## A MULTIVARIATE APPROACH TO ESTIMATION IN FINITE POPULATION SAMPLING WHEN NONIGNORABLE NONRESPONSE IS PRESENT

Stephen M. Woodruff, Bureau of Labor Statistics Room 2128, 441 G St. N.W., Washington D.C. 20212

Key Words: Superpopulation Model,

Nonignorable Nonresponse.

## 1. INTRODUCTION

This paper describes a technique for estimation from *typical* sample survey data. Such data is riddled with nonresponse due to processes unknown and, in addition, may result from poorly implemented or documented sample designs. The multivariate structure usually found in sample survey data is used to compensate for these data deficiencies by building all available relevant information into the estimator. This estimator minimizes average error with respect to the stochastic structure defined by this available information.

The estimator derived here is a best linear unbiased estimator (BLUE). It is a generalization of the ordinary regression estimator, Cochran (1977). Regression transformations are used to condition on known outcomes of auxiliary variables and covariates. This is similar to superpopulation prediction theory since both produce estimators which are robust against extreme samples by

conditioning on auxiliary data; see Cassel, Särndal, and Wretman (1977), Royall and Cumberland (1981a) (1981b), Royall and Herson (1973). The close tie between estimating first moments and estimating second moments parallels the similar tie in the EM-algorithm; see Little and Rubin (1987), Pfeffermann (1988), Srivastava and Carter (1986). The estimation problem that this paper solves is very similar to the one that Pfeffermann (1988) discusses. The solution is different from his in two major respects. The covariance matrix is estimated using an identity derived from a superpopulation model and conditioning on sample outcomes of covariate/auxiliary variables to remove bias is done by response group (to be defined).

Estimating finite population means is usually done by weighting sample responses inversely to their probability of selection. This is done item by item resulting in the Horvitz--Thompson estimator; see Cochran (1977), Raj (1968). In cases where auxiliary variables are available, these auxiliaries may be used to derive improved estimators for a single item mean. Ratio and regression estimators are examples of this, but these estimators remain essentially univariate.

Many different but related items of information are collected on sample survey questionnaires. These relationships suggest that survey estimators should be derived under a multivariate framework, which describes these item dependencies. This is particularly true when data items are missing due to either nonresponse or designed nonavailability, as in the case of rotation sampling.

The finite population to be sampled is described in terms of an N×k matrix, Z=(A,C,T), of random variables where the submatrix of auxiliary variables, A is N×m<sub>a</sub>, the submatrix of covariates, C is N×m<sub>c</sub>, and the submatrix of target variables, T is N×m<sub>t</sub>. Let the realization of Z be denoted,  $\mathcal{Z}$ . Let,  $\mathcal{Z}$ =(a,c,t) where each row of  $\mathcal{Z}$ ,  $\mathcal{Z}_i$ = (a<sub>i</sub>,c<sub>i</sub>,t<sub>i</sub>) gives the outcomes (realizations) of the k=m<sub>a</sub>+m<sub>c</sub>+m<sub>t</sub> random variables attached to the i<sup>th</sup> member of the finite population for 1≤i≤N.

To illustrate the estimation methodology consider a simple example in which N=6 and k=4. The rows of this example of Z are normally and independently, identically distributed (iid), each with a mean vector,  $(\mu_a, \mu_c, \mu_t)$ . Some of the columns of  $\mathcal{Z}$  are known for each member of the finite population. These are referred to as auxiliary variables. A sampling distribution which is usually a function of these auxiliary variables generates a sample of n<N rows of  $\mathcal{Z}$  (n=3 in the example). The statistician then attempts to observe the remaining components (non auxiliary variables) of  $\mathcal{Z}_i$  for all i in the

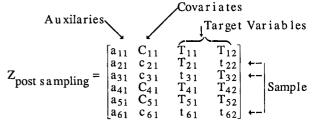
sample. Unfortunately, nonresponse may occur, and only some subset of the nonauxiliary components of  $\mathcal{Z}_i$  may be

observed. Usually, this subset will vary from sample member to sample member.

These nonauxiliary variables may be further divided into two groups — covariates and target variables. The goal of the sample survey is to measure the finite population means of each target variable. The target variables are observed only for sample members, and they suffer the nonresponse. The covariates are observed without nonresponse for all sample members.

For this example,  $\mu_t$  is to be estimated from the following four pieces of information: 1) all the auxiliary variable outcomes in  $\mathbb{Z}$ , together with  $\mu_a$  and  $\mu_c$ , 2) all the covariate outcomes in  $\mathbb{Z}$  for the sample members, 3) a subset of the target outcomes in  $\mathbb{Z}$  for the sample members. This subset is the result of nonresponse and varies from sample member to sample member, and 4) the available information on the distribution of the random matrix Z.

Suppose that rows 2, 3, and 6 were sampled from Z. In the interests of simplicity the sampling distribution will be ignored in this example. The indicator functions are used in the derivation of the BLUE in subsequent sections. In row 2,  $T_{21}$  was a nonresponse, as was  $T_{32}$  in row 3, and in row 6 no nonresponse occurred. The available sample and auxiliary data for this example is summarized in the following matrix of random variables and realizations of random variables:



The second source of information is the distribution of the random matrix Z. Since mean square error (MSE) is the estimation criterion, the covariance matrix of Z will suffice and since the rows of Z are iid the covariance matrix of a  $Z_i$  will suffice. Let  $\sum_{z}$  denote the covariance matrix of a  $Z_i$ . When Z contains all available correlated variables, its rows may be modeled as iid random vectors.

