

## **Small Area Estimation under Informative Probability Sampling of Areas and Within the Selected Areas**

Danny Pfeffermann, Hebrew University and University of Southampton

Michail Sverchkov, Bureau of Labor Statistics and BAE Systems IT

### ABSTRACT

In this article we show how to predict small area means and obtain valid MSE estimators and confidence intervals when the areas represented in the sample are sampled with unequal probabilities that are possibly related to the true (unknown) area means, and the sampling of units within the selected areas is with probabilities that are possibly related to the outcome values. Ignoring the effects of the sampling process on the distribution of the observed outcomes in such cases may bias the inference very severely. Classical design based inference that uses the randomization distribution of probability weighted estimators cannot be applied for predicting the means of nonsampled areas. We propose simple test statistics for testing the informativeness of the selection of the areas and the sampling of units within the selected areas. The proposed procedures are illustrated by a simulation study and a real application of estimating mean body mass index in counties of the U.S.A, using data from the NHANES III survey.

**Key Words:** Body mass index, Bootstrap, Design based inference, Sample distribution, Sample-complement distribution, Sampling weights.

**Acknowledgement:** Opinions expressed in this paper are of the authors and do not constitute a policy of the Bureau of Labor Statistics. The authors thank the referee for very thoughtful comments that improved the article very significantly and Lester Curtin from the National Center of Health Statistics in the U.S. for providing the data used for the empirical application and his helpful advice.

## 1. INTRODUCTION

The problem of small area estimation is how to predict the area means or other quantities of interest and assess the prediction errors when the sample sizes in these areas are too small (or zero) to warrant the use of direct design-based estimators. It is generally accepted that small area estimation should be based in such cases on statistical models that define ways of borrowing information across areas or over time. See the book by Rao (2003) for a comprehensive account of available methods. However, all the models and estimators considered so far assume either that all the areas are represented in the sample or that the sampled areas are selected with equal probabilities. A few studies consider the case where the sampling of units within the selected areas is with unequal selection probabilities that are related to the outcome values, see, Kott (1990), Arora and Lahiri (1997) and Prasad and Rao (1999), but these studies only treat the case where the input data consist of the direct estimators of the area means. Malec *et al.* (1999) consider unit level observations and use marginal likelihoods and Bayesian methods for inference. We refer to this study in Section 10.

In this article we fill this gap by considering situations where the selection of the areas is with unequal probabilities that are possibly related to the true area means, and the sampling of units within the selected areas is with probabilities that are possibly related to the outcome values, even when conditioning on the model covariates. The problem with this kind of sampling designs is that the model holding for the population values no longer holds for the sample data, giving rise to what is known in the sampling literature as '*informative sampling*'. As illustrated in this article, failure to account for the effects of an informative sampling scheme biases the predictors and increases their root mean square error. For example, the NHANES III survey that is used for the empirical application in Section 10 oversamples minority groups, and if

the target variable of interest (body mass index in our application) is related to ethnicity, then any valid inference procedure should account for the sample selection.

In theory, the effect of the sample selection can be controlled by including among the model covariates all the design variables used for the sample selection. However, this is often not practical either because some or all the design variables may not be known or available at the inference stage, or because there are too many of them, making the fitting and validation of such models formidable. One could attempt to add instead to the model the sampling weights as surrogates for the design variables, but the weights may not summarize the information in the design variables adequately, and this proposition is not operational if the sampling weights are not available for the nonsampled areas or units, which is often the case in a secondary analysis. As mentioned before, direct design based estimators are highly variable in sampled areas because of the small sample sizes, and no design based theory exists for the prediction of the means of nonsampled areas because design based theory uses the randomization distribution of an estimator over repeated sampling from a fixed finite population as the basis for inference. This theory can be used therefore for estimating the population quantities of interest, but not for predicting nonsampled values.

We use relationships between the ‘population distribution’, the ‘sample distribution’ and the ‘sample-complement distribution’ of an outcome variable developed in Pfeffermann and Sverchkov (1999) and Sverchkov and Pfeffermann (2004), in order to derive approximately unbiased predictors of the means in sampled and nonsampled areas under informative sampling of areas and within the areas. We develop estimators for the variances of these predictors and propose simple test statistics for testing the informativeness of the sample selection. The proposed procedures are illustrated by a simulation study and a real application that considers the prediction of mean body mass index (BMI) for counties in the U.S.

