

# Conditional Properties of Post-Stratified Estimators Under Normal Theory

Robert J. Casady and Richard Valliant

U.S. Bureau of Labor Statistics, 2 Massachusetts Ave. N.E., Washington DC 20212-001

## 1. INTRODUCTION

### 1.1 Background

A major thrust in sampling theory in the last twenty years has been to devise ways of restricting the set of samples used for inference. In a purely design-based approach, as described in Hansen, Madow, and Tepping (1983), no such restrictions are imposed. Statistical properties are calculated by averaging over the set of all samples that might have been selected using a particular design. Although it is generally conceded that some type of design-based, conditional inference is desirable (Fuller 1981, Rao 1985) satisfactory theory has yet to be developed except in relatively simple cases. A design-based approach to conditioning was introduced by Robinson (1987) for the particular case of ratio estimates in sample surveys. Robinson applied large sample theory and approximate normality of certain statistics to produce a conditional, design-based theory for the ratio estimator.

In this paper, we extend that line of reasoning to the problem of post-stratification. Convincing arguments have been made in the past by Durbin (1969) and Holt and Smith (1979) that post-stratified samples should be analyzed conditional on the sample distribution of units among the post-strata. Model-based, conditional analyses of post-stratified samples are presented in Little (1991) and Valliant (1993). The alternative pursued here is design-based and uses large sample, approximate normality in a way similar to that of Robinson (1987) as a means studying conditional properties of estimators.

### 1.2 Basic Definitions and Notation

The **target population** is a well defined collection of elementary (or analytic) units. For many applications the elementary units are either persons or establishments. We assume the target population has been partitioned into **first stage sampling units** (FSUs). The collection of FSUs will be referred to as the **first stage sampling frame** (or just **sampling frame**). It is assumed that there are  $M$  FSUs in the sampling frame and they are labeled 1, 2, ...,  $M$ . We also assume that the population units can be partitioned into  $K$  "post-strata" which can be used for the purposes of estimation.

We let  $y$  represent the value of the characteristic of interest for an elementary unit. Associated with the  $i^{\text{th}}$  FSU are  $2K$  real numbers:

$y_{ik}$  = aggregate of the  $y$  values for the elementary units in the  $i^{\text{th}}$  FSU which are in the  $k^{\text{th}}$  post-stratum,

$N_{ik}$  = number of elementary units in the  $i^{\text{th}}$  FSU which are in the  $k^{\text{th}}$  post-stratum.

For each post-stratum we then define the aggregate of the  $y$ 's and total number of elementary units:

$$Y_k = \sum_{i=1}^M y_{ik} \text{ and } N_k = \sum_{i=1}^M N_{ik} .$$

In what follows we assume that the  $N_k$  are known. The population aggregate of the  $y$  values and the total population size are given by

$$Y = \sum_{k=1}^K Y_k \text{ and } N = \sum_{k=1}^K N_k .$$

In sections 1-3, we assume that the sampling frame provides "coverage" of the entire target population. In section 4, we consider the problem of a defective frame, i.e. one in which the coverage of the frame differs from that of the target population.

### 1.3 Sample Design and Basic Estimation

Suppose that the first stage sampling frame is partitioned into  $L$  strata and that a multi-stage, stratified design is used with a total sample of  $m$  FSUs. In the following, the subscript representing design strata is suppressed in order to simplify the notation. For the subsequent theory, it is unnecessary to explicitly define sampling and estimation procedures for second and higher levels of the design. However, for every sample FSU, we require estimators  $\hat{y}_{ik}$  and  $\hat{N}_{ik}$  so that  $E_{2+}[\hat{y}_{ik}] = y_{ik}$  and

$E_{2+}[\hat{N}_{ik}] = N_{ik}$  where the notation  $E_{2+}$  indicates the design-expectation over stages 2 and higher. Letting  $\pi_i$  be the probability that the  $i^{\text{th}}$  FSU is included in the sample and  $w_i = 1/\pi_i$ , it follows that the estimators

$$\hat{Y}_k = \sum_{i=1}^m w_i \hat{y}_{ik} \text{ and } \hat{N}_k = \sum_{i=1}^m w_i \hat{N}_{ik}$$

are unbiased for  $Y_k$  and  $N_k$ .

