# Multiple Workloads per Stratum Designs

*Lynn Weidman[1] and Lawrence R. Ernst[2]*

This article introduces an approach to expanding a stratified sample design, $D_1$, with one primary sampling unit (PSU) selected per stratum to a larger design, $D_2$. Define a workload (WL) to be the sample size in a given stratum in $D_1$. The three-stage approach selects the number of WLs for each stratum, the PSUs to receive the additional WLs in each stratum, and the ultimate sampling units. Procedures are given for selecting PSUs in the key second stage, satisfying the following conditions when a stratum in $D_2$ is to have $s \geq 2$ WLs: (i) the expected number of WLs in a PSU is $s$ times the probability that it was selected to get the single WL in $D_1$; and (ii) the actual number of WLs assigned is within one of the expected number. These conditions are a generalization of probability proportional to size, without replacement sampling. The properties and variances of this approach are compared to those from three alternative expansion procedures via application to a proposed, but since cancelled, expansion of the Current Population Survey.

*Key words:* Stratified sample design; PSU selection; workload; variance decomposition.

## 1. Introduction

The sample expansion procedure presented here was motivated by a planned (but since cancelled) expansion of the U.S. Census Bureau's Current Population Survey (CPS) that was to be selected in two phases. Phase 1 would be the traditional sample redesign of the CPS based on recently available decennial census data. This design must meet monthly variance requirements on estimates of the number of persons unemployed for the nation, the eleven largest states, New York City and Los Angeles. The remaining states have annual variance requirements. (Denote this first phase design as $D_1$.) At a later date, additional sample would be selected from these remaining states in phase 2 to meet monthly state variance requirements. (Denote this second phase design as $D_2$.) Because of the expenses incurred in recruiting and training new interviewers, and the time required to develop their skills, it is desirable to maintain as many $D_1$ sample primary sampling units (PSUs) as possible for $D_2$. (For this reason, when redesigning the U.S. Census Bureau's major household surveys each decade, procedures to maximize overlap of PSUs between the old and new designs are used when feasible.) Although the CPS application was the

motivation for the procedure to be discussed, there are potential applications to other sample expansion problems as well.

For the CPS most of the states are divided into PSUs. PSUs smaller than a specified size (population), the noncertainty PSUs, are combined into strata of approximately equal size within states. One PSU per stratum is then selected with probability proportional to size. Ultimate sampling units (USUs), which are small clusters of adjacent housing units, are then sampled from each of the stratified PSUs in sample and the non-stratified (certainty) PSUs. This is done in such a manner that each USU within a state has the same overall selection probability. (A few states are treated as one large non-stratified PSU.) Our interest is in how to increase sample in the stratified PSUs when the total sample for the state must be increased to meet the $D_2$ variance requirements. (The sample in the non-stratified PSUs is increased in a straightforward manner by selecting additional USUs.)

In formulating what we call the multiple workloads (WLs) procedure, we were guided by the desirability of obtaining a sampling scheme with the following properties:

1. All $D_1$ sample PSUs are $D_2$ sample PSUs.
2. The $D_2$ sample PSUs may be selected subsequent to the selection of the $D_1$ sample PSUs, which would allow the procedure to be used if it is necessary to expand the sample after the $D_1$ PSUs are selected.
3. The between-PSUs variance should be low.
4. The within-PSUs variance should be low.
5. Variance estimation should not be unduly complex.

The most natural approach to satisfying property 3 is to select the $D_2$ PSUs from an optimal $D_2$ stratification. Two procedures for doing this which have been investigated at the U.S. Census Bureau, independent sample (Chandhok, Weinstein, and Gunlicks 1990) and controlled selection (Ernst 1990), are discussed in Section 6. Unfortunately, independent sample violates property 1 and controlled selection violates property 2.

Consequently, we developed the multiple WLs approach. Let the number of USUs selected from each sample PSU in a $D_1$ stratum, called a workload (WL), be the number that can be efficiently handled by a single interviewer during an interview period. The multiple WLs approach selects PSUs to receive additional full WLs in $D_2$. For each stratum $h$, let the number of PSUs (not necessarily distinct) selected for $D_2$ be denoted by $n_h$. The $D_2$ sample PSUs from stratum $h$ consist of the single $D_1$ sample PSU and $n_h$-1 additional PSUs selected to satisfy:

i.  the expected number of times a PSU is selected for $D_2$ is $n_h$ times the probability that it was selected for $D_1$; and
ii. the number of times that a PSU is selected for $D_2$ is within one of the expected number in (i).

*Example.* Consider a stratum $h$ with four PSUs and $n_h = 5$. If the $D_1$ selection probabilities for these PSUs are .48, .24, .16 and .12, then by condition (i) their expected numbers of selections in $D_2$ are 2.4, 1.2, .8 and .6. By condition (ii) the number of times each PSU can be selected for $D_2$ is 2 or 3 for the largest one, 1 or 2 for the second largest, and 0 or 1 for the smallest two.

Condition (i) is the general probability proportional to size criterion. Condition (ii) is a generalization of without replacement sampling to include the case when the expected number of times a PSU is selected in $D_2$ exceeds 1; that is, when the probability of its selection in the $D_1$ sample is greater than $1/n_h$. Condition (ii) is motivated by the fact that minimizing the variability in the number of times that a PSU is selected tends to lower between-PSU variances, the same motivation for preferring without replacement sampling to with replacement sampling.

The desire to satisfy properties 1 and 2 motivated the approach of expanding the sample by selecting additional WLs from the $D_1$ strata, and property 3 motivated the use of conditions (i) and (ii) in the selection of the additional WLs. A different approach to satisfying properties 1 and 2, independent supplement (Chandhok, Weinstein, and Gunlicks 1990), has also been investigated at the U.S. Census Bureau and is described in Section 6. However, in contrast with the multiple WLs approach, under independent supplement small PSUs can receive more than one WL and, as a result, it is not as successful as the multiple WLs approach at satisfying property 3.

In most applications, including the CPS, the selection of PSUs to receive additional WLs under the multiple WLs approach constitutes the second stage of a three-stage sampling process. The final stage, that of selecting the USUs, satisfies the condition that the probability of selection of each USU in the state is the same, allowing for self-weighting estimates. This condition was motivated by property 4, since in certain situations a self-weighting estimator will minimize within-PSUs variance and total variance (Cochran 1977, Sec. 10.10). Variance formulae for this combination of USU sampling and self-weighting estimator are derived in Section 4.

The first stage of the multiple WLs procedure is the random determination of the number of additional WLs to be selected from each $D_1$ stratum for the $D_2$ design. Note that, as a result of this process, the number of WLs in each stratum is a random variable, and there will be a between stratum component of variance when our sampling process is combined with self-weighting estimates. Thus even our recommended sampling and estimation combination does not satisfy property 5. Consequently, in Section 5 we introduce an unbiased estimation procedure as an alternative to self-weighting, which in combination with multiple WLs does not have a between-strata component of variance, and allows us to satisfy property 5. We will show analytically in Section 5 and empirically in Section 6 that there is a trade-off, however, since this alternative estimator generally yields larger within-PSUs variances than the self-weighting estimator.

In the remainder of this article we look at the details of the multiple WLs procedure. Section 2 describes all three stages of the procedure for expanding from $D_1$ to $D_2$. Section 3 details the procedures for selecting PSUs to receive the additional WLs, which is the only sampling stage for which the selection methodology is not routine. Variance formulae are derived in Section 4 for the self-weighting estimator. In Section 5 variance formulae are obtained for an alternative estimator having no between-strata component of variance. Finally, Section 6 presents, as an example, the state variances of the self-weighting multiple WLs estimator for the CPS application and compares them with those from the three other options investigated at the U.S. Census Bureau and the alternative multiple WLs estimator. The multiple WLs per stratum method is the only one of the four methods that have all three of the following desirable properties: (a) the $D_2$ sample PSUs need not be selected

at the same time as the $D_1$ sample PSUs, but can be selected at a later date; (b) all $D_1$ PSUs are retained in $D_2$; and (c) only large PSUs can receive more than one WL.

## 2. Expanding an Existing Design

The presentation in this and the next three sections considers the expansion from $D_1$ to $D_2$ for stratified PSUs only. For each non-stratified PSU the expansion is obtained by selecting an appropriate number of USUs to supplement the $D_1$ sample.

Expansion from $D_1$ to $D_2$ by the multiple WLs approach usually involves a three-stage sampling process. First, the total number of WLs are allocated among the strata. Then the WLs in each stratum are allocated among its PSUs. Finally, USUs which make up the WLs are selected within the designated PSUs. We now proceed to describe each of these stages. The details of the multiple WLs approach to PSU selection will be given in Section 3.

### 2.1. Allocation of workloads to strata

Recall from Section 1 that in most states the Current Population Survey is a multi-stage design with each USU having the same overall probability of selection. The PSUs are stratified so that the strata are of approximately the same size and, as a result, the number of USUs in sample in each stratum will be close to a specified target WL. This target is the number of USUs that can be most efficiently handled by a single interviewer during an interviewing period. In order to minimize the number of interviewers needed for $D_2$ and assure that each is able to efficiently handle his/her WL, a large increase in total sample size should be made by adding an integer number of WLs to each stratum. This can only be accomplished, while simultaneously adding the minimal number of WLs required overall for $D_2$, by allowing the number of WLs to vary among strata.

