

## REPLICATION VARIANCE ESTIMATION FOR THE TWO-PHASE REGRESSION ESTIMATOR

Wayne A. Fuller<sup>1</sup>

### ABSTRACT

A replication variance estimation method for the regression estimator in two phase samples is described. The procedure reproduces the estimated variance of the first phase sample for those attributes of the first phase sample used as controls for the second phase sample. Only the second phase sample is required by the analyst for variance estimation. The procedure is applied to the National Resources Inventory (NRI), a study of land use based on a large area sample of the United States.

**Key Words:** Double sampling, Reweighted expansion estimator, Multi-phase sampling

### 1. INTRODUCTION

In two-phase sampling, also known as double sampling, observations are made on a vector of auxiliary variables in a large sample, called the first phase sample. Then observations are made on the variables of interest using a smaller sample, called the second phase sample. The smaller sample is typically a subsample of the original sample, and often the vector of auxiliary variables is a vector of indicator variables defining subgroups (strata) of the original sample used in selecting the second phase sample.

Rao (1973) and Cochran (1977) give formulas for variance estimation when the first phase is a simple random sample and the second phase is a stratified random sample. Särndal, Swensson and Wretman (1992, ch. 9) discuss two phase sampling. Kott (1990) derived a formula for variance estimation when the first phase is a stratified random sample and the second phase is a stratified random sample of the first phase sample based on strata defined by first-phase information. Rao and Shao (1992) proposed a jackknife variance estimation method in the context of hot deck imputation where the second phase strata correspond to imputation cells. Binder (1996) illustrated a “cookbook approach” to variance estimation for the two phase ratio estimator. Binder et al. (1997) derived formulas for variance estimation for various estimators for two-phase re-stratified sampling. Rao and Sitter (1995) and Sitter (1997) consider variance estimation for two-phase samples. Fuller (1998) proposed a replicate variance estimation method for the two-phase regression estimator that is particularly convenient for multipurpose surveys. Kim and Sitter (2003) extended that procedure.

Kott and Stukel (1997) pointed out that, for a stratified second phase sample with the reweighted expansion estimator, the replication method proposed by Rao and Shao (1992) produces consistent variance estimates. Kim, Navarro and Fuller (2004) give a replication method for estimating the variance of the direct expansion estimator for a second phase stratified sample. We give the corresponding replication variance estimator for the two-phase regression estimator. In the last section, an application of replication variance estimation in the Natural Resources Inventory is described.

### 2. ASYMPTOTIC PROPERTIES

Let the finite population be of size  $N$  and indexed from 1 to  $N$ . Let the parameter of interest be the population total  $Y = \sum_{i=1}^N y_i$ , where  $y_i$  is the study variable and  $N$  is assumed to be known. Let  $E\{\hat{\theta}|F\}$  and  $V\{\hat{\theta}|F\}$  denote the expectation and variance of a statistic  $\hat{\theta}$  for a particular finite population  $F$ . These quantities are sometimes called the design expectation and design variance. Suppose we have a first phase sample of size  $n_1$ . We observe  $\mathbf{x}_i$  for all  $i \in A_1$  where  $\mathbf{x}_i$  is a vector of auxiliary variables and  $A_1$  is the set of indices of the  $n_1$  first phase elements. If the second phase sample is stratified, indicators for strata are included in  $\mathbf{x}_i$  where  $x_{ig}$  takes the value one if unit  $i$  belongs to the  $g$ -th group and is zero otherwise. A second phase sample of size  $n_2$  is selected from the  $n_1$  first phase elements.

We define the two-phase regression estimator of the mean of  $y$  by

$$\bar{y}_{reg} = \bar{y}_{\pi_2} + (\bar{\mathbf{x}}_{\pi_1} - \bar{\mathbf{x}}_{\pi_2})\hat{\boldsymbol{\beta}}_2, \quad (2.1)$$

where

$$\begin{aligned} \hat{\boldsymbol{\beta}}_2 &= \left( \sum_{i \in A_2} \pi_{2i}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}_{\pi_2})' (\mathbf{x}_i - \bar{\mathbf{x}}_{\pi_2}) \right)^{-1} \\ &\times \sum_{i \in A_2} \pi_{2i}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}_{\pi_2})' (y_i - \bar{y}_{\pi_2}), \\ (\bar{\mathbf{x}}_{\pi_2}, \bar{y}_{\pi_2}) &= \left( \sum_{i \in A_2} \pi_{2i}^{-1} \right) \sum_{i \in A_2} \pi_{2i}^{-1} (\mathbf{x}_i, y_i), \\ \bar{\mathbf{x}}_{\pi_1} &= \left( \sum_{i \in A_1} \pi_{1i}^{-1} \right)^{-1} \sum_{i \in A_1} \pi_{1i}^{-1} \mathbf{x}_i, \end{aligned}$$