Z is initially a matrix of iid random vectors in which some random components are being replaced with their realizations. First, the auxiliary data arrives, then the sample data. When the auxiliary data is observed, Z becomes a mixture of this known auxiliary data, and as yet unobserved random variables. Conditional on the observed auxiliary data, the rows of Z are no longer iid since this auxiliary data will distinguish the distribution of different rows of Z. These distributional deviations between rows of Z given the auxiliary data, are generally dealt with in the sampling literature by unequal probability sampling (stratification and clustering). Then the sample target variables are weighted inversely to selection probabilities, to get an estimator for the target variable means, which is unbiased with respect to the sampling distribution. This is the Horvitz–Thompson estimator.

This example shows a different approach to estimating target means which is based on transformations of the target components of Z. These transformations harness the knowledge contained in the distribution of Z to adjust the sample target variables for the auxiliary variable and covariate observations in the sample units.

These transformations are linear regression adjustments. They are applied to averages over each response group of sample row vectors of Z. A response group is a subset of sample units that responded to exactly the same target variables. In the example, there are three response groups, one for each sample member and the response group averages are the trivial averages (averages over a single item). To describe these response group transformations, partition  $\Sigma_{r}$  as follows:

$$\Sigma_{z} = \begin{bmatrix} H_{g} & H_{gt} \\ H_{tg} & \Sigma_{t} \end{bmatrix} \text{ where } H_{g} \text{ is the covariance}$$

matrix of  $(A_{i1}, C_{i1})$ ,  $\sum_{t}$  is the covariance matrix of  $(T_{i1}, T_{i2})$ and  $H_{gt}$  is the covariance matrix between these two subvectors. The transformation to recover iid subsequent to observing auxiliaries and covariates gives us a new vector,  $(W_{i1}, W_{i2})$ , given by:

$$(W_{i1}, W_{i2}) = (T_{i1}, T_{i2}) - (A_{i1} - \mu_a, C_{i1} - \mu_c)H_g^{-1}H_{gt}$$
  
i=2.3, and 6.

for i=2,3, and 6. Note that conditional on the auxiliaries and covariates,  $E(W_{i1}, W_{i2})=(\mu_{t1}, \mu_{t2})$  and  $Cov(W_{i1}, W_{i2})=\sum_t H_{tg}H_g^{-1}H_{gt} = \sum_w$  for i=2,3, and 6. The last equality is the definition of  $\sum_w$  and it implies that  $Var(W_{ij}|(A_i, C_i))=Var(T_{ij}|(A_i, C_i))\leq Var(T_{ij})$  for i in the sample and j=1 or 2. Nonresponse resulted in  $T_{21}$  and  $T_{32}$ being unobserved, and thus,  $W_{21}$  and  $W_{32}$  are unobservable.  $W_{22}$ ,  $W_{31}$ ,  $W_{61}$ , and  $W_{62}$  remain, and they are related to  $\mu_t$  by the following expression:

$$Y = \begin{cases} W_{22} \\ W_{31} \\ W_{61} \\ W_{62} \end{cases} = \begin{cases} 0 & 1 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{cases} \mu_{t}' + \varepsilon = X\mu_{t}' + \varepsilon, \quad (*)$$

where the first and last equalities in this line are definitions of Y and X. The covariance matrix of  $\varepsilon$  is the diagonal matrix of the appropriate submatrices of  $\Sigma_{W}$ . Let  $\Sigma$  denote the covariance matrix of  $\varepsilon$ . Then the least squares BLUE for  $\mu_t'$  is  $\hat{\mu}_t' = (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}Y$  and its variance is  $(X'\Sigma^{-1}X)^{-1}$ .

The sampling distribution and the nonresponse mechanism may be included in this BLUE by including the sample indicator function in Z as an auxiliary variable and including the response indicator vector in Z as covariates. The response indicator vector is a vector of zeros and ones. It is the same dimension as the vector of target variables, and its 1<sup>th</sup> component is one if the 1<sup>th</sup> component of the target vector is observed (a response) and zero if not. The sample indicator function is defined for each member of the finite population as one if the member is selected and zero if not.

By properly designing Z the effects on the available data of stratification and clustering in the selection process are removable by conditioning the target variables on known auxiliary variables and covariates.

The reader by wonder why the regression adjustments by response group of target variables are necessary. In particular, why not delete the unobserved components of Z, line up the observed components of Z, as was done with the  $(W_{ij})$ , (see Pfeffermann 1988) to get a version of (\*) with  $\mu = (\mu_a, \mu_c, \mu_t)'$  on the right hand side? This version of (\*) implies a BLUE for  $\mu_t$  but this estimator for  $\mu_t$  is inappropriate for two reasons.

First, the linear regression equations used in later sections to describe the relationship between components of a row of Z may induce linear dependencies between components of  $\mu$ . Therefore, the desired solution solves a restricted minimization problem, but the normal equations referred to in the preceding paragraph with  $\mu$  on the right hand side solve the unrestricted minimization.

The second reason may be more intuitive. It is usually desireable to condition decisions on all available relevant information (see conditionality in Cox and Hinkley, (1974)).

Note the following pounts about this estimation methodology:

1) there is extreme flexibility for including all available relevant data and data relationships

2) the BLUE is conditioned on all known and readily estimable quantities

3) when the response mechanism is known the response indicator vector may be included as a covariate

4) collapsing the data by response group and doing the regression adjustment for auxiliary and covariate means by response group may reduce the effect of nonignorable nonresponse in cases where the nonresponse mechanism is unknown

5) the sampling distribution may be integrated into the estimator by including the sample indicator function in Z as an auxiliary variable

6) a variance estimator for the BLUE is included in its derivation and this variance estimator,  $(X \stackrel{\frown}{\Sigma}^{-1} X)^{-1}$ , measures both variance with respect to the distribution the BLUE inherits from Z and, in some situations, the variance of the BLUE with respect to repeated sampling.