### 1.4 An Analogue to Robinson's Asymptotic Result

Following Krewski and Rao (1981), we can establish our asymptotic results as  $L \rightarrow \infty$  within in the framework of a sequence of finite populations  $\{\Pi_L\}$  with  $L$  strata in  $\Pi_L$ . It should be understood

that we implicitly assume (without formal statement) the sample design and regularity conditions as specified in Krewski and Rao and more fully developed in Rao and Wu (1985). Details of proofs are omitted.

Converting to matrix notation, we let  $\mathbf{Y} = [Y_1 \cdots Y_k]'$ ,  $\mathbf{N} = [N_1 \cdots N_k]'$ ,  $\hat{\mathbf{Y}} = [\hat{Y}_1 \cdots \hat{Y}_k]'$ ,  $\hat{\mathbf{N}} = [\hat{N}_1 \cdots \hat{N}_k]'$  and  $\mathbf{V} = \text{var} \begin{bmatrix} \hat{\mathbf{Y}} \\ \hat{\mathbf{N}} \end{bmatrix}'$  where  $\hat{\mathbf{Y}} = (1/N) \hat{\mathbf{Y}}$  and  $\hat{\mathbf{N}} = (1/N) \hat{\mathbf{N}}$ . Analogous to conditions C4 and C5 of Krewski and Rao (1981), we assume that

$$\lim_{L \rightarrow \infty} Y_k/N_k = \mu_k, \text{ for } k=1, 2, \dots, K, \quad (1)$$

$$\lim_{L \rightarrow \infty} N_k/N = \phi_k > 0 \text{ for } k=1, 2, \dots, K, \quad (2)$$

$$\lim_{L \rightarrow \infty} m\mathbf{V} = \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \text{ (positive definite)} \quad (3)$$

where  $\Sigma$  is partitioned in the obvious manner. Note that we have again suppressed the subscript representing design strata. Assumptions (1)-(3) simply require that certain key quantities stabilize in large populations. Condition (2), in particular, assures that no post-stratum is empty as the population size increases. Letting

$$\mathbf{M}_1 = \lim_{L \rightarrow \infty} \hat{\mathbf{Y}} = [\phi_1 \mu_1 \quad \phi_2 \mu_2 \quad \cdots \quad \phi_K \mu_K],$$

$$\mathbf{M}_2 = \lim_{L \rightarrow \infty} \hat{\mathbf{N}} = [\phi_1 \quad \phi_2 \quad \cdots \quad \phi_K]' \text{ and}$$

$$\mathbf{V}_c = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}, \quad \text{it follows that}$$

$$m^{1/2} \begin{bmatrix} \hat{\mathbf{Y}} - \mathbf{M}_1 - \Sigma_{12} \Sigma_{22}^{-1} (\hat{\mathbf{N}} - \mathbf{M}_2) \\ \hat{\mathbf{N}} - \mathbf{M}_2 \end{bmatrix} \text{ tends in distribution}$$

$$\text{to } N \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{V}_c & \mathbf{0} \\ \mathbf{0} & \Sigma_{22} \end{bmatrix} \right). \text{ This is analogous to the result}$$

for  $K=1$  cited by Robinson (1987). It then follows, as in Robinson, that given  $\hat{\mathbf{N}}$  (more strictly, given  $\hat{\mathbf{N}}$  in a cell of size  $\epsilon m^{-1/2}$  for small  $\epsilon$ ), the conditional distribution of  $\hat{\mathbf{Y}}$  is asymptotically

$$N \left( \mathbf{M}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\hat{\mathbf{N}} - \mathbf{M}_2), m^{-1} \mathbf{V}_c \right).$$

Note that in some sample designs  $\mathbf{1}' \hat{\mathbf{N}} = N$  (such as those in which a fixed number of elementary units are selected with equal probabilities) in which case  $\Sigma_{22}^{-1}$  does not exist; in such cases only the first  $K-1$  post-strata are considered for the purpose of conditioning.