$\pi_{1i}$  is the first phase selection probability for element  $i$ ,  $\pi_{2i} = \pi_{1i} \pi_{2i|i}$  is the second phase selection probability for element  $i$ , and  $\pi_{2i|i}$  is the conditional probability of selecting element  $i$  for the second phase sample given that  $i$  is in the first phase sample. To derive the asymptotic properties of the two-phase regression estimator we assume a sequence of samples and finite populations such as that described in Fuller (1975).

**Theorem 1.** Let  $\{F_N\}$  be a sequence of finite populations of size  $N$ . Let the  $G_N$  dimensional vectors  $(1, \mathbf{x}_{N,i}, y_{Ni}) = \mathbf{z}_{Ni}$ ,  $i = 1, 2, \dots, N$ , be a random sample from an infinite population with finite  $2 + \tau$ ,  $\tau > 0$ , moments. Assume  $E\{\mathbf{z}'_{Ni} \mathbf{z}_{Ni}\}$  is positive definite for all  $N$ . Let first phase samples of size  $n_{1,N}$ ,  $n_{1,N} \geq n_{1,N-1}$ , be selected by designs such that

$$V\{(\bar{\mathbf{z}}_{\pi_{1,N}} - \bar{\mathbf{z}}_N) | F_N\} = O_p(n_{1,N}^{-1}), \quad (2.2)$$

where

$$\bar{\mathbf{z}}_{\pi_{1,N}} = \left( \sum_{i \in A_{N1}} \pi_{1,Ni}^{-1} \right)^{-1} \sum_{i \in A_{N1}} \pi_{1,Ni}^{-1} \mathbf{z}_{Ni},$$

$A_{N1}$  is the set of indices in the first phase sample selected from population  $N$ , and  $\pi_{1,Ni}$  is the probability of selecting element  $i$  for the first phase sample. Assume

$$E\{|\bar{\mathbf{x}}_{\pi_{2,N}} - \bar{\mathbf{x}}_N|^2\} = O(n_N^{-\delta}), \quad (2.3)$$

$$E\{|\bar{\mathbf{x}}_{\pi_{1,N}} - \bar{\mathbf{x}}_N|^2\} = O(n_N^{-\delta}), \quad (2.4)$$

$$E\{|\hat{\boldsymbol{\beta}}_{2,N} - \boldsymbol{\beta}_N|^2\} = O(n_N^{-\gamma}), \quad (2.5)$$

for some  $\delta > 0$ ,  $\gamma > 0$  and  $\delta + \gamma = \lambda > 0.5$ , where  $\bar{\mathbf{x}}_{\pi_1}$  and  $\hat{\boldsymbol{\beta}}_2$  are defined in (2.1), and

$$\begin{aligned} \boldsymbol{\beta}_N &= \left[ \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}_N)' (\mathbf{x}_i - \bar{\mathbf{x}}_N) \right]^{-1} \\ &\times \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}_N)' (y_i - \bar{y}_N). \end{aligned}$$

Let  $\bar{y}_{reg,N}$  be defined by (2.1). Then

$$\bar{y}_{reg,N} - \bar{y}_N = \bar{e}_{\pi_{2,N}} + (\bar{\mathbf{x}}_{\pi_{1,N}} - \bar{\mathbf{x}}_N)\boldsymbol{\beta}_N + O_p(n^{-0.5\lambda}) \quad (2.6)$$

where

$$\bar{e}_{\pi_{2,N}} = \left( \sum_{i \in A_2} \pi_{2,Ni}^{-1} \right)^{-1} \sum_{i \in A_2} \pi_{2,Ni}^{-1} e_{Ni},$$

and  $e_{Ni} = y_{Ni} - \bar{y}_N - (\mathbf{x}_{Ni} - \bar{\mathbf{x}}_N)\boldsymbol{\beta}_N$ .