How should WLs be allocated to each stratum? One fairly straightforward method is to let the number of WLs vary by at most one among the strata, and select a simple random sample of strata to receive the higher number. (This sampling of strata in combination with a self-weighting estimator has a between-strata component of variance. However, by limiting the variation in the number of WLs each stratum can be allocated, this between-strata component will be reduced.) We describe this method notationally and then present a simple example.

Let $m^*$ denote the number of $D_1$ sample USUs and $L$ the number of strata, so that $m^*/L$ is the expected average WL size in $D_2$. If $m$ is the desired number of $D_2$ sample USUs, then the minimal number of expected average WLs needed to attain this sample size is $mL/m^*$. If this is an integer, then this is the number of $D_2$ WLs, denoted by $n$; otherwise, we round up the number of $D_2$ WLs to the next integer $\lfloor mL/m^* \rfloor + 1 = n$. (For any number x, let $\lfloor x \rfloor$ denote the greatest integer not exceeding $x$.) Now $R = n/L$ is the average number of WLs per stratum in $D_2$, so each stratum is to be allocated either $\lfloor R \rfloor$ or $\lfloor R \rfloor + 1$ WLs. This is accomplished by selecting a simple random sample of $n - L\lfloor R \rfloor$ strata to receive $\lfloor R \rfloor + 1$ WLs and the rest to receive $\lfloor R \rfloor$ WLs. (If $R$ is an integer, $n - L\lfloor R \rfloor = 0$ and all strata receive $R$ WLs.) Letting $n_h$ denote the number of WLs assigned to the $h$th stratum by the first stage of sampling, note that

$$E(n_h) = R, \qquad h = 1, \ldots, L \tag{2.1}$$

*Example.* Suppose a state has $L = 4$ strata and $m^* = 160$ USUs in sample for $D_1$. If $m = 368$ USUs are needed for $D_2$, then the smallest integral number of expected average WLs of size $160/4 = 40$ that will achieve this size is $n = \lfloor (4)(368/160) \rfloor + 1 = 10$. Then $R = 10/4 = 2.5$ is the average number of WLs per stratum, and $10\text{-}(4)\lfloor 2.5 \rfloor = 2$ strata will be chosen to receive $\lfloor 2.5 \rfloor + 1 = 3$ WLs. The remaining two will receive $\lfloor 2.5 \rfloor = 2$.

## 2.2. Allocation of workloads to PSUs within a stratum

How should the $n_h$ WLs be allocated to the PSUs in stratum $h$? Let $p_{hi}$ be the probability of selection of the $i$th PSU in the $h$th stratum in $D_1$ and $n_{hi}$ be the number of WLs allocated to this PSU in $D_2$. Typically, when $n_h$ PSUs are selected without replacement, each PSU is selected either once with expected probability $n_h p_{hi}$ or not selected (Cochran 1977; Sampford 1967). Variances are smaller for these procedures than if a PSU is allowed to be selected more than once. However, the cited authors also assume that all $p_{hi} < 1/n_h$. When this is not the case, we generalize the restriction on the number of times a PSU can be selected to conditions (ii) of Section 1. The multiple WLs approach has the following features, where

$$\pi_{hi} = P(n_{hi} = \lfloor n_h p_{hi} \rfloor + 1)$$

We see that

$$n_{hi} \geq 1 \text{ if PSU } hi \text{ is in } D_1 \qquad (2.2)$$

which is equivalent to design property (1) of the previous section;

$$n_{hi} = \lfloor n_h p_{hi} \rfloor \text{ or } n_{hi} = \lfloor n_h p_{hi} \rfloor + 1 \qquad (2.3)$$

which is condition (ii)

$$\pi_{hi} = n_h p_{hi} - \lfloor n_h p_{hi} \rfloor \qquad (2.4)$$

$$\sum_{i=1}^{N_h} n_{hi} = n_h \qquad (2.5)$$

where $N_h$ = number of PSUs in stratum $h$, which simply says that the total number of WLs in stratum $h$ is $n_h$.

The combination of (2.3) and (2.4) imply condition (i), that is,

$$E(n_{hi}|n_h) = n_h p_{hi} \qquad (2.6)$$

If $n_h = 1$, let the $D_2$ sample PSU for stratum $h$ be the $D_1$ sample PSU, since this allocation satisfies (2.2)–(2.5). For $n_h \geq 2$ a procedure for allocating WLs that satisfies (2.2)–(2.5) is presented in Section 3.

## 2.3. Allocation of USUs within a PSU and total sample size

Let $M$, $M_h$, and $M_{hi}$ be the number of USUs in the total population, the $h$th stratum, and the $hi$th PSU, respectively. We assume that $p_{hi} = M_{hi}/M_h$ and that the USUs within sample PSUs are selected in a manner such that each USU in the population has the same probability of selection, that is $m/M$. Consequently, the WL size for $D_2$ in the $h$th stratum,

denoted $m_h$, is

$$m_h = \frac{mM_h}{MR} \tag{2.7}$$

(Unless $m_h$ is already an integer, we round it up to the next integer when defining WL size.) Within each PSU $hi$ in the $h$th stratum for which $n_{hi} > 0$, $n_{hi}m_h$ USUs are selected with equal probability.

For the three-stage sampling procedure just outlined, the expected number of $D_2$ sample USUs is $m$, but the actual number selected is a random variable that depends on the first stage of sampling. This is because, by (2.7), the WL size is not the same for each stratum.

*Example.* Here we attach sizes to the strata of the previous example and calculate the range in total sample size that is possible due to the random selection of strata to receive additional WLs in $D_2$. Recall that $m = 368$ and $R = 2.5$. If the $M_h$ in the four strata are 88,000, 80,000, 78,000 and 74,000, then M = 320,000 and the $m_h$s are 41, 37, 36, and 35. The number of USUs in $D_2$ can range from $(3)(35 + 36) + (2)(37 + 41) = 369$ to $(3)(37 + 41) + (2)(35 + 36) = 376$. (Note that the average WL size of 37.25 here is smaller than the expected average WL size of 40 in the previous example. This is because rounding up the number of WLs to the nearest integer allows a slight decrease in WL size to attain a specified number of USUs in $D_2$.)

An alternative three-stage procedure for which the number of sample USUs would always be $m$ would begin by selecting the number of WLs, $n_h$, assigned to the $h$th stratum so that $E(n_h) = nM_h/M$, that is proportional to size. For every sample, $n_h$ would also be within 1 of $E(n_h)$. The second and third stages would be selected as described above, except now the WL size in each stratum would be $m/n$. This alternative approach will not be discussed further in this article.

## 3.  Selecting PSUs to Receive Additional Workloads

In this section our principal goal is to obtain probabilities for the sets of WLs selected from stratum $h$ for the $D_2$ design, conditioned on the $D_1$ sample PSU from this stratum. These conditional probabilities will satisfy (2.2) and yield unconditional selection probabilities for the $D_2$ WLs which satisfy (2.3)–(2.5).

In actuality, we work backwards. We first demonstrate how to obtain unconditional selection probabilities for the sets of $D_2$ WLs satisfying (2.3)–(2.5) and then specify conditional selection probabilities which yield these unconditional selection probabilities and satisfy (2.2). To this end, let $s_h$ denote a random vector of length $N_h$ which specifies the number of $D_2$ WLs to be selected from each PSU in stratum $h$. That is, $s_h$ is of the form $s_h = (n_{h1}, n_{h2}, ..., n_{hN_h})$, where the elements in $s_h$ satisfy (2.3) and (2.5). Let

$$T_{hi} = \{s_h : n_{hi} = \lfloor n_h p_{hi} \rfloor + 1\}, \qquad i = 1, ..., N_h \tag{3.1}$$

$$n'_h = n_h - \sum_{i=1}^{N_h} \lfloor n_h p_{hi} \rfloor \tag{3.2}$$

We seek a set of probabilities $P(s_h)$ for the $s_h$s which satisfy (2.4), that is for which

$$\pi_{hi} = \sum_{s_h \in T_{hi}} P(s_h) = n_h p_{hi} - \lfloor n_h p_{hi} \rfloor \qquad i = 1, ..., N_h \tag{3.3}$$

Observe that by (3.2), $n'_h$ of the $N_h$ elements of $s_h$ are in $T_{hi}$ for each $s_h$. That is, each $s_h$ corresponds to a set of $n'_h$ PSUs out of the $N_h$ PSUs that receive the larger of the two possible numbers of WLs for the PSUs in $D_2$. Consequently, *the problem of obtaining a set of $P(s_h)$'s which satisfy (3.3) is equivalent to the problem of finding joint probabilities of selecting without replacement $n'_h$ of the stratum h PSUs, for which the probability of selection of the PSU i is $\pi_{hi}$, $i = 1, ..., N_h$. If $n'_h = 1$ this can be accomplished by simply selecting a single PSU from the stratum with probability $\pi_{hi}$, $i = 1, ..., N_h$. If $n'_h \geq 2$ then the joint selection probabilities for the sets of $n'_h$ PSUs can be obtained by using any of a number of procedures for sampling without replacement with probability proportional to size.* The best known of these are the procedures of Brewer-Durbin (Cochran 1977) for $n'_h = 2$ and its generalization by Sampford (1967) for $n'_h \geq 3$.