7) the BLUE derived here will usually outperform the Horvitz-Thompson estimator. A model based ratio estimator for the estimation problem considered in the rest of this paper is compared to the BLUE. This ratio estimator has been in use at BLS for several decades and in spite of considerable effort to improve it, nothing substantially better has been found until the BLUE derived here.

#### 2. THE LINEAR REGRESSION EQUATION FOR AN EMPLOYMENT SURVEY

The Bureau of Labor Statistics' (BLS) Current Employment Statistics Survey (CES) is a monthly survey of over 275,000 nonagricultural business establishments. This survey produces estimates of total employment, women and production workers, and hours and earnings. These estimates are made for over 1500 industries. Monthly employment level and month to month change are of primary importance to the users of this survey. An estimator for monthly employment level is derived in section 3. This estimator is an upscale version of the estimator derived in section 1 and uses the data relationships to be defined next. The linear regression equation states that expected employment in month j is proportional to actual employment during month j–1 for each member of the population. The constant of proportionality,  $\beta_j$ , is the same for each member of the finite population. A response mechanism is hypothesized, which relates the probability of response for a target variable to its realized value. This response mechanism was introduced as a way to test the robustness of the BLUE, derived under ignorable nonresponse. A second BLUE, which assumes an approximate knowledge of this response mechanism, is also derived. This second BLUE serves as a benchmark with which to compare the BLUE under ignorable nonresponse (or complete ignorance of the response mechanism).

The linear regression equation relating the components of  $Z_i$  is given next. There is a single auxiliary variable, benchmark employment, available for this employment survey. The benchmark employment is the total establishment employment for the initial month (j=0) and is known for every business establishment in the finite population. The benchmark employment for establishment i is denoted  $A_i$  for  $1 \le i \le N$ . The total employment for establishment i during each month after the benchmark month are the target variables, denoted  $T_{i1}, T_{i2}, \dots, T_{imt}$ . These are the employments in establishment i for j=1, 2,..., m months after the benchmark month.

 $m_t$  months after the benchmark month. The response indicator vector is the only set of potential covariates. Thus, this vector will be denoted,  $C_i$ , and not included in  $Z_i$  until later. Under ignorable nonresponse  $C_i$  is uncorrelated with  $(A_i, T_i)$  and the vectors,  $Z_i = (A_i, T_{i1}, T_{i2}, T_{im})$ ,  $1 \le i \le N$ , are

iid. The components of Z<sub>i</sub> are related as follows:

Let  $A_i = \beta_0 + \lambda_{i0}^{\prime}$ ,  $T_{i1} = \beta_1 a_i + \lambda_{i1}^{\prime}$ , and  $T_{ij} = \beta_j t_{ij-1} + \lambda_{ij}$  for j > 1 (2.1) where  $\{(\lambda_{i0}, \lambda_{i1}, \dots, \lambda_{im_t}): 1 \le i \le N\}$  are iid random vectors

with mean zero and diagonal covariance matrix. The  $\{\beta_j\}$  are unknown constants, the variance of  $\lambda_{i0}$  exists, and the variance of  $\lambda_{ij}$  given  $t_{ij-1}$  is proportional to  $t_{ij-1}$  (Var $(\lambda_{ij}|t_{ij-1})=r_jt_{ij-1}$ ) where the  $\{r_j:1\leq j\leq m_t\}$  are unknown constants (Var $(\lambda_{i1}|a_i)=r_1a_i$ ).

This conditional variance property for the  $\{\lambda_{ij}\}$  has been shown to be appropriate for total employment as measured in the CES (West (1981), Royall (1981)). The covariance matrix of any  $Z_i$ ,  $\Sigma_z$ , is necessarily of the form  $(b_{lj})$  where:

the diagonal entries are:  $b_{ll} = \sigma_l^2$  for  $l = 0, 1, 2, ..., m_t$  (2.2). and the offdiagonal entries are:  $b_{lj} = \sigma_l^2 \prod_{k=l+1}^{j} \beta_k$  for  $l < j \le m_t$ 

and where  $\sigma_l^2$ , for  $l=0,1,2,..., m_t$  are unknown constants (functions of the  $\{r_i\}$  and  $\{\beta_i\}$ ).

 $C_i = (C_{i1}, C_{i2}, \dots, C_{im_t})$  where  $C_{ij} = 1$  if the data for

target variable j in unit i is a response and  $C_{ij} = 0$  if not.

A nonignorable nonresponse mechanism is introduced next. The BLUE under this additional structure is used in the simulation study (section 5) as a benchmark with which to compare the estimators to be evaluated there. Conditional on  $t_i$ , it is assumed that the components of  $C_i$  are independent and that:

are independent and that:  $P(C_{ij}=1|T_{ij}=t_{ij})=f_{j}t_{ij}+g_{j} \text{ for } 1 \le i \le N \text{ and } 1 \le j \le m_{t} (2.3).$ 

& 
$$P(C_i = u_l | T_i = t_i) = \prod_{j=1}^{m} P(C_{ij} = u_{lj} | T_i = t_i),$$

where the {( $f_j, g_j$ )} are defined so that  $f_j t_{ij} + g_j$  is between zero and one for all  $t_{ij}$  and where  $(u_{lj})$  is the  $2^{m_t} \times m_t$ matrix with rows that consist of all possible distinct  $m_t$ -tuples of zeros and ones (its last row consists of all zeros).  $u_l$  denotes the  $l^{\text{th}}$  row of  $(u_{lj})$ . Note that when  $f_j=0$ for all j, the nonresponse is ignorable.  $\mu_c$  is the expected value under both (2.1) and (2.3)

 $\mu_c$  is the expected value under both (2.1) and (2.3) of C<sub>i</sub>.