Section 2 defines the three distributions and shows the relationships between them. Section 3 defines the optimal predictors in sampled and nonsampled areas and Section 4 shows the bias resulting from ignoring an informative sampling scheme. In Sections 5 and 6 we establish the theory underlying the proposed prediction procedure, with Section 5 showing step by step how to obtain the predictors of the small area means under a particular model identified for the sample data and Section 6 developing appropriate variance estimators. Section 7 extends the theory to general sample models. In Section 8 we present test statistics for testing the informativeness of the sample selection. The simulation results are studied in Section 9, which also examines the performance of confidence intervals for the unknown area means. Section 10 considers the prediction of BMI county means in the U.S. We conclude with a brief summary in Section 11.

## 2. THE SAMPLE AND SAMPLE-COMPLEMENT DISTRIBUTIONS

Consider a finite population of  $N$  units belonging to  $M$  areas, with  $N_i$  units in area  $i$ . Let  $y$  define the target variable with value  $y_{ij}$  for unit  $j$  in area  $i$ , and denote by  $\mathbf{x}_{ij}$  the values of corresponding covariates. In what follows we consider the population  $y$ -values as outcomes of the following two-level random process:

1. *First level values (random effects)*  $\{u_1 \dots u_M\}$  are generated independently from some distribution with probability density function (*pdf*)  $f_p(u_i)$  for which,  $E_p(u_i) = 0$ ,  $E_p(u_i^2) = \sigma_u^2$ , where  $E_p$  defines the expectation operator.
2. *Second level values*  $\{y_{i1} \dots y_{iN_i}\}$  are generated independently from some distribution with *pdf*  $f_p(y_{ij} | \mathbf{x}_{ij}, u_i)$ , for  $i = 1 \dots M$ .

We assume a two-stage sampling design by which in the first stage  $m$  areas are selected with probabilities  $\pi_i = \Pr(i \in s)$ , and in the second stage  $n_i$  units are sampled

from area  $i$  selected in the first stage with probabilities  $\pi_{j|i} = \Pr(j \in s_i | i \in s)$ . Note that the sample inclusion probabilities at both stages may depend in general on all the population or area values of  $y$  and  $\mathbf{x}$ , and possibly also on the population values of design variables  $\mathbf{z}$  used for the sample selection. Denote by  $I_i$  and  $I_{ij}$  the sample indicator variables for the two sampling stages ( $I_i = 1$  iff  $i \in s$  and similarly for  $I_{ij}$ ), and by  $w_i = 1/\pi_i$  and  $w_{j|i} = 1/\pi_{j|i}$  the first and second stage sampling weights.

Following Pfeffermann *et. al* (1998), we define the conditional first level *sample pdf* of  $u_i$ , that is, the *pdf* of  $u_i$  for area  $i \in s$  as,

$$f_s(u_i) \stackrel{\text{def}}{=} f(u_i | I_i = 1) = \Pr(I_i = 1 | u_i) f_p(u_i) / \Pr(I_i = 1). \quad (2.1)$$

The conditional first level *sample-complement pdf* of  $u_i$ , that is, the *pdf* of  $u_i$  for area  $i \notin s$  is defined in Sverchkov and Pfeffermann (2004) as,

$$f_c(u_i) \stackrel{\text{def}}{=} f(u_i | I_i = 0) = \Pr(I_i = 0 | u_i) f_p(u_i) / \Pr(I_i = 0). \quad (2.2)$$

Note that the *population*, *sample* and *sample-complement pdfs* of  $u_i$  are the same if,  $\Pr(I_i = 1 | u_i) = \Pr(I_i = 1) \forall i$ , in which case the area selection is *noninformative*.