## 2. CONDITIONAL PROPERTIES OF ESTIMATORS FOR THE POPULATION MEAN

### 2.1 Estimators for the Population Mean

The population mean is, by definition,

$\mu = \lim_{L \rightarrow \infty} (Y/N) = \lim_{L \rightarrow \infty} (\mathbf{1}' \mathbf{Y} / \mathbf{1}' \mathbf{N}) = \sum_{k=1}^K \phi_k \mu_k$  where  $\mathbf{1}'$  is a row vector of  $K$  ones. Four estimators of the population mean will be considered. The first three are standard ones found in the literature while the fourth is a new estimator:

$$(1) \text{ Horvitz-Thompson estimator: } \hat{Y}_{HT} = \mathbf{1}' \hat{\mathbf{Y}} / \mathbf{1}' \mathbf{N} = \mathbf{1}' \hat{\mathbf{Y}}$$

$$(2) \text{ ratio estimator: } \hat{Y}_R = \mathbf{1}' \hat{\mathbf{Y}} / \mathbf{1}' \hat{\mathbf{N}} = \mathbf{1}' \hat{\mathbf{Y}} / \mathbf{1}' \hat{\mathbf{N}}$$

(3) post-stratified estimator:

$$\hat{Y}_{ps} = N^{-1} \sum_{k=1}^K (N_k / \hat{N}_k) \hat{Y}_k = \mathbf{r}' \hat{\mathbf{Y}}$$

$$\text{where } \mathbf{r}' = [N_1 / \hat{N}_1, \dots, N_K / \hat{N}_K]$$

(4) linear regression estimator:

$$\hat{Y}_{LR} = \left[ \mathbf{1}' \left( \hat{\mathbf{Y}} - \Sigma_{12} \Sigma_{22}^{-1} (\hat{\mathbf{N}} - \mathbf{M}_2) \right) \right]$$

The linear regression estimator is motivated by the form of the large sample mean of the conditional

random variable  $\hat{\mathbf{Y}} | \hat{\mathbf{N}}$  listed at the end of section 1.4.

This estimator was also discussed by Rao (1992) and is very similar to the generalized regression estimator discussed by Särndal, Swensson and Wretman (1992). It should be noted that the estimators require varying degrees of knowledge about  $N_k$  and  $N$ . The linear regression estimator has the additional complication that the covariance matrices  $\Sigma_{12}$  and  $\Sigma_{22}$  are unknown and must be estimated from the sample. In implementing  $\hat{Y}_{LR}$ , the known finite population quantities  $(1/N) \mathbf{N}$  will be used in place of the limiting vector  $\mathbf{M}_2$ .

### 2.2 Conditional expectations and variances of the estimators

Using the asymptotic setup given earlier, the conditional expectations and variances of the four estimators can be computed. First, define the following three matrices:

$$\mathbf{H} = \Sigma_{12} \Sigma_{22}^{-1}, \mathbf{R} = \mathbf{H} - \mathbf{D}(\mu), \text{ and } \mathbf{P} = \mathbf{H} - \mathbf{D}(\mu_k)$$

$$\text{with } \mathbf{D}(\mu) = \text{diag}(\mu, \dots, \mu), \mathbf{D}(\mu_k) = \text{diag}(\mu_1, \dots, \mu_k).$$

Below, we state the mean and variance of the four estimators without providing any details of the calculations. When the sample of first-stage units is large, each of the estimators has essentially the same conditional variance. The Horvitz-Thompson, ratio, and post-stratified estimators are, however,

conditionally biased, whereas the linear regression estimator is not. Thus, the linear regression estimator has the smallest asymptotic mean square error among the four estimators considered here.

(1) Horvitz-Thompson estimator:

$$E\left[\hat{Y}_{HT} \mid \hat{N}\right] = \mu + \left[ \mathbf{1}' \mathbf{H} (\hat{N} - \mathbf{M}_2) \right]$$

$$\text{var}\left[\hat{Y}_{HT} \mid \hat{N}\right] = m^{-1} \left[ \mathbf{1}' \mathbf{V}_c \mathbf{1} \right] = V_{HT(c)}$$