Let 
$$\overline{T}_{\alpha} = (1/n_{\alpha}) \sum_{l \in S_{\alpha}} T_{l}$$
 where  $s_{\alpha}$  is an arbitrary

subset of the first N integers and  $n_{\alpha}$  is the size of  $s_{\alpha}$ . Define  $\overline{A}_{\alpha}$  and  $\overline{C}_{\alpha}$  similarly.

 $\alpha$   $\alpha$ <u>Before</u> considering either sampling or nonresponse, the stochastic structure given in (2.1) and (2.3) may be summarized for <u>arbitrary subsets</u> s<sub> $\alpha$ </sub> as follows:

$$E(\overline{A}_{\alpha}, \overline{C}_{\alpha}, \overline{T}_{\alpha}) = \mu = (\mu_{a}, \mu_{c}, \mu_{t}) \qquad \text{and}$$

$$Cov(\overline{A}_{\alpha}, \overline{C}_{\alpha}, \overline{T}_{\alpha}) = (1/n_{\alpha}) \begin{bmatrix} \sum_{a} \sum_{a, t} D \sum_{a, t} \\ D\sum_{ta} \sum_{c} \sum_{t} D \\ \sum_{ta} D\sum_{t} \sum_{t} \end{bmatrix} \text{ where:}$$
1) D is the diagonal matrix of  $(f_{1}, f_{2}, \dots, f_{m_{t}})$ 

2)  $\left[\sum_{i}^{\Delta} \sum_{i}^{\Delta} \sum_{i}^{t}\right] = (b_{i})$  as given in (2.2).  $\sum_{a}$  is the variance of the auxiliary variable,  $\sum_{i}$  is the covariance matrix of the target variables, and  $\sum_{at}$  is the covariance matrix between the auxiliary variable and the target variables  $(\sum_{a} \sum_{i} \sum_{a})$ 3)  $\sum_{c}$  is the covariance matrix of the covariates with respect to both (2.1) and (2.3) It has off diagonal

3)  $\sum_{c}$  is the covariance matrix of the covariates with respect to both (2.1) and (2.3). It has off diagonal elements, which are the same as those of  $D\sum_{i}D$ , and diagonal elements,  $\{\mu_{cj}(1-\mu_{cj})\}_{j=1}^{m_{t}}$  where  $\mu_{cj}$  is the j<sup>th</sup> component of  $\mu_{c}$ . (2.4)

Following the methodology outlined in the introduction,

 $E(\overline{T}_{\alpha}|\overline{A}_{\alpha},\overline{C}_{\alpha}) = \mu_{t} + (\overline{A}_{\alpha} - \mu_{a}, \overline{C}_{\alpha} - \mu_{c})H_{g}^{-1}H_{gt}.$  If  $\overline{W}_{\alpha}$  is defined as:

$$\begin{split} \overline{W}_{\alpha} &= \overline{T}_{\alpha} - (\overline{A}_{\alpha} - \mu_{a}, \overline{C}_{\alpha} - \mu_{c}) H_{g}^{-1} H_{gt}, \text{ then} \\ E(\overline{W}_{\alpha} | \overline{A}_{\alpha}, \overline{C}_{\alpha}) &= \mu_{t} , \text{ the vector to be estimated. The covariance matrix of } \overline{W}_{\alpha} \text{ is } (1/n_{\alpha}) \Sigma_{w} \text{ where } \Sigma_{w} &= \Sigma_{t} - H_{tg} H_{g}^{-1} H_{gt}, \quad H_{gt} = \begin{bmatrix} \Sigma_{a} & 1 \\ \Sigma_{d} & D \end{bmatrix} \text{ and } H_{g} = \begin{bmatrix} \Sigma_{a} & \Sigma_{a} \\ D \Sigma_{a} & \Sigma_{c} \end{bmatrix}. \\ \overline{W}_{\alpha} \text{ and } (\overline{A}_{\alpha}, \overline{C}_{\alpha}) \text{ are uncorrelated, therefore, the} \end{split}$$

conditional expected value of  $\overline{W}_{\alpha}$  given  $(\overline{A}_{\alpha}, \overline{C}_{\alpha}) = (\overline{a}_{\alpha}, \overline{c}_{\alpha})$ is approximately  $\mu_t$  and the conditional covariance matrix of  $\overline{W}_{\alpha}$  given  $(\overline{A}_{\alpha}, \overline{C}_{\alpha}) = (\overline{a}_{\alpha}, \overline{c}_{\alpha})$  is approximately  $(1/n_{\alpha})\Sigma_{W}$ (if  $n_{\alpha}$  is large enough so that  $(\overline{A}_{\alpha}, \overline{C}_{\alpha}, \overline{T}_{\alpha})$  is approaching normality, then these approximations approach exact equality).

Now suppose a probability sample is selected according to a sampling scheme, which is a function of the auxiliary variable. Let  $\pi_i$  denote the probability of selection

for the i<sup>th</sup> member of the finite population (i<sup>th</sup> row of Z). After the sample has been selected and the data (minus the nonresponse) has been observed for the m<sub>t</sub> target variables, the sample units are grouped according to their response patterns. Let  $\{\mathscr{C}_{l}\}$  be the set of these groups, and  $\mathscr{C}_{l} = \{i: 1 \le i \le N \text{ and } C_{i} = u_{l}\}$ .  $\mathscr{C}_{l}$  is the set of sample units which have responses for exactly those target variables corresponding to the nonzero components of  $u_l$ . For all l

such that  $\mathscr{C}_{l} \neq \phi$ , let  $\overline{W}_{l}$  be the  $\overline{W}_{\alpha}$  defined above with  $s_{\alpha} = \mathscr{C}_{l}$  and let  $n_{l}$  be the number of sample units in  $\mathscr{C}_{l}$ .