We have thus shown how to obtain unconditional probabilities $P(s_h)$ for selecting the number of $D_2$ WLs from each PSU in stratum $h$. We then let the conditional probability that $s_h$ is selected for $D_2$ given that the $hi$th PSU was in the $D_1$ sample, denoted $P(s_h|i)$ be simply

$$P(s_h|i) = \frac{n_{hi}P(s_h)}{n_h p_{hi}} \tag{3.4}$$

where $n_{hi}$ is the $i$th element of $s_h$. Note that (3.4) satisfies (2.2), since if $n_{hi} = 0$ then $P(s_h|i) = 0$.

We show that the conditional probabilities (3.4) yield $P(s_h)$ as the unconditional probability for selecting each $s_h$. This follows since

$$\sum_{i=1}^{N_h} p_{hi}P(s_h|i) = \sum_{i=1}^{N_h} \frac{n_{hi}}{n_h}P(s_h) = P(s_h)$$

Consequently, by (3.3), these conditional selection probabilities result in unconditional probabilities satisfying (2.3)–(2.5).

Finally, we must demonstrate that for each $i = 1, ..., N_h$, the conditional probabilities sum to 1, that is $\Sigma_{s_h} P(s_h|i) = 1$, or equivalently, by (3.4)

$$\sum_{s_h} n_{hi}P(s_h) = n_h p_{hi}, \qquad i = 1, ..., N_h \tag{3.5}$$

To obtain (3.5) let

$$T_{hi}^* = \{s_h : n_{hi} = \lfloor n_h p_{hi} \rfloor\}, \qquad i = 1, ..., N_h$$

and observe that by (3.3)

$$\sum_{s_h} n_{hi}P(s_h) = (\lfloor n_h p_{hi} \rfloor + 1) \sum_{s_h \in T_{hi}} P(s_h) + \lfloor n_h p_{hi} \rfloor \sum_{s_h \in T_{hi}^*} P(s_h)$$

$$= (\lfloor n_h p_{hi} \rfloor + 1)(n_h p_{hi} - \lfloor n_h p_{hi} \rfloor) + \lfloor n_h p_{hi} \rfloor (1 - n_h p_{hi} + \lfloor n_h p_{hi} \rfloor) = n_h p_{hi}$$

*Example.* This example illustrates the procedure just described. Consider a stratum $h$ with $N_h = 4$ and $p_{hi}$s given in Table 1. Let $n_h = 5$. Then the values of $n_h p_{hi}$ and $\pi_{hi}$ are as given in Table 1.

Furthermore, $n'_h = 2$ and we can, therefore, use the Brewer-Durbin procedure to obtain pairwise joint probabilities $\pi_{hij}$, $i, j = 1, ..., 4$, $i < j$, where $\pi_{hij}$ denotes the probability

*Table 1.* $p_{hi}$, $n_h p_{hi}$, and $\pi_{hi}$

| | $i$ | | | |
|---|---|---|---|---|
| | 1 | 2 | 3 | 4 |
| $p_{hi}$ | .48 | .24 | .16 | .12 |
| $n_h p_{hi}$ | 2.4 | 1.2 | .8 | .6 |
| $\pi_{hi}$ | .4 | .2 | .8 | .6 |

*Table 2.* $P(s_h)$

| | $i,j$ | | | | | |
|---|---|---|---|---|---|---|
| | 1,2 | 1,3 | 1,4 | 2,3 | 2,4 | 3,4 |
| $s_h$ | (3,2,0,0) | (3,1,1,0) | (3,1,0,1) | (2,2,1,0) | (2,2,0,1) | (2,1,1,1) |
| $\pi_{hij} = P(s_h)$ | .0277 | .2535 | .1188 | .1188 | .0535 | .4277 |

*Table 3.* $P(s_h|i)$

| | $s_h$ | | | | | |
|---|---|---|---|---|---|---|
| $i$ | (3,2,0,0) | (3,1,1,0) | (3,1,0,1) | (2,2,1,0) | (2,2,0,1) | (2,1,1,1) |
| 1 | .0346 | .3169 | .1485 | .0990 | .0446 | .3564 |
| 2 | .0462 | .2113 | .0990 | .1980 | .0892 | .3564 |
| 3 | .0000 | .3169 | .0000 | .1485 | .0000 | .5346 |
| 4 | .0000 | .0000 | .1980 | .0000 | .0892 | .7128 |

that PSUs $hi$ and $hj$ are the two PSUs that receive the maximum number of $D_2$ WLs. This results in the probabilities $P(s_h)$ given in Table 2.

Then, from Table 2 and (3.4), the conditional probabilities $P(s_h|i)$ given in Table 3 are obtained.

*Remark.* We have demonstrated how we can select any additional number of WLs for $D_2$ in a stratum satisfying (2.2)–(2.5) when initially one PSU had been selected in $D_1$. A natural question is whether an expansion to a $D_2$ design satisfying (2.2)–(2.5) can also always be obtained if the $D_1$ design was other than a one-PSU-per-stratum design. The following simple example illustrates that this is not always possible. Consider a two-PSUs-per-stratum without replacement design, $D_1$, in which $N_h = 4$ for stratum $h$, with $p_{h1} = .45$, $p_{h2} = .40$, $p_{h3} = .10$, $p_{h4} = .05$. If $n_h = 3$ for $D_2$, then $n_h p_{h1} = 1.35$, $n_h p_{h2} = 1.20$ and hence PSUs $h1$ and $h2$ must receive at least 1 WL in $D_2$. However, if PSUs $h3$ and $h4$ had been the $D_1$ sample PSUs, then since only 1 additional WL is selected for the $D_2$ design, it would not be possible for both PSUs $h1$ and $h2$ to meet the minimum WL requirement.

## 4. Variance Decomposition for a Self-Weighting Estimator

In Sections 4 and 5 we develop variance formulae for two unbiased estimators of a population characteristic total $Y$, based on the three-stage sampling procedure of this article. The estimator of this section ($\hat{Y}$) is self-weighting and has a between-strata component of variance, while in the following section we eliminate this variance component by using an alternative estimator ($\hat{Y}'$) that is not self-weighting.

Recall from Section 2 that the final stage of sampling is carried out so that all USUs have equal overall probability of selection. Irrespective of how this is done, the unbiased Horvitz-Thompson estimator $\hat{Y}$ of $Y$ is given by

$$\hat{Y} = \frac{M}{m} \sum y_{hij} \qquad (4.1)$$

where $y_{hij}$ is the total value for characteristic $Y$ for a sample USU from PSU $hi$, and the summation is over all sample USUs. Some $y_{hij}$s may possibly appear more than once in this summation if with replacement sampling is used at the final stage. Such an estimator is called self-weighting because of the constant coefficient or "weight" $M/m$ for all USUs in (4.1). We proceed to develop a formula for $V(\hat{Y})$, the variance of the estimator $\hat{Y}$. There will be three terms in $V(\hat{Y})$, reflecting the three stages of sampling.

We first obtain an alternative expression for $\hat{Y}$. Let

$$\hat{Y}_{hi} = M_{hi} \sum_j y_{hij}/(n_{hi}m_h) \text{ if } n_{hi} > 0$$

$$= 0 \text{ if } n_{hi} = 0 \qquad (4.2)$$

Given $n_{hi} > 0$, $\hat{Y}_{hi}$ is an unbiased estimator of the total, denoted $Y_{hi}$, for PSU $hi$. Then combining (2.7), (2.8), and (4.2) we obtain

$$\hat{Y} = \frac{M}{m} \sum_{hij} y_{hij} = \frac{M}{m} \sum_{h=1}^{L} m_h \sum_{i=1}^{N_h} \frac{n_{hi}\hat{Y}_{hi}}{M_{hi}}$$

$$= \frac{1}{R} \sum_{h=1}^{L} M_h \sum_{i=1}^{N_h} \frac{n_{hi}\hat{Y}_{hi}}{M_{hi}} = \sum_{h=1}^{L} \sum_{i=1}^{N_h} a_{hi}\hat{Y}_{hi} \qquad (4.3)$$

where

$$a_{hi} = \frac{n_{hi}M_h}{RM_{hi}} = \frac{n_{hi}}{Rp_{hi}} \qquad (4.4)$$

is a random variable whose value for each sample is determined by the first two sampling stages.