The conditional second level *sample pdf* and *sample-complement pdfs* of  $y_{ij}$  in a sampled area are defined similarly to (2.1) and (2.2) as,

$$\begin{aligned} f_{si}(y_{ij} | \mathbf{x}_{ij}, u_i, I_i = 1) &\stackrel{\text{def}}{=} f(y_{ij} | \mathbf{x}_{ij}, u_i, I_i = 1, I_{ij} = 1) \\ &= \frac{\Pr(I_{ij} = 1 | y_{ij}, \mathbf{x}_{ij}, u_i, I_i = 1) f_p(y_{ij} | \mathbf{x}_{ij}, u_i, I_i = 1)}{\Pr(I_{ij} = 1 | \mathbf{x}_{ij}, u_i, I_i = 1)}, \end{aligned} \quad (2.3)$$

$$\begin{aligned} f_{ci}(y_{ij} | \mathbf{x}_{ij}, u_i, I_i = 1) &\stackrel{\text{def}}{=} f(y_{ij} | \mathbf{x}_{ij}, u_i, I_i = 1, I_{ij} = 0) \\ &= \frac{\Pr(I_{ij} = 0 | y_{ij}, \mathbf{x}_{ij}, u_i, I_i = 1) f_p(y_{ij} | \mathbf{x}_{ij}, u_i, I_i = 1)}{\Pr(I_{ij} = 0 | \mathbf{x}_{ij}, u_i, I_i = 1)}. \end{aligned} \quad (2.4)$$

Here again the *population*, *sample* and *sample-complement pdfs* are the same if,  $\Pr(\mathbf{I}_{ij} = 1 | y_{ij}, \mathbf{x}_{ij}, u_i, \mathbf{I}_i = 1) = \Pr(\mathbf{I}_{ij} = 1 | \mathbf{x}_{ij}, u_i, \mathbf{I}_i = 1) \forall j$ . The model defined by (2.1) and (2.3) defines the two-level sample model that corresponds to the population model defined by  $f_p(u_i)$  and  $f_p(y_{ij} | \mathbf{x}_{ij}, u_i)$ ; see also Pfeffermann *et al.* (2006).

The following relationships between the population *pdf*, the sample *pdf* and the sample-complement *pdf* are established in Pfeffermann and Sverchkov (1999) and Sverchkov and Pfeffermann (2004) for general pairs of random variables  $\mathbf{v}_1, \mathbf{v}_2$  measured for elements  $i$  of a population P. The symbols  $E_p, E_s$  and  $E_c$  define respectively the expectations under the three distributions and  $\{\pi_i, w_i\}$  denotes the sample inclusion probabilities and the corresponding sampling weights  $w_i = 1/\pi_i$ .

$$f_s(\mathbf{v}_{1i} | \mathbf{v}_{2i}) = f(\mathbf{v}_{1i} | \mathbf{v}_{2i}, i \in s) = E_p(\pi_i | \mathbf{v}_{1i}, \mathbf{v}_{2i}) f_p(\mathbf{v}_{1i} | \mathbf{v}_{2i}) / E_p(\pi_i | \mathbf{v}_{2i}), \quad (2.5)$$

$$E_p(\mathbf{v}_{1i} | \mathbf{v}_{2i}) = E_s(w_i \mathbf{v}_{1i} | \mathbf{v}_{2i}) / E_s(w_i | \mathbf{v}_{2i}) \quad ; \quad E_p(\pi_i | \mathbf{v}_{2i}) = 1 / E_s(w_i | \mathbf{v}_{2i}), \quad (2.6)$$

$$\begin{aligned} f_c(\mathbf{v}_{1i} | \mathbf{v}_{2i}) &= f(\mathbf{v}_{1i} | \mathbf{v}_{2i}, i \notin s) = \frac{E_p[(1 - \pi_i) | \mathbf{v}_{1i}, \mathbf{v}_{2i}] f_p(\mathbf{v}_{1i} | \mathbf{v}_{2i})}{E_p[(1 - \pi_i) | \mathbf{v}_{2i}]} \\ &= \frac{E_s[(w_i - 1) | \mathbf{v}_{1i}, \mathbf{v}_{2i}] f_s(\mathbf{v}_{1i} | \mathbf{v}_{2i})}{E_s[(w_i - 1) | \mathbf{v}_{2i}]} \quad , \end{aligned} \quad (2.7)$$

$$E_c(\mathbf{v}_{1i} | \mathbf{v}_{2i}) = \frac{E_p[(1 - \pi_i) \mathbf{v}_{1i} | \mathbf{v}_{2i}]}{E_p[(1 - \pi_i) | \mathbf{v}_{2i}]} = \frac{E_s[(w_i - 1) \mathbf{v}_{1i} | \mathbf{v}_{2i}]}{E_s[(w_i - 1) | \mathbf{v}_{2i}]} \quad (2.8)$$

Defining  $\mathbf{v}_{1i} = u_i, \mathbf{v}_{2i} = \text{const}$  yields the relationships holding for the random area effects  $u_i$ . Defining  $\mathbf{v}_{1ij} = y_{ij}; \mathbf{v}_{2ij} = (\mathbf{x}_{ij}, u_i, \mathbf{I}_i = 1)$  and substituting  $\pi_{j|i}$  and  $w_{j|i}$  for  $\pi_i$  and  $w_i$  respectively, yields the relationships holding for the observations  $y_{ij}$ .