**Proof.** To simplify the notation we omit the subscript  $N$ , except to identify finite population parameters. The estimator can be written.

$$\begin{aligned} \bar{y}_{reg} - \bar{y}_N &= \bar{y}_{\pi_2} + (\bar{\mathbf{x}}_{\pi_1} - \bar{\mathbf{x}}_{\pi_2})\hat{\boldsymbol{\beta}}_2 - \bar{y}_N \\ &= \bar{e}_{\pi_2} - (\bar{\mathbf{x}}_{\pi_2} - \bar{\mathbf{x}}_{\pi_1})(\hat{\boldsymbol{\beta}}_2 - \boldsymbol{\beta}_N) \\ &\quad + (\bar{\mathbf{x}}_{\pi_1} - \bar{\mathbf{x}}_N)\boldsymbol{\beta}_N \\ &= \bar{e}_{\pi_2} + (\bar{\mathbf{x}}_{\pi_1} - \bar{\mathbf{x}}_N)\boldsymbol{\beta}_N + O_p(n^{-0.5\lambda}). \end{aligned} \quad (2.7)$$

establishing result (2.6). ■

The result permits the size of the vector  $\mathbf{x}$  to increase with sample size. Thus a design that has a relatively large number of second phase strata is covered by the theorem. For a second phase stratified sample in which  $\mathbf{x}$  contains stratum indicators, assumptions (2.3), (2.4), and (2.5) require the number of strata to increase at a slower rate than the sample size. The variance of the approximating random variable in (2.6) is

$$\begin{aligned} V\{d_{reg,N}\} &= V\{\bar{e}_{\pi_{2,N}}\} + V\{(\bar{\mathbf{x}}_{\pi_{1,N}} - \bar{\mathbf{x}}_N)\boldsymbol{\beta}_N\} \\ &\quad + 2C\{\bar{e}_{\pi_{1,N}}, (\bar{\mathbf{x}}_{\pi_{1,N}} - \bar{\mathbf{x}}_N)\boldsymbol{\beta}_N\}, \end{aligned} \quad (2.8)$$

where  $d_{reg,N} = \bar{e}_{\pi_{2,N}} + (\bar{\mathbf{x}}_{\pi_{1,N}} - \bar{\mathbf{x}}_N) \boldsymbol{\beta}_N$ . Note that  $V\{\bar{e}_{\pi_{2,N}}\}$  is the unconditional variance.

### 3. REPLICATE VARIANCE ESTIMATION

We consider variance estimation for the stratified two-phase regression estimator where replicates are created on the basis of a replication procedure for the first phase. Let a replication variance estimator of the first phase estimated total be

$$\hat{V}_1\{\hat{T}_{1x}\} = \sum_{k=1}^L c_k (\hat{T}_{1x}^{(k)} - \hat{T}_{1x})^2 \tag{3.1}$$

where  $\hat{T}_{1x}^{(k)}$  is the  $k$  th replicate of the estimated total,  $\hat{T}_{1x}$  is the estimated total based on the first phase,  $L$  is the total number of replicates, and  $c_k, k = 1, 2, \dots, L$ , are constants determined by the replication method.

The sample modification that creates replicates for the first phase is applied to the associated second-phase units to create the second phase replicate. For example, let the first phase sample be a simple random sample and let a single element be removed to form a jackknife replicate. Then the second phase replicate contains the original second phase elements less the removed element, provided the removed first phase element is in the second phase sample. The replicate estimator for the two phase sample is then

$$\hat{V}_{2p}\{\bar{y}_{reg}\} = \sum_{k=1}^L c_k (\bar{y}_{reg}^{(k)} - \bar{y}_{reg})^2, \tag{3.2}$$

where

$$\bar{y}_{reg}^{(k)} = \bar{y}_2^{(k)} + (\bar{\mathbf{x}}_1^{(k)} - \bar{\mathbf{x}}_2^{(k)}) \hat{\boldsymbol{\beta}}_2^{(k)}$$

and  $(\bar{y}_2^{(k)}, \bar{\mathbf{x}}_2^{(k)}, \hat{\boldsymbol{\beta}}_2^{(k)})$  are computed from the second phase replicate.

**Theorem 2.** Assume a sequence of finite populations in which each population is composed of  $G$  groups with proportion  $W_g$  in the  $g$  th group. Let the  $N$  th population be

$$F_N = \{(\mathbf{x}_{1N}, y_{1N}), (\mathbf{x}_{2N}, y_{2N}), \dots, (\mathbf{x}_{NN}, y_{NN})\},$$

where the first  $G - 1$  elements of  $\mathbf{x}_{iN}$  are indicator variables for membership in  $G - 1$  of the groups. Assume the finite population in the  $g$  th group is a sample from an infinite population with  $4 + \delta, \delta > 0$ , moments. Let two-phase samples be selected where the first-phase selection probabilities are  $\pi_{i,N}$ . Let the second phase sample be a stratified sample with