The variance of $\hat{Y}$ can then be written in the form

$$V(\hat{Y}) = V_1 E_2 E_3 \left( \sum_{h=1}^{L} \sum_{i=1}^{N_h} a_{hi}\hat{Y}_{hi} \right) + E_1 V_2 E_3 \left( \sum_{h=1}^{L} \sum_{i=1}^{N_h} a_{hi}\hat{Y}_{hi} \right)$$

$$+ E_1 E_2 V_3 \left( \sum_{h=1}^{L} \sum_{i=1}^{N_h} a_{hi}\hat{Y}_{hi} \right) \qquad (4.5)$$

The subscripts on the expectations denote the three stages of sampling. Since whenever $n_{hi} > 0$, $E_3(\hat{Y}_{hi}) = Y_{hi}$, then

$$V(\hat{Y}) = V_1 \left[ \sum_{h=1}^{L} \sum_{i=1}^{N_h} E_2(a_{hi})Y_{hi} \right] + E_1 \left[ \sum_{h=1}^{L} V_2 \left( \sum_{i=1}^{N_h} a_{hi}Y_{hi} \right) \right]$$

$$+ E_1 \left[ \sum_{h=1}^{L} \sum_{i=1}^{N_h} E_2\{a_{hi}^2 V_3(\hat{Y}_{hi})\} \right]. \qquad (4.6)$$

These three terms are, respectively, the between-strata, between-PSUs-within-strata, and within-PSUs components of variance.

### 4.1. Between-strata variance

Let $Y_h = \sum_{i=1}^{N_h} Y_{hi}$, the characteristic total for stratum $h$. Let $S$ denote the set of strata $h$ for which $n_h = \lfloor R \rfloor + 1$ and let $\ell = n - L \lfloor R \rfloor$, the number of such strata. If $\ell = 0$, then the between-strata component of variance is 0. Otherwise, by (2.6), (4.4), and the fact that $\lfloor R \rfloor Y_h$ is a constant, we have

$$
V_1 \left[ \sum_{h=1}^{L} \sum_{i=1}^{N_h} E_2(a_{hi}) Y_{hi} \right] = V_1 \left( \sum_{h=1}^{L} \sum_{i=1}^{N_h} \frac{n_h}{R} Y_{hi} \right) = \frac{1}{R^2} V_1 \left( \sum_{h=1}^{L} n_h Y_h \right)
$$

$$
= \frac{1}{R^2} V_1 \left( \sum_{h=1}^{L} (n_h - \lfloor R \rfloor) Y_h \right) = \left( \frac{\ell}{R} \right)^2 V_1 \left( \sum_{h \in S} \frac{Y_h}{\ell} \right) \tag{4.7}
$$

Since the set $S$ is a simple random sample of $\ell$ out of $L$ strata, (4.7) is $(\ell/R)^2$ times the variance of the mean from a simple random sample. Using Theorem 2.2 of Cochran (1977), we obtain that (4.7) reduces to

$$
\frac{(L - \ell)\ell}{R^2 L (L - 1)} \left[ \sum_{h=1}^{L} \left( Y_h - \frac{Y}{L} \right)^2 \right] \tag{4.8}
$$

### 4.2. Between-PSUs-within-strata variance

To evaluate the second bracketed term in (4.6), first let $p^*(k) = P(n_h = k)$ and observe that for all $h$

$$
p^*(\lfloor R \rfloor) = 1 + \lfloor R \rfloor - R, \quad p^*(\lfloor R \rfloor + 1) = R - \lfloor R \rfloor \tag{4.9}
$$

and $p^*(k) = 0$ for all other $k$. Consequently

$$
E_1 \left[ \sum_{h=1}^{L} V_2 \left( \sum_{i=1}^{N_h} a_{hi} Y_{hi} \right) \right] = \sum_{h=1}^{L} \sum_{k=\lfloor R \rfloor}^{\lfloor R \rfloor + 1} V \left( \sum_{i=1}^{N_h} a_{hi} Y_{hi} | n_h = k \right) p^*(k) \tag{4.10}
$$

To evaluate the variance of the term in parentheses in (4.10), let $n'_{hi} = n_{hi} - \lfloor n_h p_{hi} \rfloor$ and expand

$$
V \left( \sum_{i=1}^{N_h} a_{hi} Y_{hi} | n_h = k \right) = \frac{1}{R^2} V \left( \sum_{i=1}^{N_h} n_{hi} \frac{Y_{hi}}{p_{hi}} \bigg| n_h = k \right) = \frac{1}{R^2} V \left( \sum_{i=1}^{N_h} n'_{hi} \frac{Y_{hi}}{p_{hi}} \bigg| n_h = k \right)
$$

$$
= \frac{1}{R^2} \left[ \sum_{i=1}^{N_h} V(n'_{hi} | n_h = k) \frac{Y_{hi}^2}{p_{hi}} \right.
$$

$$
\left. + \sum_{\substack{i=1 \\ i \neq j}}^{N_h} \sum_{j=1}^{N_h} Cov(n'_{hi}, n'_{hj} | n_h = k) \frac{Y_{hi}}{p_{hi}} \frac{Y_{hj}}{p_{hj}} \right] \tag{4.11}
$$

where the substitution of $n'_{hi}$ for $n_{hi}$ is justified by the fact that conditional on $n_h$ these two variables differ by a constant.

To evaluate $V(n'_{hi}|n_h = k)$ and $Cov(n'_{hi}, n'_{hj}|n_h = k)$, first observe that $n'_{hi}$ is a 0-1 variable and that $n'_{hi} = 1$ if and only if $n_{hi} = \lfloor n_h p_{hi} \rfloor + 1$. We then let

$$\pi_{hi}(k) = P(n'_{hi} = 1|n_h = k) = k p_{hi} - \lfloor k p_{hi} \rfloor \tag{4.12}$$

$$\pi_{hij}(k) = P(n'_{hi} = 1, \ n'_{hj} = 1|n_h = k) \tag{4.13}$$

where the last equation in (4.12) follows from (2.4). This is equivalent to the notation used in Section 3, except now we have converted $\pi_{hi}$, $\pi_{hij}$ to functions to indicate their dependence on the value of $n_h$. Furthermore, to compute $\pi_{hij}(k)$, observe that by (3.2), if $n_h = k$ then $n'_h = k - \Sigma_{i=1}^{N_h} \lfloor k p_{hi} \rfloor$. If $n'_h = 1$, then $\pi_{hij}(k) = 0$ for all $i,j$, while if $n'_h \geq 2$ then $\pi_{hij}(k)$ depends on the sampling method employed to select the number of WLs in stratum $h$. For example, if $n'_h = 2$ and the Brewer-Durbin method is used, then the $\pi_{hij}(k)$ are the joint probabilities for selecting pairs of units given in Cochran (1977, p. 262), while if $n'_h \geq 3$ and Sampford's (1967) method is used then the $\pi_{hij}(k)$ are the joint probabilities for pairs given in that reference.

From (4.12) and (4.13) we have, since $n'_{hi}$ is a 0-1 variable, that

$$V(n'_{hi}|n_h = k) = E(n'_{hi}|n_{hi} = k) - [E(n'_{hi}|n_{hi} = k)]^2 = \pi_{hi}(k) - \pi_{hi}^2(k) \tag{4.14}$$

and

$$Cov(n'_{hi}, n'_{hj}|n_h = k) = E(n'_{hi}n'_{hj}|n_h = k) = \pi_{hij}(k) - \pi_{hi}(k)\pi_{hj}(k). \tag{4.15}$$

Finally, we combine (4.10), (4.11), (4.14), and (4.15) to conclude

$$E_1 \left[ \sum_{h=1}^{L} V_2 \left( \sum_{i=1}^{N_h} a_{hi} Y_{hi} \right) \right] = \frac{1}{R^2} \sum_{h=1}^{L} \sum_{k=\lfloor R \rfloor}^{\lfloor R \rfloor + 1} \left[ \sum_{i=1}^{N_h} (\pi_{hi}(k) - \pi_{hi}^2(k)) \frac{Y_{hi}^2}{p_{hi}^2} \right.$$

$$\left. + \sum_{\substack{i=1 \\ i \neq j}}^{N_h} \sum_{j=1}^{N_h} (\pi_{hij}(k) - \pi_{hi}(k)\pi_{hj}(k)) \frac{Y_{hi} Y_{hj}}{p_{hi} p_{hj}} \right] p^*(k) \tag{4.16}$$

### 4.3. Within-PSUs variance

Finally, we evaluate the third bracketed term in (4.6) under the assumption that all sample WLs within a PSU are selected by simple random sampling, either with replacement for large $M_{hi}$ or without replacement with a negligible finite population correction factor. Appropriate modifications are necessary for other within-PSU selection procedures.

From (4.2) it follows that if $n_{hi} > 0$, then

$$V_3(\hat{Y}_{hi}) \doteq - \frac{M_{hi}^2 S_{3hi}^2}{n_{hi} m_h} \tag{4.17}$$

where $S_{3hi}^2$ is the population variance of the $y_{hij}$ in PSU $hi$.