### 3. OPTIMAL SMALL AREA PREDICTORS

The target population parameters are the small area means  $\bar{Y}_i = \sum_{j=1}^{N_i} y_{ij} / N_i$  for  $i = 1 \dots M$ , (the means in sampled and nonsampled areas). Let  $D_s = \{(y_{ij}, w_{j|i}, w_i),$

$(i, j) \in s; \mathbf{x}_{kl}, (k, l) \in U\}$  define the known data. Note that we do not assume knowledge of the sampling weights of nonsampled units or areas. The MSE of a predictor  $\hat{Y}_i$  with respect to the *population pdf*, given  $D_s$  and  $I_i$  is,

$$MSE(\hat{Y}_i | D_s, I_i) = E_p[(\hat{Y}_i - \bar{Y}_i)^2 | D_s, I_i] = [\hat{Y}_i - E_p(\bar{Y}_i | D_s, I_i)]^2 + V_p(\bar{Y}_i | D_s, I_i). \quad (3.1)$$

The variance  $V_p(\bar{Y}_i | D_s, I_i)$  does not depend on the form of the predictor and hence the MSE is minimized when  $\hat{Y}_i = E_p(\bar{Y}_i | D_s, I_i)$ .

In what follows we make the following mild assumption (see Remark 1 below):

*Ass.1-*  $f_{ci}(y_{il} | D_s, u_i, I_i = 1) = f_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1)$ , implying that unobserved outcomes in a sampled area are independent of the observed outcomes and their sampling weights when conditioning on the area random effect and the covariates.

Ass.1 is satisfied under the following two conditions:

$$(a) f(y_{il}, y_{ij} | u_i, \mathbf{x}_{il}, \mathbf{x}_{ij}, I_{il} = 0, I_{ij} = 1) = f(y_{il} | u_i, \mathbf{x}_{il}, I_{il} = 0) f(y_{ij} | u_i, \mathbf{x}_{ij}, I_{ij} = 1),$$

$$(b) f(\pi_{ji} | u_i, y_{il}, y_{ij}, \mathbf{x}_{il}, \mathbf{x}_{ij}, I_{il} = 0, I_{ij} = 1) = f(\pi_{ji} | u_i, y_{ij}, \mathbf{x}_{ij}, I_{ij} = 1).$$

The first condition is very mild since the population outcomes are independent given the random effect, and the area selection probability under informative sampling is related to the area mean  $E(\bar{Y}_i | u_i, \mathbf{X}_i)$ , where  $\mathbf{X}_i$  denotes the area covariates, and not to individual deviations from the mean, such that by conditioning on the random effect the independence of the outcomes is preserved. The second condition also seems mild for the common situation in small area estimation of large areas and small samples.

**Remark 1.** In the paper we only use the weaker assumption,  $E_{ci}(y_{ik} | D_s, u_i, I_i = 0) = E_{ci}(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 0)$ . Note also that  $E_{si}[E_{ci}(y_{ik} | D_s, u_i, I_i = 0) | \mathbf{x}_{ik}, u_i, I_i = 0] = E_{ci}(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 0)$ , such that the second expectation ‘predicts’ the first even if the weaker assumption is not satisfied.

As shown above, the optimal predictor for a given area  $i$  is,  $\hat{Y}_i = E_p(\bar{Y}_i | D_s, I_i)$ .