A modest generalization of the methodology outlined in the introduction will produce a BLUE for the target finite population means. This BLUE is derived in section 3 together with estimates of the unknown regression equation parameters.

#### 3. THE BEST LINEAR UNBIASED ESTIMATOR AND ESTIMATORS OF THE UNKNOWNREGRESSION PARAMETERS

Let  $G_l$  be the response indicator matrix for the  $l^{th}$ response group,  $\mathscr{C}_{l}$ , formed by deleting from the  $m_t \times m_t$ identity matrix all columns, j, such that  $u_{li}=0$ . Then  $y_l = \overline{W}_l G_l$  consists of only those components of  $\overline{W}_l$ , which are observable (responses). The observed data and data relationships are summarized:  $\bar{Y}=X\mu_{t}'+\epsilon$ (3.1)

where  $X' = (G_1, G_2, ..., G_2^{m_t}, G_2^{m_t}),$   $Y' = (y_1, y_2, ..., y_2^{m_t}, -1).$ and  $\varepsilon$  is the random vector with mean, zero and covariance matrix,  $\Sigma$  given as the block diagonal matrix of the  $\Sigma_l = (1/n_l)G_l \Sigma_w G_l$  for all l with  $n_l > 0$  from l=1 in the upper

left to  $l=2^{m_t}-1$  in the lower right. X, Y, and  $\varepsilon$  contain only those arrays  $\{G_l, y_l, and (\varepsilon_l, \Sigma_l)\}$  such that  $n_l > 0$ .

The BLUE, 
$$\hat{\mu}_{l}$$
 is:  $(X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}Y$   
= $\left\{\sum_{l \in \mathscr{I}} G_{l}\Sigma_{l}^{-1}G_{l}\right\}^{-1}\sum_{l \in \mathscr{I}} G_{l}\Sigma_{l}^{-1}y_{l}$  (3.2)

 $\mathscr{I} = \{l:n_j > 0\}.$  $\Sigma_z$ , D,  $\mu_a$ , and  $\mu_c$  are necessary to derive the  $\{y_l\}$ ,

 $\Sigma_{w}$ , and  $\hat{\mu}_{t}$ . D is known.  $\mu_{c}$ , the vector of population response indicator means, contains the means of the only covariates for this problem and is estimable with the Horvitz-Thompson estimator (recall the covariates suffer no nonresponse). Denote this estimator,  $\mu_c$ .  $\Sigma_z$  is estimated by inverting  $\sum_{z}^{1}$ , where  $\sum_{z}^{1}$  is estimated by inserting estimates of  $\beta_{j}$  and  $\tau_{j}$  for  $0 \le j \le m_{t}$  into  $\sum_{z}^{-1}$ . Recall the  $\{\beta_{j}\}$ are defined in (2.1).  $\tau_0 = \sigma_0^2$ , and  $\tau_j = \sigma_j^2 - \beta_j^2 \sigma_{j-1}^2$  for  $1 \le j \le m_t$ , where the  $\{\sigma_j^2\}$  are defined in (2.2). (2.2) also implies that the inverse of  $\Sigma_{r}$  is the tridiagonal matrix with diagonal elements given by:

 $(1/\tau_j) + (\beta_{j+1}^2/\tau_{j+1})$  for  $0 \le j < m_t$  and with the last diagonal entry,  $1/\tau_m$ . The  $(j,j+1)^{th}$  and  $(j+1,j)^{th}$  off-diagonal elements of  $\sum_{z}^{1}$  for  $0 \le j \le m_t$  are  $-(\beta_{j+1}/\tau_{j+1})$ . All other entries of  $\sum_z^{-1}$  are zero.  $\beta_0$  is estimated with the finite population mean,

 $\bar{a}=\beta_0$ , of the auxiliary variables  $\{a_i\}_{i=1}^N$ .  $\beta_1,\beta_2,...,\beta_{m_t}$  are estimated with the ratio of matched sums of target variables. For example,  $\hat{\beta}_j = (\sum_{i \in s_j} t_{ij} / \sum_{i \in s_j} t_{ij-1})$ , where  $s_j$  is the set of sample units with responses for both target variable j and target variable j-1. Let  $\hat{\theta}_j = \prod_{l=0}^{J} \hat{\beta}_l$  for j=0,1,...,m<sub>t</sub>. Initially, let  $\hat{\mu}_a = \hat{\theta}_0$  and  $\hat{\mu}_t^o = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_m)$ . (3.3)  $\tau_i$  is the expected value of the conditional variance

of  $T_{ij}$  given  $(a_i, t_{i1}, \dots, t_{ij-1})$  for  $j \ge 1$ .  $\tau_j = r_j \mu_{tj-1}$  for  $j = 1, 2, \dots$  $m_t$ . Let  $\hat{\tau}_j = \hat{r}_j \hat{\theta}_{j-1}$  for j=1,2,...,  $m_t$ . where:  $\hat{r}_j = (1/[n(s_j)-1]) \sum_{i \in S_i} (t_{ij} - \hat{\beta}_j t_{ij-1})^2 / t_{ij-1}$  and  $n(s_j)$  is

the size of s<sub>1</sub>.

$$\hat{\tau}_{0} = \hat{\sigma}_{0}^{2} = (1/N-1) \sum_{i=1}^{N} (a_{i} - \bar{a})^{2}.$$

This completes the estimation of  $\sum_{z}$ , and with this  $\hat{\Sigma}_{z}$ ,  $\hat{\mu}_{c}$  and  $\hat{\mu}_{a}$  the  $\{W_{l}\}$  and  $\hat{\Sigma}_{w}$  may be calculated and in turn,  $\hat{\mu}_{t}$  as given in (3.2). With the exception of one additional covariate, this is how the BLUE,  $\hat{\mu}_{t}$ , was derived in the simulation study to be described in section 5. This BLUE was calculated in two versions. The first version assumes D=0 and would be used when the nonresponse mechanism is either ignorable or unknown. The second version assumes an approximate knowledge of D to derive the BLUE. This second BLUE provides a benchmark with which to compare the BLUE under the weakened model (assuming D=0).