$G$  strata and  $n_{2g}$  elements selected in the  $g$  th stratum. Assume a set of fixed probabilities,  $\pi_{2i|i} = \kappa_{2i}$ , that are constant within a group, is used to select the second-phase stratified random sample. Assume

$$C_{\pi S} < Nn_{1N}^{-1}\pi_{i,N} < C_{\pi B} \tag{3.3}$$

for all  $N$ , where  $C_{\pi S}$  and  $C_{\pi B}$  are fixed positive constants. Assume

$$E\{(\hat{\mathbf{T}}_{1x}, \hat{T}_{1y}) | F_N\} = (\mathbf{T}_x, T_y) \tag{3.4}$$

where

$$(\hat{\mathbf{T}}_{1x}, \hat{T}_{1y}) = \sum_{i \in A_1} w_{i,N}(\mathbf{x}_{iN}, y_{iN})$$

and  $w_{i,N} = \pi_{i,N}^{-1}$ . Assume

$$V\{\hat{T}_{1y} | F_N\} \leq K_M V\{\hat{T}_{y,SRS} | F_N\}, \tag{3.5}$$

for a fixed  $K_M$ , for any variable with fourth moments, where  $V\{\hat{T}_{y,SRS} | F\}$  is the variance of the Horvitz-Thompson estimator of the total for a simple random sample of size  $n_{1N}$ . Assume that the variance of a first-phase linear estimator of the mean is a symmetric quadratic function, and that

$$n_N V\left\{N^{-1} \sum_{i \in A_{1N}} \pi_{i,N}^{-1} y_{iN} | F_N\right\} = \sum_{i=1}^N \sum_{j=1}^N \omega_{ij,N} y_{iN} y_{jN} \tag{3.6}$$

for coefficients  $\omega_{ij,N}$ , satisfying

$$\sum_{i=1}^N |\omega_{ij,N}| = O(N^{-1}). \tag{3.7}$$

Let  $\hat{V}_1\{\hat{\boldsymbol{\theta}}\}$  be a first-phase replicate variance estimator of  $\hat{\boldsymbol{\theta}} = \sum_{i \in A_1} w_{i,N} y_{iN}$  and assume

$$E\left\{\left[\left(V[\hat{\boldsymbol{\theta}} | F_N]\right)^{-1} \hat{V}_1\{\hat{\boldsymbol{\theta}}\} - 1\right]^2 | F_N\right\} = o(1) \tag{3.8}$$

for any variable with bounded fourth moments. Assume the replicates for the first-phase sample estimator of a total,  $\hat{T}_1$ , satisfy

$$E\left\{\left[c_{kN} (\hat{T}_1^{(k)} - \hat{T}_1)\right]^2 | F_N\right\} < K_\gamma L_N^{-2} [V\{\hat{T}_1 | F_N\}]^2 \tag{3.9}$$

uniformly in  $N$  for any variable with fourth moments, where  $K_\gamma$  is a fixed constant.

Then the variance estimator (3.2) satisfies

$$\begin{aligned} \hat{V}_{2p}\{\bar{y}_{reg}\} &= V\{\bar{y}_{reg} | F_N\} \\ &- N^{-2} \sum_{i=1}^N \pi_{2i|i}^{-1} (1 - \pi_{2i|i}) e_{iN}^2 + o_p(n_1^{-1}), \end{aligned}$$

where  $e_{iN} = y_{iN} - \bar{y}_N - (\mathbf{x}_{iN} - \bar{\mathbf{x}}_N) \boldsymbol{\beta}_N$  and

$$\boldsymbol{\beta}_N = \left[ \sum_{i=1}^N (\mathbf{x}_{iN} - \bar{\mathbf{x}}_N)' (\mathbf{x}_{iN} - \bar{\mathbf{x}}_N) \right]^{-1} \times \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}_N)' (y_i - \bar{y}_N). \tag{3.10}$$