Then by (4.4), (2.6), (2.7), and the relation $p_{hi} = M_{hi}/M_h$

$$E_2[a_{hi}^2 V_3(\hat{Y}_{hi})] \doteq \frac{E_2(n_{hi})M_{hi}^2 S_{3hi}^2}{R^2 p_{hi}^2 m_h} = \frac{n_h M_{hi}^2 S_{3hi}^2}{R^2 p_{hi} m_h} = \frac{n_h M M_{hi} S_{3hi}^2}{Rm} \tag{4.18}$$

which we combine with (2.1) to conclude

$$E_1\left[\sum_{h=1}^{L}\sum_{i=1}^{N_h}E_2(a_{hi}^2 V_3(\hat{Y}_{hi}))\right] \doteq \frac{M}{m}\sum_{h=1}^{L}\sum_{i=1}^{N_h}M_{hi}S_{3hi}^2 \tag{4.19}$$

Finally, if the $S_{3hi}^2$ are the same for all $hi$, with common value denoted $S_3^2$, then by summing (4.19) over all $hi$ we obtain that the within-PSUs variance is approximately $M^2 S_3^2/m$. This is approximately the sampling variance for the standard estimator of population total from a simple random sample with replacement, for a sample of size $m$ selected from a population of size $M$, for $M$ large, with variance $S_3^2$. Similar assumptions lead to the same approximate within-PSUs variance for the other options investigated for the two-phase sampling application, a result which will be used in some of the comparisons in Section 6.

## 5.  An Alternative Estimator of Y

The estimator $\hat{Y}$ defined by (4.1) was selected because it is an unbiased estimator of $Y$ and because it assigns each USU the same coefficient $M/m$. However, for the $D_2$ sample design, $\hat{Y}$ has the disadvantage that it results in a non-zero between-strata component of variance, which makes variance estimation much more complex. In this section we present an alternative unbiased estimator $\hat{Y}'$ of $Y$ for which the between-strata variance is 0. We will also obtain expressions for the between-PSUs-within-strata component and the within-PSUs component of $V(\hat{Y}')$, and compare the within-PSUs components of $V(\hat{Y})$ and $V(\hat{Y}')$.

Let

$$\hat{Y}' = \sum_{h=1}^{L}\frac{M_h}{n_h m_h}\sum_{i=1}^{N_h}\sum_{j}y_{hij} \tag{5.1}$$

Since, by (2.7)

$$\frac{M_h}{n_h m_h} = \frac{MR}{mn_h} \tag{5.2}$$

$\hat{Y}'$ assigns weights to the USUs which are a function of the number of $D_2$ WLs, $n_h$, in each stratum. Also let

$$a'_{hi} = \frac{Ra_{hi}}{n_h} \tag{5.3}$$

Then by (4.2), (5.1), (5.2), and (5.3)

$$\hat{Y}' = \sum_{h=1}^{L}\sum_{i=1}^{N_h}a'_{hi}\hat{Y}_{hi} \tag{5.4}$$

and $V(\hat{Y}')$ can be written in the form (4.6) with $a_{hi}$ replaced by $a'_{hi}$.

Furthermore, by (2.6), (4.4), and (5.3)

$$E_2(a'_{hi}) = \frac{R}{n_h}E_2(a_{hi}) = 1$$

and hence

$$E_2[E_3(\hat{Y}')] = \sum_{h=1}^{L} \sum_{i=1}^{N_h} E_2(a'_{hi})Y_{hi} = Y \tag{5.5}$$

From (5.5) it follows that $\hat{Y}'$ is an unbiased estimator of $Y$, and the between-strata component of $V(\hat{Y}')$ is 0.

The between-PSUs-within-strata component of $V(\hat{Y}')$ can be obtained by proceeding as in Section 4.2 with $a'_{hi}$ substituted for $a_{hi}$ and, by (5.3), the final expression in (4.11) modified by multiplying it by $R^2/k^2$. Hence, with these changes, we obtain, analogously to (4.16), that

$$E_1\left[\sum_{h=1}^{L} V_2\left(\sum_{i=1}^{N_h} a'_{hi}Y_{hi}\right)\right] = \sum_{h=1}^{L} \sum_{k=\lfloor R \rfloor}^{\lfloor R \rfloor+1} \left[\sum_{i=1}^{N_h}(\pi_{hi}(k) - \pi_{hi}^2(k))\frac{Y_{hi}^2}{p_{hi}^2}\right.$$
$$\left.+ \sum_{\substack{i=1 \\ i \neq j}}^{N_h} \sum_{j=1}^{N_h}(\pi_{hij}(k) - \pi_{hi}(k)\pi_{hj}(k))\frac{Y_{hi}}{p_{hi}}\frac{Y_{hj}}{p_{hj}}\right]\frac{p^*(k)}{k^2} \tag{5.6}$$

There is no clear relation between the values of (4.16) and (5.6). The expression within the brackets is multiplied by $1/R^2$ in (4.16) and $1/k^2$ in (5.6), and $1/k^2 \geq 1/R^2$ for $k = \lfloor R \rfloor$ while $1/k^2 < 1/R^2$ for $k = \lfloor R \rfloor + 1$. However, the expression within the brackets is a complex function of $k$. It is possible, for example, for this term to be larger when $k = \lfloor R \rfloor$ than when $k = \lfloor R \rfloor + 1$, and also for the opposite to be true.

To obtain the within-PSUs component of $V(\hat{Y}')$, first observe that by (5.3) and (4.18)

$$E_2[a'^2_{hi}V_3(\hat{Y}_{hi})] = \frac{R^2}{n_h^2}E_2[a_{hi}^2 V_3(\hat{Y}_{hi})] \doteq \frac{RM M_{hi}S_{3hi}^2}{n_h m} \tag{5.7}$$

Furthermore, by (4.9)

$$E_1\left(\frac{R}{n_h}\right) = R\sum_{k=\lfloor R \rfloor}^{\lfloor R \rfloor+1} \frac{p^*(k)}{k} = R\left(\frac{1+\lfloor R \rfloor - R}{\lfloor R \rfloor} + \frac{R - \lfloor R \rfloor}{\lfloor R \rfloor + 1}\right) = f(R) \tag{5.8}$$

where

$$f(R) = \frac{R(1 - R + 2\lfloor R \rfloor)}{\lfloor R \rfloor(\lfloor R \rfloor + 1)}. \tag{5.9}$$

We then combine (5.7) and (5.8) to obtain

$$E_1\left[\sum_{h=1}^{L} \sum_{i=1}^{N_h} E_2(a'^2_{hi}V_3(\hat{Y}_{hi}))\right] \doteq E_1\left(\frac{R}{n_h}\right)\left[\frac{M}{m}\sum_{h=1}^{L}\sum_{i=1}^{N_h} M_{hi}S_{3hi}^2\right] = f(R)\left[\frac{M}{m}\sum_{h=1}^{L}\sum_{i=1}^{N_h} M_{hi}S_{3hi}^2\right] \tag{5.10}$$

and that, using the same reasoning as in Section 4.3, if $S_{3hi}^2 = S_3^2$ for all $hi$, then the within-PSUs variance of $\hat{Y}'$ is approximately $f(R)M^2S_3^2/m$.

Finally, it follows from (5.10) and (4.19) that the ratio of the within-PSUs component of $V(\hat{Y}')$ to the within-PSUs component of $V(\hat{Y})$ is approximately $f(R)$. Thus, unlike the

between-PSUs-within-strata component, there is a relatively simple relationship between the within-PSUs component of variance for these two estimators.

The function $f$ behaves as follows: $f(R) = 1$ whenever $R = \lfloor R \rfloor$, that is, when $R$ is an integer. Furthermore, as can be shown by means of elementary calculus, within the interval bounded by any two consecutive integers, $f(R)$ has a single relative optimum, a relative maximum, at $\lfloor R \rfloor + 1/2$, with

$$f(\lfloor R \rfloor + 1/2) = 1 + \frac{1}{4\lfloor R \rfloor (\lfloor R \rfloor + 1)} \tag{5.11}$$

Thus $f(1) = 1$; $f$ increases in the interval $[1, 3/2]$, with $f(3/2) = 9/8$ by (5.11); $f$ decreases in the interval $[3/2, 2]$ with $f(2) = 1$; $f$ increases in the interval $[2, 5/2]$, with $f(5/2) = 25/24$; and so forth.

One of the comparisons in the next section will be of $V(\hat{Y})$ and $V(\hat{Y}')$ for the CPS expansion.

## 6. Comparison of Methods for the CPS Expansion

In this section variances for the multiple WLs per stratum method estimator $\hat{Y}$ are first compared to variances for three other methods for selecting the $D_2$ sample for the formerly planned CPS expansion, then $V(\hat{Y})$ is compared to $V(\hat{Y}')$. These three other methods are the independent sample, the independent supplement (both described in Chandhok, Weinstein, and Gunlicks 1990), and controlled selection (Ernst 1990). For each of these four methods for the study to be described, the number of USUs selected from the $D_1$ noncertainty PSUs for the $D_2$ design is the same, $m$. As was the case for the multiple WLs method, the expansion from $D_1$ to $D_2$ for the other three methods is presented for the $D_1$ noncertainty PSUs only.