(2) ratio estimator:

$$E\left[\hat{Y}_R \mid \hat{N}\right] = \mu + \left[ \mathbf{1}' \mathbf{R} (\hat{N} - \mathbf{M}_2) \right] + o(m^{-1})$$

$$\text{var}\left[\hat{Y}_R \mid \hat{N}\right] = (N / \hat{N})^2 V_{HT(c)}$$

$$= V_{HT(c)} + o(m^{-3/2})$$

(3) post-stratified estimator:

$$E\left[\hat{Y}_{PS} \mid \hat{N}\right] = \mu + \left[ \mathbf{1}' \mathbf{P} (\hat{N} - \mathbf{M}_2) \right] + o(m^{-1})$$

$$\text{var}\left[\hat{Y}_{PS} \mid \hat{N}\right] = m^{-1} \left[ \mathbf{r}' \mathbf{V}_c \mathbf{r} \right] = V_{HT(c)} + o(m^{-3/2})$$

(4) linear regression estimator:

$$E\left[\hat{Y}_{LR} \mid \hat{N}\right] = \mu, \quad \text{var}\left[\hat{Y}_{LR} \mid \hat{N}\right] = V_{HT(c)}$$

Note that some minor modifications of the above formulas are necessary for designs, such as simple random sampling, in which  $\mathbf{1}' \hat{N} = N$ .

The large-sample biases of the first three estimators depend on  $\hat{N} - \mathbf{M}_2$ , a measure of how well the sample estimates the population distribution among the post-strata. In some special cases each of the first three can be conditionally unbiased. The post-stratified estimator, for example, will be approximately unbiased if  $\mathbf{1}'(\mathbf{H} - \mathbf{D}(\mu_k)) = \mathbf{0}'$ . This occurs in simple random sampling and is possible, though certainly not generally true, in more complex designs. The matrix  $\mathbf{H}$  can be interpreted as the slope in a multivariate regression of  $\hat{Y}$  on  $\hat{N}$ , or of  $\bar{Y}$  on  $\bar{N}$  when the sample estimates are close to the population values. Thinking heuristically in superpopulation terms, if  $E_\xi(y_k) = \mu_k N_k$ , as in Valliant (1993), with  $E_\xi$  denoting an expectation with respect to the model, then  $E_\xi(Y_k) = \mu_k N_k$ . The slope of the regression of  $Y_k$  on  $N_k$  is then  $\mu_k$ . In the unusual case in which the  $\hat{Y}_k$ 's are independent,  $\mathbf{H} = \mathbf{D}(\mu_k)$  and the conditional design-bias of the post-stratified estimator would be

zero. If, on the other hand, the model has an intercept, i.e. if  $E_\xi(Y_k) = \alpha_k + \mu_k N_k$ , then the post-stratified estimator may have a substantial conditional design-bias.

Similar model-based thinking can be applied to the Horvitz-Thompson and ratio estimators to show that restrictive and unrealistic models are required in order for the conditional design-biases to vanish.

### 3. SIMULATION RESULTS

The theory developed in the preceding sections was tested in a set of simulation studies using three separate populations; the results for two of these populations are given below. The population size and basic sample design parameters for the two studies are listed in Table 1.

The first population consists of a subset of the persons included in the first quarter sample of the 1985 National Health Interview Survey (NHIS). The variable of interest is the number of restricted activity days in the two weeks prior to the interview. Four post-strata were formed on the basis of age and sex in order to create population sub-groups that were homogenous with respect to the variable of interest.

The second population is artificial; it was created with the intention of producing a substantial conditional bias in the post-stratified estimator of the mean. As noted in section 2.2,  $\hat{Y}_{PS}$  will be conditionally biased if the FSU post-stratum totals for the variable of interest, conditional on the number of units in each FSU/post-stratum, follow a model with a non zero intercept. With this in mind, we generated the population in such a way that

$$E_\xi(y_k | N_k) = \alpha_k + \beta N_k + \gamma N_k^2 \quad (4)$$

where  $N_k$  is the number of units in the  $k^{\text{th}}$  post-stratum for the  $i^{\text{th}}$  FSU and  $\alpha_k$ ,  $\beta$ , and  $\gamma$  are constants. Five post-strata were used with  $\alpha_k = 100k$  ( $k=1, \dots, 5$ ),  $\beta = 10$ , and  $\gamma = -.05$ . Two thousand FSUs were generated with the total number of units in the  $i^{\text{th}}$  FSU, say  $N_i$ , being a Poisson random variable with mean 10. Conditional on  $N_i$ , the numbers of units in the five post-strata (i.e.,  $N_{i1}, N_{i2}, \dots, N_{i5}$ ) for the  $i^{\text{th}}$  FSU were determined using a multinomial distribution with parameters  $N_i$  and  $p_k = .20$  for  $k = 1, 2, \dots, 5$ .