**Proof.** To simplify the notation we omit the  $N$  subscript. We write the regression estimator as

$$\begin{aligned} \bar{y}_{reg} &= \bar{y}_{2\pi} - (\bar{\mathbf{x}}_{2\pi} - \bar{\mathbf{x}}_N) \hat{\boldsymbol{\beta}}_{2\pi} + (\bar{\mathbf{x}}_{1\pi} - \bar{\mathbf{x}}_N) \hat{\boldsymbol{\beta}}_{2\pi} \\ \bar{y}_{reg} - \bar{y}_{reg} &= \bar{y}_{2\pi}^{(k)} - \bar{y}_{2\pi} - (\bar{\mathbf{x}}_{2\pi}^{(k)} - \bar{\mathbf{x}}_{2\pi}) \hat{\boldsymbol{\beta}}_{2\pi} \\ &+ (\bar{\mathbf{x}}_{1\pi}^{(k)} - \bar{\mathbf{x}}_{1\pi}) \hat{\boldsymbol{\beta}}_{2\pi} - (\bar{\mathbf{x}}_{2\pi}^{(k)} - \bar{\mathbf{x}}_{1\pi}^{(k)}) (\hat{\boldsymbol{\beta}}_{2\pi}^{(k)} - \hat{\boldsymbol{\beta}}_{2\pi}) \\ &= \bar{e}_{2\pi}^{(k)} - \bar{e}_{2\pi} + (\bar{\mathbf{x}}_{1\pi}^{(k)} - \bar{\mathbf{x}}_{1\pi}) \boldsymbol{\beta}_N + O_p(L^{-1/2} n^{-1}), \end{aligned} \tag{3.11}$$

where we used (3.9) and (3.5).

Let

$$\begin{aligned} a_i &= 1 && \text{if element } i \text{ is selected for the second} \\ & && \text{phase sample} \\ &= 0 && \text{otherwise.} \end{aligned}$$

Under the assumption that the second-phase sampling rates,  $\pi_{2i|i} = \kappa_{2i}$  are fixed, we can conceptualize the sample selection process as composed of two steps. First an  $a$  is generated for every element in the population and then a first-phase sample is selected from the created population of  $(a_i, a_i y_i, \mathbf{x}_i)$  vectors.

Observe that

$$\begin{aligned} \bar{e}_{2\pi} &= \left( \sum_{i \in A_1} \pi_{1i}^{-1} \pi_{2i|i}^{-1} a_i \right)^{-1} \sum_{i \in A_1} \pi_{1i}^{-1} \pi_{2i|i}^{-1} a_i e_i \\ &= \frac{N^{-1} \sum_{i \in A_1} \pi_{1i}^{-1} \kappa_{2i}^{-1} a_i e_i}{N^{-1} E \left\{ \sum_{i \in A_1} \pi_{2i}^{-1} a_i \right\}} + O_p(n^{-1}) \\ &= N^{-1} \sum_{i \in A_1} \pi_{2i}^{-1} a_i e_i + O_p(n^{-1}) \\ &=: \bar{e}_{2,HT} + O_p(n^{-1}). \end{aligned} \tag{3.12}$$

Therefore we consider the variance of  $\bar{e}_{2,HT}$ , which we write in terms of conditional expectations,

$$\begin{aligned} V_1 \left\{ \sum_{i \in A_1} \pi_{2i}^{-1} a_i e_i \mid F \right\} &= E \left\{ V_1 \left[ \bar{e}_{2,HT} \mid (\mathbf{a}_N, F) \right] \mid F \right\} \\ &+ V \left\{ E \left[ \bar{e}_{2,HT} \mid (\mathbf{a}_N, F) \right] \mid F \right\}, \end{aligned} \tag{3.13}$$

where  $\mathbf{a}_N$  denotes the  $N$  dimensional vector  $(a_1, a_2, \dots, a_N)$ . By assumption (3.8), the first-phase variance estimator is consistent and, hence, the estimator  $\hat{V}_1 \{ \bar{e}_{2,HT} \mid (\mathbf{a}_N, F_N) \}$  is consistent for  $V_1 \{ \bar{e}_{2\pi} \mid (\mathbf{a}_N, F_N) \}$ .

To show that  $\hat{V}_1 \{ \bar{e}_{2,HT} \mid (\mathbf{a}_N, F_N) \}$  is consistent for  $E \{ V_1 [ \bar{e}_{2\pi} \mid (\mathbf{a}_N, F_N) \mid F ] \}$  we must also show that  $V_1 \{ \bar{e}_{2,HT} \mid (\mathbf{a}_N, F_N) \}$  converges to  $E \{ V_1 [ \bar{e}_{2,HT} \mid (\mathbf{a}_N, F_N) \mid F ] \}$ . Assume first that the second-phase sample is a Poisson sample so that the  $a_j$  are independent Bernoulli random variables. Then

$$\begin{aligned} Cov(a_i a_j, a_k a_m \mid F) &= \kappa_{2i} \kappa_{2j} - \kappa_{2i}^2 \kappa_{2j}^2 \\ &\text{if } ij = km \text{ or } ij = mk \text{ and} \\ Cov(a_i a_j, a_k a_m \mid F) &= 0 \quad \text{otherwise.} \end{aligned} \tag{3.14}$$