The independent sample method is used as a benchmark for evaluating the variances of the other methods. For $D_2$, PSUs are stratified for a one-PSU-per-stratum design of size $m$ USUs, and PSUs are selected to be in the $D_2$ sample independently of the $D_1$ sample PSUs. This stratification procedure attempts to optimize an objective function that includes estimates of both unemployment and civilian labor force, and in the rest of this section we will refer to a design that uses it as optimal. This is the method that would be used if the $D_2$ sample was selected for a straightforward redesign of a survey which did not require an intervening $D_1$ sample.

The independent supplement method starts with the $D_1$ sample. It then restratifies all the PSUs for a one-PSU-per-stratum design of size $m$-$m^*$ USUs, and selects PSUs from the new strata independently of the $D_1$ sample PSUs. The resulting sample is added to that of $D_1$ to complete $D_2$.

In general, two-dimensional controlled selection (Goodman and Kish 1950; Causey, Cox, and Ernst 1985) is a technique for selecting sampling units simultaneously satisfying two criteria of stratification, while preserving the selection probabilities for both stratifications. In this particular application, the $D_1$ design consists of $L$ strata and the $D_2$ design consists of (say) $L'$ strata, with both designs one-PSU-per-stratum. The $D_2$ stratification for controlled selection is the same as for the independent sample method. To satisfy the $D_2$ stratification criterion, the technique selects $L'$ PSUs, with one PSU from each of the $D_2$ strata. To satisfy the $D_1$ stratification criterion, the method ensures that $L$ of these $L'$ PSUs will be from separate $D_1$ strata and will constitute the $D_1$ sample PSUs, while the

remaining $L'$-$L$ PSUs will fall in a leftover "stratum" of PSUs that will be in the $D_2$ sample but not in the $D_1$ sample. See Ernst (1990) for further details.

### 6.1.   Properties of the methods

The independent sample method uses an optimal $D_2$ stratification, but is the only one of the four methods that does not assure that $D_1$ sample PSUs are retained in $D_2$. Controlled selection is the only method that selects PSUs from optimal stratifications for both designs while retaining all $D_1$ sample PSUs in $D_2$. However, unlike the other three methods, controlled selection requires that the $D_1$ and $D_2$ sample PSUs be selected simultaneously, and consequently cannot be used for an expansion planned after $D_1$ is in place. (This technique also has the disadvantage of introducing a between-strata component of variance in both the $D_1$ and $D_2$ designs.) If we want to have $D_1 \subset D_2$, but are unable to select $D_2$ at the same time as $D_1$, then multiple WLs and the independent supplement are the only two methods of the four compared that are operationally feasible. However, they will probably have larger variances than the other two methods because they do not use an optimal $D_2$ stratification. The decision as to which of the four methods should be used in any situation depends on the importance attached to the different properties and the variances of the methods relative to each other.

Tables 4 through 7 present the ratios of variances for controlled selection, independent supplement, and multiple WLs methods (both estimators, $\hat{Y} = $ MW and $\hat{Y}' = $ AL), to the independent sample method. Tables 5 and 7 show between-PSUs variances and Tables 4 and 6 show total variances. These total variances include the within-PSUs component from both certainty and noncertainty PSUs. For all four methods, 1980 census data were used to obtain the stratifications, since 1990 census data were unavailable at the time these computations were done. The variables used were number of unemployed persons and number of persons in the civilian labor force (CLF). The ratios were computed for 31 states. Averages of these ratios over these 31 states were also computed. The remaining states were omitted for various reasons, as described in Ernst (1990).

For each of the four methods (in this and the next three paragraphs we refer to the multiple WL estimator $\hat{Y}$ only), the within-PSUs variances were obtained by computing the simple random sampling with replacement variance for sample size $m$ and multiplying by a design factor to account for the fact that clustered, systematic sampling was actually used within each PSU. For the multiple WLs method, this approach to computing the within-PSUs variances is at least partially justified by the results at the end of Section 4. The within-PSUs component of each variance is thus computed to be the same for all four methods and the differences among the methods are due solely to differences in the between-PSUs component for all methods, and also the between-strata component for the controlled selection and multiple WLs methods, which are the only methods among the four methods to have such a component.

Tables 4 and 5 compare the methods at the time of stratification using 1980 data. Tables 6 and 7 use 1970 data to simulate a ten-year lag between stratification and the collection of the survey data, which would be roughly the average lag time for the two-phase CPS. Because the between-PSUs components are such a small proportion of the total variances in these designs, the ratios in Tables 4 and 6 are all extremely close

*Table 4. 1980 ratios of total variances for other options to the independent sample*

| State | Unemployed | | | | Civilian labor force | | | |
|---|---|---|---|---|---|---|---|---|
| | CS | IS | MW | AL | CS | IS | MW | AL |
| Alabama | 1.000 | 1.009 | 1.000 | 1.021 | 0.999 | 1.028 | 1.032 | 1.043 |
| Arizona | 0.999 | 1.017 | 0.996 | 1.003 | 0.998 | 1.015 | 1.035 | 1.030 |
| Arkansas | 0.998 | 1.014 | 1.005 | 1.023 | 1.003 | 1.008 | 1.037 | 1.037 |
| Colorado | 0.998 | 1.051 | 0.996 | 1.008 | 0.998 | 1.260 | 1.057 | 1.051 |
| Georgia | 1.000 | 1.004 | 0.999 | 0.999 | 0.998 | 1.012 | 1.025 | 1.025 |
| Idaho | 1.000 | 1.165 | 0.994 | 1.003 | 1.000 | 1.096 | 1.077 | 1.044 |
| Indiana | 1.000 | 1.058 | 1.000 | 1.055 | 0.995 | 1.068 | 1.058 | 1.086 |
| Iowa | 1.000 | 1.018 | 0.997 | 1.004 | 1.000 | 1.021 | 1.012 | 1.013 |
| Kansas | 0.999 | 1.008 | 1.001 | 1.018 | 1.001 | 1.055 | 1.046 | 1.047 |
| Kentucky | 0.999 | 1.022 | 0.999 | 1.021 | 0.998 | 1.039 | 1.050 | 1.030 |
| Louisiana | 1.001 | 1.016 | 0.997 | 1.012 | 1.000 | 1.016 | 1.051 | 1.065 |
| Maryland | 1.000 | 1.013 | 0.999 | 0.999 | 1.000 | 1.110 | 1.052 | 1.052 |
| Minnesota | 1.004 | 1.008 | 0.984 | 0.991 | 1.000 | 1.061 | 1.025 | 1.005 |
| Mississippi | 1.001 | 1.014 | 1.007 | 1.014 | 1.001 | 1.008 | 1.044 | 1.047 |
| Missouri | 1.000 | 1.017 | 1.000 | 1.021 | 1.001 | 1.092 | 1.046 | 1.052 |
| Montana | 1.005 | 1.114 | 0.990 | 0.995 | 0.993 | 1.033 | 1.031 | 1.030 |
| Nebraska | 0.998 | 1.019 | 0.999 | 1.008 | 1.000 | 1.023 | 1.057 | 1.060 |
| Nevada | 1.000 | 1.021 | 0.998 | 1.002 | 1.002 | 1.139 | 1.093 | 1.070 |
| New Mexico | 1.001 | 1.015 | 0.993 | 1.011 | 1.005 | 1.106 | 1.107 | 1.105 |
| North Dakota | 0.999 | 1.060 | 0.996 | 0.996 | 1.000 | 1.135 | 1.046 | 1.046 |
| Oklahoma | 0.998 | 1.013 | 0.997 | 0.997 | 0.999 | 1.077 | 1.045 | 1.045 |
| Oregon | 1.000 | 1.025 | 0.998 | 1.012 | 1.001 | 1.034 | 1.053 | 1.058 |
| South Carolina | 1.000 | 1.010 | 1.008 | 1.009 | 1.000 | 1.014 | 1.043 | 1.011 |
| South Dakota | 0.995 | 1.061 | 0.987 | 0.987 | 0.998 | 1.065 | 1.048 | 1.048 |
| Tennessee | 0.999 | 1.013 | 0.999 | 1.007 | 1.000 | 1.031 | 1.036 | 1.036 |
| Utah | 1.000 | 1.012 | 1.000 | 1.009 | 1.000 | 1.048 | 1.068 | 1.057 |
| Virginia | 1.000 | 1.014 | 1.004 | 1.004 | 1.003 | 1.080 | 1.069 | 1.069 |
| Washington | 0.997 | 1.018 | 0.989 | 1.002 | 1.000 | 1.059 | 1.049 | 1.049 |
| West Virginia | 1.000 | 1.001 | 1.003 | 1.021 | 0.999 | 0.997 | 1.020 | 1.025 |
| Wisconsin | 1.000 | 1.017 | 1.007 | 1.016 | 1.000 | 1.042 | 1.058 | 1.048 |
| Wyoming | 1.000 | 1.093 | 1.000 | 1.018 | 0.999 | 1.135 | 1.089 | 1.037 |
| Mean | 1.000 | 1.030 | 0.998 | 1.009 | 1.000 | 1.062 | 1.050 | 1.046 |

CS = Controlled Selection
IS = Independent Supplement
MW = Multiple Workloads ($\hat{Y}$)
AL = Alternative MW Estimator ($\hat{Y}'$)

to 1. Thus, in the CPS case, there is not much difference between the methods. In order to compare the methods more generally, we concentrate on the between-PSUs variances.