If  $\hat{\mu}_t$  varies much from the initial estimate  $\hat{\mu}_{t}^{0} = (\hat{\theta}_{1}, \hat{\theta}_{2}, \dots, \hat{\theta}_{m}), \text{ which is defined above, then this } \hat{\mu}_{t}$ may be used in place of  $\hat{\mu}_t^o$  to reestimate  $\sum_{\mathbf{z}}$  (and the  $\{\overline{W}_l\}$ ) and thus to reestimate  $\hat{\mu}_t$ . This may be continued until convergence. The simulation results described in section 5 do not use this iterative reestimation process, which characterizes the EM-algorithm.

4. THE EFFECT OF IGNORING PERTINENT DATA

The effect of disregarding either useful data or known data relationships is quantified next. This is done by computing the MSE under the full (or strong) model of the BLUE derived under a model (weak) that assumes only some proper subset of the full models data and data relationships. This MSE is always larger than the variance of the BLUE derived under the full model. The bias and variance components of this additional MSE are derived below.

Suppose the same subset of the auxiliary variables and covariates is deleted from each Z<sub>i</sub>. Let these random variables be denoted by,  $Z_{i1}$ . Let the remaining auxiliary variables and covariates be  $Z_{i2}$ . Then the two random vectors,  $(Z_{i1}, Z_{i2})$  and  $(A_i, C_i)$ , contain exactly the same random variables but they may be ordered differently.

Let: 
$$E(Z_{11}, Z_{12}, T_1) = (\mu_1, \mu_2, \mu_1)$$
 and  
 $Cov(Z_{11}, Z_{12}, T_1) = \begin{bmatrix} \sum_{1}^{1} \sum_{2}^{12} \sum_{2}^{11} \sum_{1} \sum_{2} \sum_{1}^{12} \sum_{1} \sum_{1}$ 

Let  $\mu_t^{c}$  denote the BLUE under the model where

the rows of the matrix, Z, are  $(Z_{i2}, T_i)$ . In this case,  $\overline{W}_l^c = \overline{T_l} - (Z_{l2} - \mu_2) \Sigma_2^{-1} \Sigma_{21} = \overline{W_l} + b_l$ , where  $\overline{W_l}$  is as defined in section 2 and  $\overline{W}_l^c$  is its analogue for this weak model. The difference between  $\overline{W}_l^c$  and  $\overline{W}_l$  is  $b_l$  and given by the expression:

 $b_l = O_{l1}L\Sigma_{1t} - O_{l1}LM\Sigma_{2t} - O_{l2}M'L\Sigma_{1t} + O_{l2}M'LM\Sigma_{2t}, \text{ where:}$ 

$$O_{l1} = (\overline{Z}_{l1} - \mu_1), \quad O_{l2} = (\overline{Z}_{l2} - \mu_2),$$
  

$$L = (\Sigma_1 - \Sigma_{12} \Sigma_2^{-1} \Sigma_{21})^{-1}, \text{ and } M = \Sigma_{12} \Sigma_2^{-1}.$$

If  $y_l^c$  denotes the weak model version of  $y_l$ , then

 $y_l^c = W_l^c G_l = W_l G_l + b_l G_l = y_l + b_l G_l$ . Recall that since all auxiliary variables and covariates are being conditioned upon, the conditional variance of  $b_l$  is zero.  $Cov(y_l) = Cov(y_l) = Cov(y_l) + Cov(y_l) = Cov(y_l) + Cov(y_l) +$ 

 $\operatorname{Cov}(\mathbf{y}_l^{\mathsf{c}}) = (1/n_l) \mathbf{G}_l \Sigma_{\mathsf{w}} \mathbf{G}_l.$ 

The BLUE under the weak model is:  $\hat{\mu}_t^{C'}$ =

 $(X'\Sigma^{-1}X) \stackrel{-1}{} X'\Sigma^{-1}Y^{c}$ , where X and  $\Sigma$  are exactly as given in section 3 and  $Y^{c} = (y_{1}^{c}, y_{2}^{c}, \dots, y_{2}^{c}m_{t-1})^{t}$ . This weak model BLUE is related to the BLUE under the full model by:  $\hat{\mu}_{t}^{c} = \hat{\mu}_{t}^{t} + HB$  where

 $H = (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1} \text{ and } B = (b_1G_1, b_2G_2, \dots, b_2m_{-1}G_2m_{-1})'. \text{ Since HB is}$ 

constant under the full model,  $\hat{\mu}_t^c$  and  $\hat{\mu}_t$  differ, only by this constant bias, HB.

