areas in the SIPP. Consider the variance over one set of either four or three strata grouped for the reduction described in Section 3.3. Let π_s represent the unconditional probability of inclusion for stratum s ; π_{st} , the joint inclusion probabilities for strata s and t ; π_{is} , the inclusion probability for PSU i in stratum s , conditional on inclusion

of stratum s ; π_{ijs} , the joint inclusion probabilities of PSU's i and j in stratum s , conditional on inclusion of stratum s ; Y_{is} , the weighted estimate of a characteristic from PSU i in stratum s ; and $\text{Var}(Y_{is})$, an estimate of the within-PSU variance of Y_{is} . Then, an estimate of the variance of the weighted estimate, Y^* , over this group of strata is given by

$$\begin{aligned} \text{Var}^*(Y^*) = & \sum_{s < t} (\pi_s \pi_t / \pi_{st} - 1) (Y_{s+}^* - Y_{t+}^*)^2 + \\ & \sum_s \left\{ 1 - \left(\sum_{s \neq t} (\pi_s \pi_t / \pi_{st} - 1) \right) \right\} \\ & \left\{ (\pi_{is}\pi_{js}/\pi_{ijs}-1)(Y_{is}/\pi_{is}-Y_{js}/\pi_{js})^2 \right. \\ & \left. + f_{is}\text{Var}^*(Y_{is}/\pi_{is}) + f_{js}\text{Var}^*(Y_{js}/\pi_{js}) \right\} \end{aligned} \quad (3.8)$$

where the f_{is} appeared earlier in (3.7) and summations are only over elements in sample, and Y_{s+}^* denotes the weighted sum over stratum s . Separate eigenvectors are associated with each of the terms in (3.8) involving Y 's: either two eigenvectors corresponding to the first term of (3.8) for reductions of four strata to three or one eigenvector otherwise; a separate eigenvector within each non-self representing stratum for the between-PSU component, and eigenvectors representing the within-PSU variance, which was estimated by dividing each PSU into appropriate half-samples.

A few strata included in the SIPP sample were represented by a single PSU. Such strata were paired and an eigenvector corresponding to a collapsed-stratum estimator, with a single degree of freedom, was created. Variance estimation within self-representing PSU's was also comparatively simple, and most were divided into a pair of half-samples. Additional half-samples were created for the largest PSU's by dividing them into substrata for purposes of the variance calculation.

All together, there were several hundred possible degrees of freedom, substantially more than intended to be represented by replicates. Analogous situations frequently arise in application of replication when many more potential degrees of freedom may be present than are required to derive a reasonable estimate of variance. When the degrees of freedom present substantially exceed the number of replicates to be produced, creation of replicates through independent randomization results in only modest losses in

efficiency. Rather than forgo the benefits of partial balancing, however, a Hadamard matrix (for example, Wolter 1985) of order 100 was employed, and a different approach was used, namely to employ the notion of confounding in experimental design to associate dissimilar degrees of freedom. Sets of dissimilar degrees of freedom grouped together for assignment to each of the rows of the matrix. For example, a row of the matrix might be assigned to one degree of freedom from within-PSU variation in a self-representing PSU, to one degree of freedom for between-stratum variance in another region, and to another degree of freedom for between-PSU variance in a third region. The confounding was such that not only would variances of estimates for national estimates preserve approximately the full 100 degrees of freedom, but estimates for important subdomains, including large central cities, regions, SMA's, and other such areas would also be represented by effectively this many degrees of freedom as well.

4. Multiple Imputation

Although a common practice in survey research is to impute for missing values and to treat the imputed values as if observed, it is well known (Little and Rubin 1987) that standard design-based estimators typically either fail to include or understate the resulting uncertainty. Multiple imputation (Rubin 1987) offers a possible approach to estimate the uncertainty arising from the imputation procedures, as well as to improve the precision of the estimates in some applications.

In the multiple imputation approach, variation among different sets of estimates are employed to infer a component of variance due to uncertainty from the missing data. This component can be added to standard estimates of variance for the whole sample.

A direct approach to incorporate multiple imputations into a survey data file is to include additional space on computer records to hold the values of different imputations. This solution presents some processing difficulties, since most computer software is not prepared to deal easily with such complexities. As Rubin and others have observed, however, it is possible to identify those survey cases requiring multiple imputation and to create multiple records for each separate imputation, allocating the total survey weight among the different records. With this representation, the generalized replication methods discussed here become effective means to complete the calculation of the variances. Assuming that the design-based expression for the sampling variance takes the general form (2.1), the resulting expression for total variance again takes the form (2.1), although with a different $C(s)$. In simple applications,

the eigenvectors will divide simply into those associated with variation due to missingness and those associated with the sample design. Consequently, replication methods of this form represent a simple and direct way to implement the variance calculations in the analysis of multiple imputations.

Replication also appears to offer a likely means to extend this methodology to problems not fully dealt with by current theory, including representing the effect of the complex design on the estimation of parameters in the missing data model. The generality of replication to represent complex variance and covariance properties would appear to offer the flexibility to address these problems in a computationally feasible manner.

¹ This paper reports the general results of research undertaken by Census Bureau staff. The views expressed are attributable to the author and do not necessarily reflect those of the Census Bureau. William Bell provided helpful comments on an earlier draft.

REFERENCES

- Cochran, W.G. (1977), *Sampling Techniques*, New York: John Wiley & Sons.
- Dippo, C.S., R.E. Fay, and D.H. Morganstein (1984), "Computing Variances from Complex Samples with Replicate Weights," *Proceedings of the Section on Survey Research Methods*, American Statistical Association, Washington, DC, pp. 489-494.
- Durbin, J. (1967), "Design of Multi-Stage Surveys for the Estimation of Sampling Errors," *Applied Statistics*, 16, 152-164.
- Efron, B. (1982), *The Jackknife, the Bootstrap and Other Resampling Plans*, Philadelphia, PA: Society for Industrial and Applied Mathematics.
- Ernst, L.R., V.J. Huggins, and D.E. Grill (1986), "Two New Variance Estimation Techniques," *Proceedings of the Survey Research Methods Section*, American Statistical Association, Washington, DC, pp. 400-405.
- Fay, R.E. (1984), "Some Properties of Estimators of Variance Based on Replication Methods," *Proceedings of the Section on Survey Research Methods*, American Statistical Association, Washington, DC, pp. 495-500.
- Krewski, D. and J.N.K. Rao (1981), "Inference from Stratified Samples: Properties of the Linearization, Jackknife, and Balanced Repeated Replication Methods," *Annals of Statistics*, 9, 1010-1019.
- Little, R.J.A. and D.B. Rubin (1987), *Statistical Analysis with Missing Data*, New York: John Wiley & Sons.
- Rubin, D.B. (1987), *Multiple Imputation for Nonresponse in Surveys*, New York: John Wiley & Sons.
- Wolter, K. M. (1985), *Introduction to Variance Estimation*, New York: Springer-Verlag.