For 1980 the between-PSUs variances of unemployed for $\hat{Y}$ are less than the benchmark for 21 of 31 states. The influence of the large ratio for South Carolina is the primary reason for the average ratio of the 31 states being greater than 1. The CLF ratios are $\geq$ 2.25 for all but one state. (The result for the unemployed is surprising. Since the 1980 $D_2$ stratification used for the benchmark is optimal while $\hat{Y}$ is based upon a stratification optimal for the smaller $D_1$, it might be expected that the variances of the benchmark would be smaller. The fact that they are not might be due to the attempt to choose the $D_2$ stratification to optimize a function of

Table 5.  1980 ratios of between PSU variances for other options to the independent sample

| State | Unemployed | | | | Civilian labor force | | | |
|---|---|---|---|---|---|---|---|---|
| | CS | IS | MW | AL | CS | IS | MW | AL |
| Alabama | 1.01 | 2.40 | 0.99 | 0.92 | 0.91 | 3.57 | 4.00 | 2.96 |
| Arizona | 0.89 | 4.53 | 0.12 | 0.08 | 0.76 | 2.54 | 4.51 | 3.26 |
| Arkansas | 0.70 | 3.14 | 1.80 | 1.33 | 1.12 | 1.35 | 2.61 | 1.69 |
| Colorado | 0.85 | 5.44 | 0.64 | 0.57 | 0.96 | 6.68 | 2.25 | 1.83 |
| Georgia | 1.02 | 1.96 | 0.83 | 0.83 | 0.87 | 1.77 | 2.65 | 2.65 |
| Idaho | 0.99 | 12.38 | 0.58 | 0.33 | 1.05 | 14.87 | 12.15 | 5.42 |
| Indiana | 1.05 | 9.48 | 0.94 | 0.89 | 0.65 | 6.14 | 5.39 | 3.36 |
| Iowa | 1.01 | 3.06 | 0.69 | 0.63 | 0.98 | 4.54 | 3.10 | 1.97 |
| Kansas | 0.89 | 1.55 | 1.10 | 1.48 | 0.92 | 1.71 | 0.85 | 0.64 |
| Kentucky | 0.91 | 4.48 | 0.81 | 0.87 | 0.90 | 3.50 | 4.19 | 1.54 |
| Louisiana | 1.12 | 3.37 | 0.63 | 0.57 | 1.02 | 2.53 | 5.86 | 5.72 |
| Maryland | 1.00 | 5.13 | 0.59 | 0.59 | 1.00 | 43.15 | 20.81 | 20.81 |
| Minnesota | 1.14 | 1.30 | 0.39 | 0.38 | 1.02 | 4.21 | 2.33 | 0.84 |
| Mississippi | 1.26 | 4.59 | 2.76 | 2.52 | 1.08 | 1.46 | 3.44 | 3.15 |
| Missouri | 1.06 | 3.98 | 1.06 | 1.12 | 1.05 | 4.74 | 2.85 | 2.32 |
| Montana | 1.18 | 5.37 | 0.60 | 0.57 | 0.86 | 1.66 | 1.61 | 1.47 |
| Nebraska | 0.79 | 3.38 | 0.88 | 0.84 | 0.97 | 3.05 | 6.13 | 5.60 |
| Nevada | 0.98 | 5.91 | 0.42 | 0.48 | 1.07 | 4.87 | 3.58 | 2.84 |
| New Mexico | 1.10 | 2.12 | 0.50 | 0.45 | 1.19 | 4.73 | 4.77 | 4.05 |
| North Dakota | 0.95 | 3.86 | 0.79 | 0.79 | 0.99 | 8.55 | 3.60 | 3.60 |
| Oklahoma | 0.85 | 2.27 | 0.70 | 0.70 | 0.95 | 5.55 | 3.66 | 3.66 |
| Oregon | 1.05 | 4.00 | 0.71 | 0.76 | 1.17 | 5.31 | 7.78 | 6.68 |
| South Carolina | 1.43 | 19.29 | 15.38 | 14.53 | 1.00 | 4.89 | 12.74 | 3.56 |
| South Dakota | 0.87 | 2.72 | 0.63 | 0.63 | 0.92 | 3.53 | 2.84 | 2.84 |
| Tennessee | 0.89 | 3.81 | 0.83 | 0.86 | 1.01 | 3.06 | 3.37 | 2.91 |
| Utah | 0.93 | 3.56 | 1.04 | 1.07 | 1.00 | 4.52 | 5.96 | 4.55 |
| Virginia | 1.06 | 4.57 | 2.02 | 2.02 | 1.25 | 8.91 | 7.85 | 7.85 |
| Washington | 0.88 | 1.84 | 0.47 | 0.42 | 1.04 | 6.85 | 5.86 | 4.38 |
| West Virginia | 0.59 | 2.14 | 4.11 | 3.91 | 0.87 | 0.41 | 5.02 | 2.46 |
| Wisconsin | 1.07 | 4.68 | 2.59 | 2.61 | 1.00 | 3.77 | 4.77 | 3.55 |
| Wyoming | 1.04 | 26.75 | 1.04 | 0.80 | 0.82 | 31.75 | 21.40 | 5.20 |
| Mean | 0.99 | 5.26 | 1.50 | 1.43 | 0.98 | 6.64 | 5.81 | 4.04 |

CS = Controlled Selection
IS = Independent Supplement
MW = Multiple Workloads $(\hat{Y})$
AL = Alternative MW Estimator $(\hat{Y}')$

both unemployed and CLF.) For the ten-year lag of 1970 the situation is reversed, with all but four unemployed ratios greater than 1.0 and the majority of CLF ratios less than 1.0. In all cases the mean between-PSU variances are smaller for the benchmark than for $\hat{Y}$.

At the time of PSU size measurement in 1980, the mean between-PSU variances are less for multiple WLs than for the independent supplement for both unemployed and CLF, but in the CLF case each method has smaller variance for about half the states. The between-PSU variances for the 1970 data are smaller for $\hat{Y}$ than for IS in 22 states for unemployed and 29 states for CLF. The mean is less for $\hat{Y}$ in each case, so the multiple WL method is clearly preferred to the independent supplement for this CPS example.

*Table 6.   1970 ratios of total variances for other options to the independent sample*

| State | Unemployed | | | | Civilian labor force | | | |
|-------|------|------|------|------|------|------|------|------|
|       | CS | IS | MW | AL | CS | IS | MW | AL |
| Alabama | 1.001 | 1.019 | 1.016 | 1.051 | 1.000 | 1.028 | 1.064 | 1.025 |
| Arizona | 1.000 | 1.003 | 1.014 | 1.043 | 0.998 | 1.069 | 1.016 | 1.018 |
| Arkansas | 1.001 | 1.032 | 0.991 | 1.020 | 0.999 | 1.036 | 0.957 | 0.976 |
| Colorado | 1.000 | 1.031 | 1.024 | 1.054 | 1.000 | 1.212 | 1.065 | 1.059 |
| Georgia | 1.000 | 1.010 | 1.002 | 1.002 | 1.000 | 1.024 | 0.990 | 0.990 |
| Idaho | 1.002 | 1.259 | 1.041 | 1.048 | 1.003 | 1.155 | 1.045 | 1.002 |
| Indiana | 0.999 | 1.028 | 1.033 | 1.127 | 1.018 | 1.112 | 1.003 | 1.090 |
| Iowa | 0.998 | 1.030 | 1.001 | 1.011 | 0.997 | 1.055 | 1.001 | 1.005 |
| Kansas | 1.000 | 1.026 | 1.009 | 1.034 | 0.996 | 1.035 | 0.993 | 1.009 |
| Kentucky | 1.001 | 1.020 | 1.028 | 1.062 | 0.997 | 1.050 | 1.040 | 1.044 |
| Louisiana | 1.000 | 1.007 | 1.012 | 1.041 | 0.999 | 1.054 | 1.037 | 1.045 |
| Maryland | 1.000 | 1.020 | 1.000 | 1.001 | 1.000 | 1.196 | 0.952 | 0.952 |
| Minnesota | 1.002 | 1.011 | 1.014 | 1.029 | 0.996 | 1.057 | 1.013 | 1.011 |
| Mississippi | 0.999 | 1.047 | 1.003 | 1.012 | 1.001 | 1.132 | 0.980 | 0.981 |
| Missouri | 1.000 | 1.026 | 1.010 | 1.059 | 1.005 | 1.120 | 0.979 | 1.018 |
| Montana | 0.997 | 1.032 | 1.043 | 1.036 | 1.009 | 1.261 | 0.979 | 0.990 |
| Nebraska | 1.000 | 1.020 | 1.006 | 1.020 | 1.000 | 1.078 | 0.983 | 0.998 |
| Nevada | 1.000 | 1.035 | 0.998 | 1.040 | 0.992 | 1.855 | 0.960 | 0.994 |
| New Mexico | 0.999 | 1.013 | 1.040 | 1.056 | 1.011 | 1.084 | 1.053 | 1.049 |
| North Dakota | 0.999 | 1.071 | 1.011 | 1.011 | 0.986 | 1.106 | 0.978 | 0.978 |
| Oklahoma | 1.000 | 1.012 | 1.009 | 1.009 | 0.993 | 1.117 | 0.995 | 0.995 |
| Oregon | 0.999 | 1.015 | 1.021 | 1.054 | 0.999 | 1.216 | 0.981 | 1.005 |
| South Carolina | 1.003 | 1.054 | 0.993 | 1.016 | 1.002 | 1.135 | 1.009 | 1.023 |
| South Dakota | 0.996 | 1.054 | 1.020 | 1.020 | 0.997 | 1.092 | 0.977 | 0.977 |
| Tennessee | 1.001 | 1.047 | 1.005 | 1.026 | 1.000 | 1.060 | 0.993 | 1.005 |
| Utah | 1.000 | 1.019 | 1.002 | 1.039 | 0.995 | 1.128 | 0.980 | 0.987 |
| Virginia | 0.999 | 1.031 | 0.998 | 0.998 | 1.018 | 1.112 | 0.970 | 0.970 |
| Washington | 0.999 | 1.034 | 1.012 | 1.051 | 1.008 | 1.148 | 1.015 | 1.030 |
| West Virginia | 1.000 | 1.002 | 1.012 | 1.041 | 1.001 | 0.994 | 1.002 | 1.032 |
| Wisconsin | 0.999 | 1.028 | 1.015 | 1.025 | 0.999 | 1.086 | 0.989 | 0.988 |
| Wyoming | 0.999 | 1.014 | 1.012 | 1.053 | 1.000 | 1.212 | 1.067 | 1.050 |
| Mean | 1.000 | 1.034 | 1.013 | 1.035 | 1.001 | 1.130 | 1.002 | 1.010 |