If the area is sampled ( $I_i = 1$ ), then by (2.7) and Ass.1,

$$\begin{aligned} E_p(\bar{Y}_i | D_s, I_i = 1) &= \frac{1}{N_i} E_p \{ [\sum_{l=1}^{N_i} E_p(y_{il} | D_s, u_i, I_i = 1)] | D_s, I_i = 1 \} \\ &= \frac{1}{N_i} \{ \sum_{j \in s_i} y_{ij} + \sum_{l \notin s_i} E_s [E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s] \} \end{aligned} \quad (3.2)$$

For a nonsampled area ( $I_i = 0$ ), if  $E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 0) = E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 1)$  (see below), then,

$$\begin{aligned} E_p(\bar{Y}_i | D_s, I_i = 0) &= \frac{1}{N_i} \sum_{k=1}^{N_i} E_p [E_p(y_{ik} | D_s, u_i, I_i = 0) | D_s, I_i = 0] \\ &= \frac{1}{N_i} \sum_{k=1}^{N_i} E_c [E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 0) | D_s] = \frac{1}{N_i} \sum_{k=1}^{N_i} E_c [E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 1) | D_s], \end{aligned} \quad (3.3)$$

with the first expression on the second line following from the fact that the outcomes  $y_{ik}$  are in a nonsampled area. The condition  $E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 0) = E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 1)$  is not restrictive since as discussed with regard to Ass.1, the area selection probabilities are related to the area mean and are not dependent on individual deviations from the mean.

#### 4. BIAS OF SMALL AREA PREDICTORS WHEN IGNORING AN INFORMATIVE SAMPLING SCHEME

Consider first a sampled area. Ignoring the sampling scheme within a selected area implies an implicit assumption that the *sample-complement* model in the area is the same as the *sample* model such that,

$$\hat{Y}_{i,IGN} = \frac{1}{N_i} \{ \sum_{j \in s_i} y_{ij} + \sum_{l \notin s_i} E_s [E_{si}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s] \} \text{ (compare with 3.2). Hence,}$$

$$\begin{aligned} Bias(\hat{Y}_{i,IGN}) &= E_p [(\hat{Y}_{i,IGN} - \bar{Y}_i) | D_s, I_i = 1] \\ &= \frac{1}{N_i} \sum_{l \notin s_i} E_s [E_{si}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s] - \frac{1}{N_i} \sum_{l \notin s_i} E_s [E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s] \end{aligned}$$



$$= -\frac{1}{N_i} E_s \left[ \sum_{l \in s_i} \frac{\text{Cov}_{si}(y_{il}, w_{li} | \mathbf{x}_{il}, u_i, \mathbf{I}_i = 1)}{E_{si}[(w_{li} - 1) | \mathbf{x}_{il}, u_i, \mathbf{I}_i = 1]} \mid D_s \right], \quad (4.1)$$

with the last equality following from (2.8). Thus, if the outcomes  $y_{il}$  and the sampling weights  $w_{li}$  are correlated given the covariates and the random effect, ignoring the sampling scheme yields biased predictors.

Next consider a non-sampled areas. By (3.3),

$$\begin{aligned} \text{Bias}(\hat{Y}_{i,IGN}) &= E_p[(\hat{Y}_{i,IGN} - \bar{Y}_i) \mid D_s, \mathbf{I}_i = 0] \\ &= \frac{1}{N_i} \sum_{k=1}^{N_i} E_s[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, \mathbf{I}_{ik} = 1) \mid D_s] - \frac{1}{N_i} \sum_{k=1}^{N_i} E_c[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, \mathbf{I}_i = 1) \mid D_s]. \end{aligned} \quad (4.2)$$

Adding and subtracting  $\frac{1}{N_i} \sum_{k=1}^{N_i} E_c[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, \mathbf{I}_{ik} = 1) \mid D_s]$  and applying (2.8)

and (2.6) yields,

$$\begin{aligned} \text{Bias}(\hat{Y}_{i,IGN}) &= -\frac{1}{N_i} \sum_{k=1}^{N_i} \frac{\text{Cov}_s(E_p(y_{ik} | \mathbf{x}_{ik}, u_i, \mathbf{I}_{ik} = 1), w_i \mid D_s)}{E_s[(w_i - 1) \mid D_s]} \\ &\quad - \frac{1}{N_i} E_c \left[ \sum_{k=1}^{N_i} \frac{\text{Cov}_{si}(y_{ik}, w_{ki} | \mathbf{x}_{ik}, u_i, \mathbf{I}_i = 1)}{E_{si}[w_{ki} | \mathbf{x}_{ik}, u_i, \mathbf{I}_i = 1]} \mid D_s \right]. \end{aligned} \quad (4.3)$$

The first covariance reflects the bias induced by the informative selection of areas. The second covariance reflects the bias induced by the informative sampling within the selected areas (compare with 4.1). In Section 8 we propose simple tests for testing whether the covariances in (4.1) and (4.3) are zero, such that ignoring the sample selection does not bias the predictors. See also the simulation results in Section 9.