Finally, the value of the variable of interest for the  $j^{\text{th}}$  unit in the  $k^{\text{th}}$  post-stratum for the  $i^{\text{th}}$  FSU was a realization of the random variable

$$y_{ijk} = \alpha_k / N_k + \beta + \gamma N_k + \varepsilon_{i1} + \varepsilon_{2k} + \varepsilon_{3jk} N_i$$

where  $\varepsilon_{i1}$ ,  $\varepsilon_{2k}$ , and  $\varepsilon_{3jk}$  are three independent standardized chi-square (6 d.f.) random variables.

This structure implies that  $E_{\xi}(Y_{\star}|N_{\star})$  is given by (4).

A single-stage stratified design was used for the NHIS population with "households" being the FSUs. Ten design strata were used and an approximate 10% simple random sample of households was selected without replacement from each stratum. Each sample consisted of 115 households and each sample household was enumerated completely.

A two-stage stratified sample design was used for the artificial population. One hundred design strata were created with each stratum having approximately the same number of FSUs and a systematic sample of  $m=2$  FSUs was selected with probabilities proportional to size; thus, 200 FSUs were selected. The within FSU sample size was set at 15 which resulted in the complete enumeration of most sample FSUs.

A total of 5,000 samples was selected from each of the populations for the simulation study. In each sample, we computed  $\hat{Y}_{HT}$ ,  $\hat{Y}_R$ ,  $\hat{Y}_{PS}$ , and two versions of  $\hat{Y}_{LR}$ . For the first version of the regression estimator, denoted  $\hat{Y}_{LR}(emp)$  in the tables,  $H$  was estimated separately from each sample as would be required in practice. Each component of  $\Sigma_{12}$  and  $\Sigma_{22}$  was estimated using the ultimate cluster estimator of covariance, appropriate to the design. The second version, denoted  $\hat{Y}_{LR}(theo)$ , used the same value of  $H$  in each sample, which was an estimate more nearly equal to the theoretical value of the  $H$  matrix.

Table 2 lists unconditional results summarized over all 5,000 samples from each population. Empirical root mean square errors (*rmse*'s) were

calculated as  $rmse(\hat{Y}) = \left[ \sum_{s=1}^S (\hat{Y}_s - \bar{Y})^2 / S \right]^{1/2}$  with  $S =$

5,000 and  $\hat{Y}_s$  being one of the estimates of the population mean from sample  $s$ . In the artificial population, results for the Horvitz-Thompson and the ratio estimators were nearly identical so that only the former is shown. Across all samples, the bias of each of the estimators was negligible. As anticipated by the theory,  $\hat{Y}_{LR}(theo)$  was the most precise of the choices, although the largest gain compared to  $\hat{Y}_{PS}$  was only 4.7% in the artificial population. The need to estimate  $H$  destabilizes the regression estimator as shown in the results for  $\hat{Y}_{LR}(emp)$ . For the NHIS population,  $\hat{Y}_{LR}(emp)$  has a larger root *mse* than both  $\hat{Y}_{LR}(theo)$  and  $\hat{Y}_{PS}$ .

Figures 1 and 2 present conditional simulation results. The 5,000 samples were sorted by the theoretical bias factors presented in section 2.2. The sorting was done separately for each of the estimators of the population mean. In the cases of the two regression estimators, which are theoretically

unbiased in large samples, the bias factor for  $\hat{Y}_{PS}$  was used for sorting. The sorted samples were then put into 25 groups of 200 samples each and empirical biases and root *mse*'s were computed within each group. The group results were then plotted versus theoretical bias factors in the figures. The upper sets of points in each figure are the empirical root *mse*'s of the groups, while the lower sets are empirical biases. The two regression estimators are conditionally unbiased as expected. The other estimators, however, have substantial conditional biases that, in the most extreme sets of samples, are important parts of the *mse*'s. In the neighborhood of the balance point,  $\hat{N} = \bar{N}$ , all estimators perform about the same, but, because of a lack of data at the design stage, we have no control on how close to balance a particular sample may be. The safest choice for controlling conditional bias is, thus,  $\hat{Y}_{LR}(emp)$ .