CS = Controlled Selection
IS = Independent Supplement
MW = Multiple Workloads ($\hat{Y}$)
AL = Alternative MW Estimator ($\hat{Y}'$)

Finally, we compare the variances of the two multiple WL estimators $\hat{Y}$ and $\hat{Y}'$. In 1980 $\hat{Y}'$ has smaller between-PSUs variances in 23 states for unemployed and all states for CLF, and the means are also less for $\hat{Y}'$. The 1970 results only have two unemployed cases for which $\hat{Y}'$ has slightly larger between-PSU variances. (As the proportion of between-PSUs variance due to between-strata variance decreases for a state (not shown), the ratio of $\hat{Y}$ to $\hat{Y}'$ between-PSUs variances increases, as suggested by the absence of a between-strata component for $\hat{Y}'$.) On the other hand, total variances are smaller in most of these cases for $\hat{Y}$. This is because the between-PSUs component makes up a small fraction of the total

*Table 7.   1970 ratios of between PSU variances for other options to the independent sample*

| State | Unemployed | | | | Civilian labor force | | | |
|---|---|---|---|---|---|---|---|---|
| | CS | IS | MW | AL | CS | IS | MW | AL |
| Alabama | 1.12 | 3.36 | 3.00 | 2.62 | 1.03 | 2.87 | 5.33 | 1.26 |
| Arizona | 1.07 | 4.21 | 18.63 | 16.44 | 0.89 | 4.68 | 1.84 | 0.34 |
| Arkansas | 1.06 | 2.67 | 0.53 | 0.32 | 0.99 | 1.57 | 0.34 | 0.14 |
| Colorado | 1.01 | 4.32 | 3.59 | 2.57 | 1.00 | 5.31 | 2.33 | 1.42 |
| Georgia | 1.04 | 2.42 | 1.25 | 1.25 | 1.02 | 1.91 | 0.64 | 0.64 |
| Idaho | 1.10 | 13.47 | 2.95 | 1.67 | 1.06 | 4.64 | 2.05 | 0.26 |
| Indiana | 0.89 | 4.00 | 4.52 | 2.30 | 1.47 | 3.90 | 1.08 | 0.44 |
| Iowa | 0.78 | 3.63 | 1.10 | 0.97 | 0.83 | 4.29 | 1.09 | 0.66 |
| Kansas | 0.95 | 3.49 | 1.90 | 1.48 | 0.92 | 1.71 | 0.85 | 0.64 |
| Kentucky | 1.08 | 3.35 | 4.22 | 4.00 | 0.84 | 3.67 | 3.13 | 1.43 |
| Louisiana | 0.95 | 2.65 | 3.89 | 3.30 | 0.97 | 2.63 | 2.14 | 1.46 |
| Maryland | 1.00 | 4.62 | 1.09 | 1.09 | 1.00 | 4.92 | 0.04 | 0.04 |
| Minnesota | 1.12 | 1.51 | 1.69 | 1.44 | 0.87 | 2.75 | 1.39 | 0.73 |
| Mississippi | 0.93 | 4.84 | 1.28 | 1.08 | 1.03 | 3.70 | 0.59 | 0.39 |
| Missouri | 0.95 | 3.99 | 2.17 | 2.16 | 1.09 | 3.25 | 0.61 | 0.45 |
| Montana | 0.88 | 2.16 | 2.54 | 1.70 | 1.19 | 6.24 | 0.58 | 0.49 |
| Nebraska | 0.99 | 2.83 | 1.52 | 1.36 | 1.00 | 2.78 | 0.62 | 0.59 |
| Nevada | 1.02 | 5.81 | 0.78 | 0.82 | 0.86 | 15.02 | 0.35 | 0.25 |
| New Mexico | 0.91 | 3.06 | 7.26 | 3.27 | 1.41 | 4.13 | 2.97 | 1.32 |
| North Dakota | 0.91 | 6.34 | 1.85 | 1.85 | 0.72 | 3.09 | 0.57 | 0.57 |
| Oklahoma | 1.02 | 2.72 | 2.27 | 2.27 | 0.83 | 3.86 | 0.88 | 0.88 |
| Oregon | 0.89 | 2.25 | 2.72 | 2.15 | 0.98 | 5.40 | 0.61 | 0.32 |
| South Carolina | 1.17 | 3.64 | 0.68 | 0.41 | 1.10 | 9.67 | 1.56 | 0.67 |
| South Dakota | 0.90 | 2.53 | 1.56 | 1.56 | 0.95 | 2.50 | 0.63 | 0.63 |
| Tennessee | 1.10 | 8.16 | 1.77 | 1.81 | 1.00 | 3.01 | 0.76 | 0.50 |
| Utah | 0.97 | 2.66 | 1.18 | 1.18 | 0.94 | 2.67 | 0.73 | 0.39 |
| Virginia | 0.95 | 2.87 | 0.87 | 0.87 | 1.33 | 3.07 | 0.44 | 0.44 |
| Washington | 0.94 | 3.29 | 1.78 | 1.68 | 1.25 | 5.86 | 1.48 | 0.66 |
| West Virginia | 0.93 | 1.60 | 4.23 | 2.05 | 1.12 | 0.37 | 1.22 | 0.52 |
| Wisconsin | 0.96 | 2.17 | 1.62 | 1.38 | 0.98 | 2.94 | 0.74 | 0.37 |
| Wyoming | 0.67 | 5.69 | 5.04 | 5.03 | 1.03 | 35.91 | 12.08 | 2.60 |
| Mean | 0.98 | 3.88 | 2.89 | 2.33 | 1.02 | 5.11 | 1.60 | 0.69 |

CS = Controlled Selection
IS = Independent Supplement
MW = Multiple Workloads ($\hat{Y}$)
AL = Alternative MW Estimator ($\hat{Y}'$)

variance ($< .07$) for every one of these estimates and $f(R) > 1$ except when $R$ is an integer. In cases when the between-PSUs component is a larger proportion of the total variance, we would expect to see the comparisons of $\hat{Y}$ and $\hat{Y}'$ be more mixed, moving toward favoring $\hat{Y}'$ more often as the proportion increases.

## 7.   References

Causey, B.D., Cox L.H., and Ernst, L.R. (1985). Applications of Transportation Theory to Statistical Problems. Journal of the American Statistical Association, 80, 903–909.

Chandhok, P., Weinstein, R., and Gunlicks, C. (1990). Augmenting a Sample to Satisfy Subpopulation Reliability Requirements. Proceedings of the Section on Survey Research Methods, American Statistical Association, 696–701.

Cochran, W.G. (1977). Sampling Techniques. New York: John Wiley and Sons.

Ernst, L.R. (1990). Simultaneous Selection of Primary Sampling Units for Two Designs. Proceedings of the Section on Survey Research Methods, American Statistical Association, 688–693.

Goodman, R. and Kish, L. (1950). Controlled Selection — A Technique in Probability Sampling. Journal of the American Statistical Association, 45, 350–372.

Sampford, M.R. (1967). On Sampling Without Replacement With Unequal Probabilities of Selection. Biometrika, 54, 499–513.