## 5. PREDICTION OF SMALL AREA MEANS

Our approach requires specifying the two-level sample model,  $f_s(u_i) = f(u_i | \mathbf{I}_i = 1)$  (Eq. 2.1) and  $f_{si}(y_{il} | \mathbf{x}_{il}, u_i, \mathbf{I}_i = 1)$  (Eq. 2.3), and the conditional sample expectations,  $E_{si}(w_{ji} | \mathbf{x}_{ij}, y_{ij}, u_i, \mathbf{I}_i = 1)$ , all of which can be identified and tested using the observed data since they refer to sample models. We do not assume any model for the population data or the unobserved data, and make no

assumptions regarding the selection of areas or the model holding for the area sampling weights. In order to facilitate the presentation of our approach, we consider in Sections 5 and 6 a particular sample model and sampling scheme within the selected areas. In Section 7 we outline the basic steps of computing the predictors under a general model fitted to the sample data with continuous or discrete outcomes and fixed and random effects, and a general sampling scheme within the selected areas.

The first step of our approach is therefore to fit a model to the sample data, which of course is a necessary step in any small area estimation application. Note that although we consider informative sampling, the sample model can be identified and estimated from the sample data using standard techniques, see Rao (2003) for small area model identification and diagnostic methods. In this and the next section we assume that the sample model for the outcome values is the ‘nested error regression model’,

$$y_{ij} = \mathbf{x}'_{ij}\boldsymbol{\beta} + u_i + e_{ij}; \quad u_i | I_i = 1 \overset{ind}{\sim} N(0, \sigma_u^2), \quad e_{ij} | I_{ij} = 1 \overset{ind}{\sim} N(0, \sigma_e^2). \quad (5.1)$$

Suppose that the sampling weights,  $w_{j|i} = (1/\pi_{j|i})$  within the selected areas satisfy,

$$E_{st}(w_{j|i} | \mathbf{x}_{ij}, y_{ij}, u_i, I_i = 1) = E_{st}(w_{j|i} | \mathbf{x}_{ij}, y_{ij}, I_i = 1) = k_i \exp(\mathbf{x}'_{ij}\mathbf{a} + by_{ij}), \quad (5.2)$$

where  $k_i = N_i n_i^{-1} \sum_{j=1}^{N_i} \exp(-\mathbf{x}'_{ij}\mathbf{a} - by_{ij}) / N_i$  (follows from (2.6)), and  $\mathbf{a}$  and  $b$  are fixed (unknown) constants. (If  $x_{i0} = const$ , we assume  $a_0 = 1$  for uniqueness). Note that for large areas,  $\sum_{j=1}^{N_i} \exp(-\mathbf{x}'_{ij}\mathbf{a} - by_{ij}) / N_i \cong E_\rho[\sum_{j=1}^{N_i} \exp(-\mathbf{x}'_{ij}\mathbf{a} - by_{ij}) / N_i] = const$ , such that  $k_i \cong (N_i / n_i) \times const$ . As becomes evident below, for sufficiently small sampling fractions the predictors for sampled and nonsampled areas do not depend on  $\mathbf{a}$  and  $k_i$ .

Remark 2: It follows from Pfeffermann *et al.* (1998) that under the sampling scheme (5.2) the population model is also of the form (5.1) but with different parameters, if the areas are selected with probabilities  $\pi_i$  satisfying  $E(\pi_i | \theta_i) \propto \exp[\gamma_0 \theta_i + \mathbf{z}'_i \boldsymbol{\gamma}]$ , where

$\theta_i = \bar{\mathbf{X}}_i' \boldsymbol{\beta} + u_i$  are the area means,  $\mathbf{z}_i$  represents area level design variables and  $(\gamma_0, \boldsymbol{\gamma})$  are fixed coefficients. The model (5.1) is in common use for small area estimation under noninformative sampling (in which case the population and sample models coincide), see, e.g., Battese *et al.* (1988). However, as emphasized above, we only assume knowledge of the forms of the sample model (5.1) and the conditional expectations in (5.2), but not the form of the population model or the relationship between the area selection probabilities and the area means.