 $MSE(\overset{\land c}{\mu_t}) = Cov(\overset{\land}{\mu_t}) + HBB'H'$  and this implies

that the MSE of each component of  $\hat{\mu}_t^c$  is larger than the

# MSE (variance) of the corresponding component of $\hat{\mu}_t$ .

Just as ignoring pertinent covariates or auxiliary variables will add only bias to the BLUE, it can be shown that ignoring pertinent target variables adds only variance to the BLUE. 5. SIMULATION STUDY

The empirical MSEs of four estimators are compared here. The four estimators are:

1)  $\hat{\mu}_{t}$  under (2.1) and (2.3), D=0

D=0)

2) the BLUE, 
$$\hat{\mu}_t^c$$
, under (2.1) alone ( =  $\hat{\mu}_t^c$  with

3) HT, the Horvitz–Thompson estimator for the target variables adjusted for nonresponse. The nonresponse adjustment factor is the ratio of total sample weights to total responding unit weights for each target variable

4) LR, the link relative estimator given by (3.3),  $(\hat{\theta}_2, \hat{\theta}_3, \hat{\theta}_4, \hat{\theta}_5)$ . This is the estimator currently used by the BLS for the employment survey referred to in the introduction.

The populations used for this study were derived from a BLS data base with data for six consecutive months from 1110 establishments.  $Z_i = (A_i, T_{i1}, T_{i2}, \dots, T_{i5})$ , where the components of  $Z_i$  are related by (2.1) and lie in the interval (0,1000). For this simulation  $T_{i1}$  is treated as a covariate (no nonresponse) and conditioned upon by using the Horvitz-Thompson estimator to estimate  $\mu_{t1}$ . The vector,  $(\bar{t}_2, \bar{t}_3, \bar{t}_4, \bar{t}_5)$ , is to be estimated where,  $\bar{t}_j$  is the finite population mean of the j<sup>th</sup> target variable, j=2,3,4,5. A sample of size n=200 is selected for each of the 500 replications. Sampling is by probability proportionate to size as measured by the  $\{a_i\}$ .

The target variables are subjected to three different response mechanisms. A separate table is given for each of these response mechanisms. For each target variable and response mechanism an interval is tabulated. The lower limit is the probability of response when the target variable is 0, and the upper limit is the probability of response when the target variable is 1000. For a target value between 0 and 1000 the probability of response is the linear interpolation between these two endpoints.

**Response Mechanisms** 

RM\Targets	2	•	3	4	5
1	(.5, 1.0)		(.4,1.0)	(.2,.9)	(.2,1.0)
2	(.9,.9)		(.8,.8)	(.65,.65)	(.5,.5)
3	(.8, 1.0)		(.75,1.0)	(.6,1.0)	(.5,.7)
					.1 . 11

For each target variable and estimator in the tables, the estimated MSE is the average squared difference between the target variable estimator and the finite population mean for that target variable. This average is

		1 MSEs, Bias		nces.	
MSE /(Bias) for Response Mechanism 1.					
Target Variable	2	3	4	5	
<u>Estimator</u>					
HT	66.9	83.0	213.1	228.7	
	(6.6)	(7.3)	(12.6)	(13.3)	
LR	4.1	11.5	15.1	22.3	
	(.55)	(1.17)	(1.5)	(1.89)	
$\hat{\mu}_t^c$	3.4	7.6	7.8	10.9	
<sup>μ</sup> t					
	(.38)	(.71)	(.80)	(1.07)	
μ̂,	3.2	7.1	7.1	9.7	
r-t	(10)	( 20)	( 00)	( 00)	
	(.18)	(.39)	(.09)	(06)	
-1 $-1$ $-1$					

The Variance Estimator of  $\hat{\mu}_t^c$  from  $(X'\Sigma^{-1}X)^{-1}$ ; Its Estimated Expected Value and Its Estimated Standard Error.

Louindied Expect	ieu failue and			
$\hat{E}(\hat{V}(\hat{\mu}_{t}^{c}))$	3.8	5.8	9.4	12.6
σ	.5	.7	1.2	1.5
Table	2. Estimated	MSEs, Biase	s. and Varian	nces.
		esponse Mech		
Target Variable	2	3	4	5
Estimator				-
HT	9.4	17.2	34.2	55.8
	(.13)	(12)	(.59)	(31)
LR	3.0	6.9	9.4	15.3
	(.12)	(.10)	(0.0)	(14)
µ <sup>c</sup> t	2.8	6.1	7.0	10.7
r't				1011
	(.15)	(.15)	(.09)	(04)
The Variance Es	timator of $\hat{\mu}_t^c$	from $(X'\hat{\Sigma}^{-1})$	X) <sup>-1</sup> ; Its	
Estimated Expec	ted Value and	d Its Estimated	d Standard E	rror.
$\hat{E}(\hat{V}(\hat{\mu}_{t}^{c}))$	3.3	5.2	8.4	12.5
~ <b>L</b>	•	-		

.5

0.9

1.4

.2

Table 3. Estimated MSEs, I	Biases, and Variances.
MSE /(Bias) for Response I	Mechanism 3.

Target Variable Estimator	2	3	4	5
HT	16.1 (2.3)	21.3 (2.69)	48.0 (4.81)	49.9 (2.75)
LR	3.1 (.25)	7.3 (.44)	8.0 (.55)	11.7 (.51)
μ <sup>c</sup> t	2.9	6.3	6.4	9.3
t	(.19)	(.27)	(.24)	(.11)
μ̂ <sub>t</sub>	2.8	6.3	6.3	9.3
•	(.14)	(.20)	(.07)	(09)

The Variance Estimator of  $\hat{\mu}_t^c$  from  $(X'\hat{\Sigma}^{-1}X)^{-1}$ ; Its

Estimated Expected Value and Its Estimated Standard Error.

$\hat{E}(\hat{V}(\hat{\mu}_{t}^{c}))$	3.3	5.1	7.9	11.5
σ	.4	.5	.3	.5

over the 500 replications of the sampling and estimation process. This number of replications gives the MSE estimates (upper entry – unbracketed), relative errors of around 5% and nearly always less than 10%.