#### 4. DEFECTIVE FRAMES

##### 4.1 The Basic Problem of Defective Frames

In most real world applications not all of the elementary units in the population are included in the sampling frame. In household surveys, it is not unusual for some demographic subgroups, especially minorities, to be poorly covered by the sampling frame. Bailar (1989), for example, notes that in 1985 the sample estimate from the CPS of the total number of Black males, ages 22-24, was only 73% of an independent estimate of the total population of that group.

To formalize the discussion of this type of coverage problem, suppose that  $N_{\star}$  now refers to the number of elementary units in the frame and that  $\dot{N}_{\star}$  is the actual number of population elements in the  $k^{\text{th}}$  post-stratum. In the discussion below terms with a dot on the top are population values while terms with no dot are frame values. Letting  $\dot{Y}_k$  be the aggregate of the  $y$  values over all population elements in the  $k^{\text{th}}$  post-stratum, then it follows that the true population mean is given by

$$\dot{\mu} = \lim_{L \rightarrow \infty} \frac{\sum_{k=1}^K \dot{Y}_k}{\sum_{k=1}^K \dot{N}_k} = \lim_{L \rightarrow \infty} \sum_{k=1}^K \frac{\dot{N}_k}{\dot{N}_k} \frac{\dot{Y}_k}{\dot{N}_k} = \sum_{k=1}^K \dot{\phi}_k \dot{\mu}_k$$

Obviously, all four of the estimators of the mean given in section 2 are biased (both conditionally and

unconditionally) for  $\hat{\mu}$ . The additional bias term is  $\mu - \hat{\mu}$  for all of the estimators, and being  $o(1)$ , it will dominate the other bias terms listed in section 2.2 as the number of FSUs increases. A more basic problem is that the individual frame values  $N_k$  are usually unknown so only the ratio estimator is well defined.

On the other hand, the  $\hat{N}_k$  (or least the proportions  $\hat{\phi}_k$ ) may be known from independent sources and hence be available for the purposes of estimator construction.

Before attempting to construct unbiased estimators for  $\hat{\mu}$  it should be noted that

$$\mu - \hat{\mu} = \sum_{k=1}^K (\phi_k - \hat{\phi}_k)(\mu_k - \hat{\mu}_k) + \sum_{k=1}^K (\phi_k - \hat{\phi}_k)\hat{\mu}_k + \sum_{k=1}^K \hat{\phi}_k(\mu_k - \hat{\mu}_k)$$

So, if we assume that for each post-strata the mean of the units in the frame is equal to the true population mean, (i.e.  $\mu_k = \hat{\mu}_k$  for every  $k$ ) then the bias term reduces to

$$\mu - \hat{\mu} = \sum_{k=1}^K (\phi_k - \hat{\phi}_k)\mu_k = \sum_{k=1}^K (\phi_k - \hat{\phi}_k)\hat{\mu}_k$$

This is very strong and expedient assumption; however, addressing the problem of defective frame bias without such a condition is virtually impossible.

#### 4.2 Alternative Estimators

The basic strategy is to construct an estimator for the defective frame bias,  $\mu - \hat{\mu}$ , and then subtract this estimator from the estimators studied earlier. Two cases need to be considered: (1) The frame parameters  $\{\phi_k, 1 \leq k \leq K\}$  are unknown, and (2) The frame parameters  $\{\phi_k, 1 \leq k \leq K\}$  are known.