Remark 3: As with the sample model (5.1), the expectation in (5.2) refers to the sample distribution within the sampled areas. The relationship in the sample between the sampling weights and the outcomes can be identified and estimated therefore from the sample data, see Skinner (1994) and Pfeffermann and Sverchkov (1999, 2003) for discussion and examples. On the other hand, the relationship between the sampling weights  $w_i$  and the area means is more difficult to identify since the area means are not observable, and we do not model this relationship. See Pfeffermann *et al.* (2006) for examples of modeling the area selection probabilities. Kim (2003) assumes the model (5.1) for the population values and a similar model to (5.2) for the sampling probabilities within the areas, but his article assumes implicitly that all the population areas are sampled.

The analysis that follows assumes known model parameters. In practice, the unknown model parameters are replaced under the frequentist approach by sample estimates, yielding the corresponding ‘empirical predictors’. Maximum likelihood estimation of the model parameters has to be based in the present case on the sample distribution of the sample outcomes. Alternatively, the model parameters can possibly be estimated by the ‘method of moments’, depending on the underlying model. See the empirical study in Sections 9 and 10.

As established in Section 3, the optimal predictor for a sampled area  $i$  is,  $E_p(\bar{Y}_i | D_s, I_i = 1) = \{\sum_{j \in S_i} y_{ij} + \sum_{l \notin S_i} E_s[E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s]\} / N_i$ . In order to compute  $E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1)$  note that by (2.7), (5.1) and (5.2),

$$\begin{aligned} f_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) &= \frac{[E_{si}(w_{li} | \mathbf{x}_{il}, y_{il}, u_i, I_i = 1) - 1] f_{si}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1)}{E_{si}(w_{li} | \mathbf{x}_{il}, u_i, I_i = 1) - 1} \\ &= \frac{\lambda_{il}}{\lambda_{il} - 1} \frac{1}{\sigma_e} \phi\left(\frac{y_{ij} - u_{il} - b\sigma_e^2}{\sigma_e}\right) - \frac{1}{\lambda_{il} - 1} \frac{1}{\sigma_e} \phi\left(\frac{y_{il} - u_{il}}{\sigma_e}\right), \end{aligned} \quad (5.3)$$

where  $u_{il} = \mathbf{x}'_{il}\boldsymbol{\beta} + u_i$ ,  $\lambda_{il} = k_i \exp[(b^2\sigma_e^2/2) + \mathbf{x}'_{il}\mathbf{a} + bu_{il}] = E_{si}(w_{li} | \mathbf{x}_{il}, u_i, I_i = 1)$  and  $\phi$  is the standard normal *pdf*. In the special case where  $b=0$  (the selection probabilities within the sampled areas only depend on the  $\mathbf{x}$ -values so that the sampling is noninformative), the *pdf* in (5.3) reduces to the sample normal density (5.1). By computing the expectation under the sample-complement *pdf* (5.3) we find,

$$E_s[E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s] = E_s\left[u_{il} + \frac{\lambda_{il}}{\lambda_{il} - 1} b\sigma_e^2 \mid D_s\right]. \quad (5.4)$$

The expectation in the right hand side of (5.4) is with respect to the sample distribution of  $u_i | D_s, I_i = 1$ . Under the sample model (5.1),

$$u_i | D_s, I_i = 1 \sim N(\hat{u}_i, \sigma_i^2 \gamma_i), \quad (5.5)$$

where  $\hat{u}_i = \gamma_i[\bar{y}_i - \bar{\mathbf{x}}'_i\boldsymbol{\beta}]$ ;  $(\bar{y}_i, \bar{\mathbf{x}}_i) = \sum_{j=1}^{n_i} (y_{ij}, \mathbf{x}_{ij}) / n_i$  are the sample means of  $(y, \mathbf{x})$  in sampled area  $i$ ,  $\gamma_i = \sigma_u^2 / [\sigma_u^2 + \sigma_i^2]$  and  $\sigma_i^2 = \sigma_e^2 / n_i = \text{Var}_s(\bar{y}_i | u_i)$ . Thus, the expectation  $E_s[E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s]$  is obtained by computing the expectation in the right hand side of (5.4) with respect to the normal distribution (5.5). We get,

$$E_s[E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s] = (\mathbf{x}'_{il}\boldsymbol{\beta} + \hat{u}_i) + b\sigma_e^2 E_s[(1 - \lambda_{il}^{-1})^{-1} | D_s]. \quad (5.6)$$