Now by (3.6)

$$\begin{aligned} n_N V_1 \left\{ N^{-1} \sum_{i \in A_1} \pi_{1i}^{-1} \kappa_{2i}^{-1} a_i e_i \mid (\mathbf{a}_N, F) \right\} &= \sum_{i=1}^N \sum_{j=1}^N \omega_{ij} \kappa_{2i}^{-1} \kappa_{2j}^{-1} a_i a_j e_i e_j \\ \text{and letting } G_{ij} &= \kappa_{2i}^{-2} \kappa_{2j}^{-2} (\kappa_{2i} \kappa_{2j} - \kappa_{2i}^2 \kappa_{2j}^2), \\ V \left\{ n_N V_1 \left[ N^{-1} \sum_{i \in A_1} \pi_{1i}^{-1} \kappa_{2i}^{-1} a_i e_i \mid (\mathbf{a}_N, F) \right] \mid F \right\} &= \sum_{i=1}^N \sum_{j=1}^N 2 \omega_{ij}^2 G_{ij} e_i^2 e_j^2 \\ &\leq \max |\omega_{ij}| \sum_{i=1}^N \sum_{j=1}^N 2 |\omega_{ij}| G_{ij} e_i^2 e_j^2 \\ &= O_p(n^{-1}) \end{aligned} \tag{3.15}$$

by (3.15) and assumption (3.7). Therefore  $\hat{V}_1 \{ \bar{e}_{2,HT} \mid (\mathbf{a}_N, F) \}$  is consistent for  $E \{ V_1 [ \bar{e}_{2,HT} \mid (\mathbf{a}_N, F) \mid F ] \}$ . Because  $\bar{\mathbf{x}}_{1\pi}$  is not a

function of  $\mathbf{a}_N$ ,  $\hat{V}_1\{\bar{e}_{2,HT} + \bar{\mathbf{x}}_{1\pi}\boldsymbol{\beta}_N\}$  is consistent for  $E\{V_1[\bar{e}_{2,HT} + \bar{\mathbf{x}}_{1\pi}\boldsymbol{\beta}_N | (\mathbf{a}_N, F)] | F\}$ .

To evaluate the second term on the right side of equality (3.14) we have

$$E\left\{N^{-1} \sum_{i \in A_1} \pi_{1i}^{-1} \kappa_{2i}^{-1} a_i e_i | (\mathbf{a}_N, F)\right\} = N^{-1} \sum_{i=1}^N \kappa_{2i}^{-1} a_i e_i$$

and, for Poisson sampling,

$$V\left\{N^{-1} \sum_{i=1}^N \kappa_{2i}^{-1} a_i e_i | F\right\} = N^{-2} \sum_{i=1}^N \kappa_{2i}^{-1} (1 - \kappa_{2i}^{-1}) e_i^2 \tag{3.17}$$

Combining the consistency of  $\hat{V}_1\{\bar{e}_{2,HT} | (\mathbf{a}_N, F)\}$  with (3.17) we have conclusion (3.10) for second-phase Poisson sampling. By the arguments of Hájek (1960) the result also holds for stratified samples. ■

If the first-phase sampling rate is small, the second term of (3.10) is small relative to the first term, and estimator (3.2) can be used for the two-phase regression estimator. If the second term of (3.10) is judged to be important the term can be estimated directly or with replicates.

The regularity conditions of the theorem guarantee that the errors in  $\bar{\mathbf{x}}_1$  and  $\hat{\boldsymbol{\beta}}_2$  are small. Assumption (3.10) requires the replication procedure to be such that the remainder in a Taylor expansion be small relative to the other terms. Theorem 2 is a result for the regression estimator that is closely related to the result of Kim, Navarro and Fuller (2004) for the stratified estimator.

**4. JACKKNIFE VARIANCE ESTIMATION FOR THE NRI**

The U.S. National Resources Inventory (NRI) is conducted by the U.S. Natural Resources Conservation Service in cooperation with Center for Survey Statistics and Methodology of Iowa State University. The survey was conducted as a panel survey in 1982, 1987, 1992, and 1997. Data were collected on soil characteristics, land use, land cover, wind erosion, water erosion, and conservation practices. The sample is a stratified area sample of all states and Puerto Rico. The primary sampling units are areas of land called *segments*. Data are collected for the entire segment on items such as

urban land, roads and water. Detailed data on soil properties and land use are collected at a random sample of points within the segment, where the typical number of points is three. The sample for 1997 contained about 300,000 segments and about 800,000 points. See Nusser and Goebel (1997) for a more complete description of the survey and Fuller (1999) for estimation procedures.