**Case 1.** For this case only the ratio estimator is well defined and the only obvious candidate for an estimator of the bias is

$$\hat{B}_1 = \sum_{k=1}^K \left( \frac{\hat{N}_k}{\hat{N}} - \hat{\phi}_k \right) \frac{\hat{Y}_k}{\hat{N}_k} = \hat{Y}_R - \sum_{k=1}^K \hat{\phi}_k \frac{\hat{Y}_k}{\hat{N}_k}$$

Using the strategy given above, the resulting estimator for  $\hat{\mu}$  is

$$\hat{Y}_1 = \hat{Y}_R - \hat{B}_1 = \sum_{k=1}^K \hat{\phi}_k \frac{\hat{Y}_k}{\hat{N}_k}$$

This is the "post-stratified" estimator usually found in practice. It is straightforward to verify the following properties of  $\hat{Y}_1$ :

$$E\left[\hat{Y}_1 \mid \hat{N}\right] = \hat{\mu} + \left[\mathbf{p}'\mathbf{P}(\hat{N} - \mathbf{M}_1)\right] + o(m^{-1})$$

$$\text{var}\left[\hat{Y}_1 \mid \hat{N}\right] = m^{-1}[\mathbf{p}'\mathbf{V}_c\mathbf{p}] + o(m^{-3/2})$$

$$E\left[\hat{Y}_1\right] = \hat{\mu} + o(m^{-1}) \text{ and}$$

$$\text{var}\left[\hat{Y}_1\right] = m^{-1}\left[\mathbf{p}'\left[\Sigma_{11} - 2\mathbf{D}(\mu_k)\Sigma_{21} + \mathbf{D}(\mu_k)\Sigma_{22}\mathbf{D}(\mu_k)\right]\mathbf{p}\right] + o(m^{-3/2})$$

$$\text{where } \mathbf{p}' = \left[\hat{\phi}_1/\phi_1, \hat{\phi}_2/\phi_2, \dots, \hat{\phi}_K/\phi_K\right].$$

The attempt to correct for the defective frame bias is successful in the sense that  $\hat{Y}_1$  is unconditionally unbiased for  $\hat{\mu}$ . However, the conditional bias is still present.

**Case 2.** For this case it can be verified that the estimator

$$\hat{B}_2 = (\mathbf{1} - \mathbf{p})' \left[ \hat{Y} - \Sigma_{12}\Sigma_{22}^{-1}(\hat{N}/\hat{N} - \mathbf{M}_2) \right]$$

is approximately, conditionally unbiased for  $\mu - \hat{\mu}$

and, as  $\hat{Y}_{LR}$  is conditionally unbiased for  $\mu$ , it follows directly that the estimator

$$\hat{Y}_2 = \hat{Y}_{LR} - \hat{B}_2 = \mathbf{p}' \left[ \hat{Y} - \Sigma_{12}\Sigma_{22}^{-1}(\hat{N}/\hat{N} - \mathbf{M}_2) \right]$$

is both conditionally and unconditionally, approximately unbiased for  $\hat{\mu}$ . It can also be verified that

$$\text{var}\left[\hat{Y}_2 \mid \hat{N}\right] = \text{var}\left[\hat{Y}_2\right] = m^{-1}[\mathbf{p}'\mathbf{V}_c\mathbf{p}].$$

In addition to the problems of the linear regression estimator cited earlier, this estimator is usually not even well defined as the frame parameters  $\{\phi_k, 1 \leq k \leq K\}$  are rarely, if ever, known when the frame is defective.

#### 5. CONCLUSION

This study has generalized the asymptotic techniques suggested by Robinson (1987) to study the problem of post-stratification from a design-based, conditional point-of-view. From a conditional point of view the linear regression estimator is preferable among the four studied here. Only the regression estimator is conditionally unbiased. The post-stratified estimator is no better (or worse) than either the Horvitz-Thompson or the ratio estimator; all have conditional bias terms of order  $m^{-(1/2)}$ . All of the estimators have the same conditional variance to terms of order  $m^{-1}$ ; furthermore, the conditional variance does not depend on  $\hat{N}$ , the vector of estimated proportions in the post-strata. Consequently, because of its conditional unbiasedness, the regression estimator has the smallest conditional mean square error.

The problem of a defective frame introduces complications not found otherwise. Each of the

estimators of the mean studied here is biased both conditionally and unconditionally. Bias adjustments are possible only under the restrictive assumption that the mean of units within each post-stratum is the same for all population units whether they are included or excluded from the frame.

### 6. REFERENCES

Bailar, B. (1989), "Information Needs, Surveys, and Measurement Errors," in *Panel Surveys*, eds D. Kasprzyk, G. Duncan, G. Kalton, and M.P. Singh. New York: Wiley.