The lower portion of each table evaluates

 $(X' \sum^{-1} X)^{-1}$  as a variance estimator. The results of section 4 say that this estimated covariance matrix may be used to estimate the variance of both  $\hat{\mu}_t^c$  and  $\hat{\mu}_t$  (recall that  $\hat{\mu}_t^c$  and  $\hat{\mu}_t$  differ only by the covariates in C<sub>i</sub>). The average of the diagonals of the  $\{(X' \sum^{-1} X)^{-1}\}$  from the 500 replications is displayed in the first row. The estimators of standard error for the diagonal of a single  $(X' \sum^{-1} X)^{-1}$  is given below each of these estimated means.

In Table 2, the estimated MSE of  $\hat{\mu}_t^c$  for target variable 5 is 10.7, and the variance estimator for the estimator of target variable 5 is in the  $2\sigma$ -interval,

12.5 ± 2(1.4), with high probability, (if normality holds,  $\approx$ .95). For Table 2, MSE and variance of  $\hat{\mu}_t^c$  are the same because nonresponse is ignorable and  $\hat{\mu}_t^c = \hat{\mu}_t$ .

In Table 1, where the nonignorability is most extreme,  $\hat{\mu}_t^c$  still does relatively well. As predicted in section 4, the differences in MSE between  $\hat{\mu}_t$  and  $\hat{\mu}_t^c$  appear to be mostly bias. Squaring the biases (bracketed entries) for  $\hat{\mu}_t^c$  and subtracting these from the corresponding MSEs gives results that are very close to the MSEs of  $\hat{\mu}_t^c$ . For example, the MSE of  $\hat{\mu}_t^c$  for target variable 5 is 10.9, its bias is 1.07, and 10.9 - (1.07)<sup>2</sup> = 9.8  $\cong$  MSE for target variable 5 of  $\hat{\mu}_t^c$  (9.7).

The performance of  $\hat{\mu}_t^c$  in all of the tables is only slightly inferior to  $\hat{\mu}_t$ . Other simulations were done on test populations with different amounts of correlation between adjacent months. The results on these populations were similar to those tabled here.

## 6. CONCLUSIONS

The BLUE,  $\hat{\mu}_{t}^{c}$ , is a generalization of the ordinary

regression estimator (Cochran (1977)). The close tie in this methodology between estimating first and second moments parallels the similar tie in the EM-algorithm. This dependency between estimates of first and second moments implies that estimates of variance for the BLUE are produced as a fallout of these procedures.

The potential computional complexity of this BLUE presents few problems. Computer languages like SASPROC MATRIX, SAS IML, GAUSS, and APL can make quick work of the matrix arithmetic. The simulation results in section five cost about \$50.00 per table for 500 replications of the sampling and estimation on an IBM main frame.

For the example presented here the sampling was pps. This allowed the use of the Horvitz—Thompson estimator as a second benchmark for comparison with the BLUE. Although its relative error was small, it still provided a fairly gross upper bound of MSE compared to

the other estimators. The precision of  $\hat{\mu}^c_t$  may be further

improved by selecting a best sample with respect to the distribution of Z. This type of sampling may be dangerous in cases where the distribution of Z is erroneously specified, but the procedures suggested here are to be used when a long history of sampling experience confirms the correct form of this distribution. For reasons beyond statistics, probability sampling will always be a necessary part of applied sampling theory. Therefore, it may be useful to think of the sampling distribution as the unique case of a *known* nonignorable nonresponse mechanism acting on the sampling universe.

#### ing universe. REFERENCES

Cassel C., Sarndal C., and Wretman J.H. (1977), Foundations of Inference in Survey Sampling, John Wiley & Sons.

Cochran W.C. (1977), Sampling Techniques, John Wiley & Sons.

Cox D.R. and Hinkley D.V. (1974), *Theoretical Statistics*, Chapman and Hall.

Little J.A. and Rubin D.B. (1987), *Statistical Analysis with Missing Data*, John Wiley & Sons.

Pfeffermann D. (1988), "The Effect of Sampling Design and Response Mechanism on Multivariate Regression Based Predictors", Journal of the American Statistical Association, 83, 824–833.

Raj D. (1968), Sampling Theory, McGraw-Hill.

Rao C.R. (1973), Linear Statistical Inference and Its Applications, John Wiley & Sons.

Royall R.M. (1981), "Study of the Role of Probability Models in 790 Survey Design and Estimation", Bureau of Labor Statistics contract Report 80–98.

Royall R.M. and Cumberland W.G. (1981a), "An Empirical Study of the Ratio Estimator and Estimators of Its Variance", *Journal of the American Statistical Association*, 76, 66–77.

Royall R.M. and Cumberland W.G. (1981b), "The Finite–Population Linear Regression Estimator and Estimators of its Variance – An Empirical Study", *Journal of the American Statistical Association*, 76, 924–930.

of the American Statistical Association, 76, 924–930. Royall R.M. and Herson J.H. (1973), "Robust Estimation in Finite Populations I", Journal of the American Statistical Association, 68, 880–889.

Rubin D.B. (1987), Multiple Imputation for Nonresponse in Surveys, John Wiley & Sons.

Srivastava M.S. and Carter E.M. (1986), "The Maximum Likelihood Method for Nonresponse in Sample Surveys", *Survey Methodology*, Vol 12, Number 1, 61–72.

Survey Methodology, Vol 12, Number 1, 61–72. West S.A. (1981), "Linear Models for Monthly All Employment Data", Bureau of Labor Statistics report.