The 2002 sample is a subsample of about 102,000 segments selected from the 1997 sample. The subsample is a stratified sample of the 48 coterminous states, where the segments were placed in twelve strata in each state. The replication method used in NRI variance estimation is a form of “delete-a-group jackknife”. See Kott (2001) for a description. The 1997 NRI is the phase one sample, and the 2002 NRI is the phase two sample. For both the 1997 and the 2002 NRI, we suppose that we have completed an estimation process in which the estimation weights are calibrated for the auxiliary information at the population level and, in the case of the 2002 NRI, for a set of estimates from the 1997 NRI. The information at the population level is the total acres of federal land, acres in census water (water bodies greater than 40 acres and streams greater than 1/8 mile wide) and the acres in a federal program called the Conservation Reserve Program (CRP).

Let  $w_{(1)i}$  and  $w_{(2)i}$  represent the original phase one design weight (1997 NRI) and phase two (2002 NRI) design weight, respectively, for point  $i$ . Let  $w_{(1)i}^*$  and  $w_{(2)i}^*$  represent the fully calibrated regression weights. For any variable,  $y_i$ , the estimate for a 2002 total is computed as

$$\hat{T}_{y(2)} = \sum_{i \in A_2} w_{(2)i}^* y_i$$

For a first phase variable,  $x$ , used to calibrate the second phase,

$$\sum_{i \in A_1} w_{(1)i}^* x_i = \sum_{i \in A_2} w_{(2)i}^* x_i$$

The goal of the variance estimation procedure for the 2002 NRI is to construct  $L$  new sets of weights  $w_{(2)i}^{*(k)}$ ,  $k = 1, 2, \dots, L$ , from which a user of the NRI data is able to compute replicate estimates

$\hat{T}_{y(2)}^{(k)} = \sum_{i \in A_2} w_{(2)i}^{*(k)} y_i$  for any variable  $y_i$ . The user can then calculate a variance estimate for  $\hat{T}_{y(2)}$  as

$$\hat{V}\{\hat{T}_{y(2)}\} = \sum_{k=1}^L c_k \left(\hat{T}_{y(2)}^{(k)} - \bar{T}_{y(2)}\right)^2,$$

where  $c_k$  is a constant determined by the replication procedure and

$$\bar{T}_{y(2)} = L^{-1} \sum_{k=1}^L \hat{T}_{y(2)}^{(k)}.$$

As a first step,  $L$  replicate sets of phase-one weights  $w_{(1)i}^{*(k)}$  are constructed for the 1997 NRI. The sample is divided into  $L$  groups  $G_{1,k}$ ,  $k = 1, \dots, L$ , by ordering the segments geographically and then using systematic sampling to create 29 groups of approximately equal size. The program is flexible with  $L = 29$  in the current version. The geographic ordering is used as an approximation to the stratification of the original sample.

$$w_{(1)i,0}^{*(k)} = 0 \quad i \in G_{1,k}$$

$$= L(L-1)^{-1} w_{(1)i}^* \quad i \notin G_{1,k}.$$

Data are not collected on federal land and census water, but the points falling on federally owned land and falling on census water are kept in the data set so that the surface area of the country is fully represented. The acres in these classifications are known from external sources. Therefore we set  $w_{(1)i,0}^{*(k)} = w_{(1)i}^*$  for every replicate if the point  $i$  is federal land or census water.

For the points that do not fall on federal or census water, the initial replicate weights are calibrated (through a raking procedure) to acreages for CRP available from administrative sources for the previous years, resulting in the final replicate weights  $w_{(1)i}^{*(k)}$ .

Next, initial replicate weights  $w_{(2)i,0}^{*(k)}$  are constructed for the phase two sample using a raking procedure and the first phase controls. Since the 2002 NRI sample was selected as a stratified sample from the 1997 NRI, the  $k$ th 2002 NRI replicate is the portion of the 2002 sample that is in the  $k$ th first phase replicate. The initial replicate weights are

$$w_{(2)i,0}^{*(k)} = 0 \quad i \in G_{2,k}$$

$$= L(L-1)^{-1} w_{(2)i}^* \quad i \notin G_{2,k}$$

except if point  $i$  is federal or census water. For federal and census water the final replicate weight  $w_{(2)i}^{*(k)} = w_{(2)i}^*$ , as for the first phase sample.