Durbin, J. (1969), "Inferential Aspects of Randomness of Sample Size in Survey Sampling," in *New Developments in Survey Sampling*, N.L. Johnson and H. Smith, eds. New York: Wiley.

Fuller, W. A. (1981), "Comment" on "An Empirical Study of the Ratio Estimator and Estimators of its Variance," by R.M. Royall and W.G. Cumberland, *Journal of the American Statistical Association*, 76, 78-80.

Hansen, M.H., Madow, W.G., and Tepping, B.J. (1983), "An Evaluation of Model-Dependent and Probability-Sampling Inferences in Sample Surveys," *Journal of the American Statistical Association*, 78, 776-796.

Holt, D. and Smith, T.M.F. (1979), "Post Stratification," *Journal of the Royal Statistical Society A*, 142, 33-46.

Krewski, D. and Rao, J.N.K. (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. (1991), "Post-Stratification: A Modeler's Perspective," *Proceeding of the Section on Survey Methods Research*, Washington: American Statistical Association, in press.

Rao, J.N.K. (1985), "Conditional Inference in Survey Sampling," *Survey Methodology*, 11, 15-31.

\_\_\_\_\_ (1992), "Estimating Totals and Distribution Functions Using Auxiliary Information at the Estimation Stage," presented at Workshop on Uses of Auxiliary Information in Surveys, Statistics Sweden.

Rao, J.N.K. and Wu, C.F.J. (1985), "Inference from Stratified Samples: Second Order Analysis of Three Methods for Nonlinear Statistics," *Journal of the American Statistical Association*, 80, 620-630.

Robinson, J. (1987), "Conditioning Ratio Estimates Under Simple Random Sampling," *Journal of the American Statistical Association*, 82, 826-831.

Särndal C.-E., Swensson, B., and Wretman, J. (1992), *Model Assisted Survey Sampling*, New York: Springer-Verlag.

Valliant, R. (1993), "Post-stratification and Conditional Variance Estimation," *Journal of the American Statistical Association*, 88, in press.

Table 1. Population size and basic sample design parameters for three simulation studies.

Population	Pop. Size $N$	No. of FSUs $M$	No. of sample FSUs $m$
HIS	2,934	1,100	115
Artificial	22,001	2,000	200

Table 2. Simulation results for three populations.

Estimator	Rel-bias $\hat{Y}$ (%)	$rms_e(\hat{Y})$	$100 \cdot \left[ \frac{rms_e(\hat{Y})}{rms_e(\hat{Y}_{PS})} - 1 \right]$
<b>HIS population</b>			
$\hat{Y}_{HT}$	.12	.141	.05
$\hat{Y}_R$	.10	.141	.02
$\hat{Y}_{PS}$	.11	.141	0
$\hat{Y}_{LR}(\text{emp})$	.19	.162	14.71
$\hat{Y}_{LR}(\text{theo})$	.08	.140	-9.6
<b>Artificial population</b>			
$\hat{Y}_{HT}$	.02	2.30	-2.93
$\hat{Y}_{PS}$	.12	2.37	0
$\hat{Y}_{LR}(\text{emp})$	.04	2.31	-2.41
$\hat{Y}_{LR}(\text{theo})$	.02	2.26	-4.70

Figure 1. NHIS simulation,  $m=115$

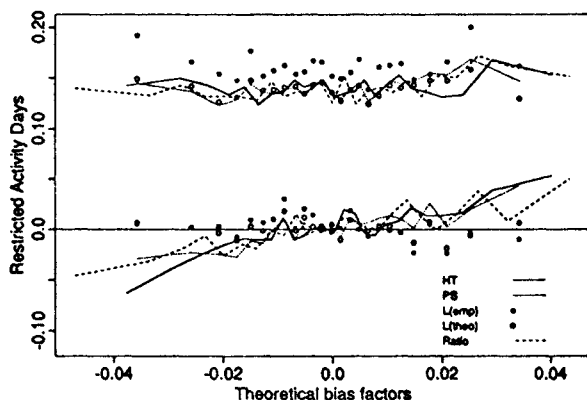


Figure 2. Artificial population simulation,  $m=200$