For the points that do not fall on federal or census water, the initial 2001 NRI replicate weights  $w_{(2)i,0}^{*(k)}$  for each  $k$  are raked to a set of state-level estimates constructed with the corresponding final 1997 NRI replicate weights  $w_{(1)i}^{*(k)}$ . The categories for which the phase two estimates are calibrated to the phase one are the same as those used in the construction of the original weights  $w_{(2)i}^*$ , and

include several broad landuse categories as well as a wetland classification. At the end of this raking step, we obtain the final replicate weights  $w_{(2)i}^{*(k)}$ , which are appended to the NRI 2002 dataset for the purpose of variance calculation.

At this writing the results of the 2002 NRI have not been released. Table 1 contains some illustrative results. The entries are not true estimates, but the entries give realistic relative values for a midwestern state. The 1997 values are the first phase sample estimates. The values for 2002 are based on the second phase sample which is about one third the size of the 1997 sample. Because of the correlation between 1997 and 2002 values, the standard errors for 2002 are less than the 1997 values multiplied by the square root of three. The computed standard error for large water is zero because large water is an auxiliary variable assumed to be known without error.

Table 1. Illustrative Results (million acres)

Land use	Year	
	1997	2002
Cultivated	240.0	230.0
Cropland	(1.2)	(1.6)
Urban	2.5 (0.6)	2.8 (0.8)
Large Water	20 (0)	20 (0)

### Acknowledgement

This research was funded in part by Cooperative Agreement 68-3A75-14 between the USDA Natural Resources Conservation Service and Iowa State University.

### References

- Binder, D. A. (1996). "Linearization Methods for Single Phase and Two-Phase Samples: A Cookbook Approach", *Survey Methodology* 22, pp. 17-22.
- Binder, D. A., Babyak, C., Brodeur, M., Hidiroglou, M., and Jocelyn, W. (1997), "Variance Estimation For Two-Phase Stratified Sampling", *ASA Proceedings of the Section on Survey Research Methods*, pp. 267-272.

- Cochran, W. G. (1977), *Sampling techniques (3<sup>rd</sup> edition)*, New York: Wiley.
- Fuller, W. A. (1975), "Regression analysis for sample survey", *Sankhyā, Series C, Indian Journal of Statistics* 37, pp. 117-132.
- Fuller, W. A. (1998), "Replication variance estimation for two-phase samples", *Statistica Sinica* 8, pp. 1153-1164.
- Fuller, W. A. (1999), "Estimation procedures for the United States National Resources Inventory", In *Proceedings of the Survey Methods Section of the Statistical Society of Canada*, pp. 39-44.
- Hájek, J. (1960). "Limiting distributions in simple random sampling from a finite population", *Publications of the Mathematical Institute of the Hungarian Academy* 5, pp. 361-374.
- Kim, J. K., Navarro, A., and Fuller, W. A. (2004), "Replication variance estimation for two-phase stratified sampling" unpublished manuscript. Iowa State University, Ames, IA.
- Kim, J. K., and Sitter, R. R. (2003), "Efficient replication variance estimation for two-phase sampling", *Statistica Sinica* 13, pp.641-653.
- Kott, P. S. (1990), "Variance estimation when a first-phase area sample is re-stratified", *Survey Methodology* 16, pp. 99-103.
- Kott, P. S., and Stukel, D. M. (1997), "Can the Jackknife Be Used With a Two-Phase Sample?", *Survey Methodology* 23, pp. 81-89.
- Kott, P. S. (2001), "The delete-a-group jackknife", *Journal of Official Statistics* 17, pp. 521-526.
- Nusser, S. M., and Goebel, J. J. (1997), "The national resources inventory: a long-term multi-resource monitoring programme", *Environmental and Ecological Statistics* 4, pp. 181-204.
- Rao, J. N. K. (1973), "On double sampling for stratification and analytical surveys", *Biometrika* 60, pp. 125-133.
- Rao, J. N. K., and Shao, J. (1992), "Jackknife variance estimation with survey data under hot deck imputation", *Biometrika* 79, pp. 811-822.
- Rao, J. N. K., and Sitter, R. R. (1995), "Variance estimation under two-phase sampling with application to imputation for missing data", *Biometrika* 82, pp. 453-460.
- Särndal, C.-E., Swensson, B., and Wretman, J. (1992), *Model Assisted Survey Sampling*, Springer-Verlag.
- Sitter, R. R. (1997), "Variance estimation for the regression estimator in two-phase sampling", *Journal of the American Statistical Association* 92, pp. 780-787.

---

<sup>1</sup> Department of Statistics, Iowa State University, Ames, IA, 50011, U.S.A.