Note that if  $b=0$  (noninformative sampling within the area),

$$E_s[E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s] = \mathbf{x}'_{il}\boldsymbol{\beta} + \hat{u}_i, \text{ which is the standard result.}$$

The expectation  $E_s[(1-\lambda_{il}^{-1})^{-1} | D_s]$  can be computed numerically. Alternatively, in the practical case where the sampling fractions within the selected areas are small,  $\lambda_{il} = E_s(w_{li} | \mathbf{x}_{il}, u_i, I_i=1)$  is typically much larger than 1 and hence we may approximate,  $E_s[(1-\lambda_{il}^{-1})^{-1} | D_s] \cong 1$ , in which case by (5.6),

$$E_s[E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i=1) | D_s] \cong \hat{u}_{il} + b\sigma_e^2; \quad \hat{u}_{il} = \mathbf{x}'_{il}\boldsymbol{\beta} + \hat{u}_i. \quad (5.7)$$

It follows from (3.2), (5.6) and (5.7) that for given parameters  $\{\boldsymbol{\beta}, b, \sigma_u^2, \sigma_e^2\}$ , the mean  $\bar{Y}_i$  of sampled area  $i$  can be predicted as,

$$E_p(\bar{Y}_i | D_s, I_i=1) = N_i^{-1} \{(N_i - n_i)\hat{\theta}_i + n_i[\bar{y}_i + (\bar{\mathbf{X}}_i - \bar{\mathbf{x}}_i)'\boldsymbol{\beta}] + (N_i - n_i)b\sigma_e^2\}, \quad (5.8)$$

where  $\hat{\theta}_i = \hat{u}_i + \bar{\mathbf{X}}_i'\boldsymbol{\beta}$  is the optimal predictor of the sample model mean  $\theta_i = \bar{\mathbf{X}}_i'\boldsymbol{\beta} + u_i = E_{si}(\bar{Y}_i | \mathbf{X}_i, u_i)$ . The last term in (5.7) corrects for the sample selection effects, that is, the difference between the sample-complement expectation and the sample expectation in sampled areas. For  $b=0$ , the predictor (5.8) reduces to the optimal predictor under noninformative sampling (Rao, 2003, Eq. 7.2.37). Note that the predictor (5.8) does not depend on  $K_i$  and  $\mathbf{a}$  featuring in the expectation (5.2).

The optimal predictor for *nonsampled* areas is defined in (3.3) to be,

$$\begin{aligned} E_p(\bar{Y}_i | D_s, I_i=0) &= \frac{1}{N_i} \sum_{k=1}^{N_i} E_c[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i=1) | D_s]. \text{ By (2.8) and (2.6),} \\ E_c[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i=1) | D_s] &= E_s\left[\frac{(w_i - 1)E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i=1)}{E_s(w_i | D_s) - 1} \mid D_s\right] \\ &= \{E_s[(w_i - 1) \frac{E_{si}(w_{k|i}y_{ik} | \mathbf{x}_{ik}, u_i, I_i=1)}{E_{si}(w_{k|i} | \mathbf{x}_{ik}, u_i, I_i=1)} \mid D_s]\} / [E_s(w_i | D_s) - 1]. \end{aligned} \quad (5.9)$$

Computing the expectations  $E_{si}(\cdot)$  in the numerator and the denominator, using (5.1) and (5.2) yields after some algebra,

$$E_c[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i=1) | D_s] = \mathbf{x}'_{ik}\boldsymbol{\beta} + b\sigma_e^2 + E_s\left[\frac{(w_i - 1)u_i}{E_s(w_i | D_s) - 1} \mid D_s\right]. \quad (5.10)$$

Estimating the two sample expectations in the right hand side of (5.10) by the corresponding sample means and substituting  $\hat{u}_i = \gamma_i[\bar{y}_i - \bar{\mathbf{x}}_i' \boldsymbol{\beta}]$  for  $u_i$  yields the following estimator for  $E_{ik} = E_c[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, \mathbf{I}_i = 1) | D_s]$ ,

$$\hat{E}_{ik} = \mathbf{x}'_{ik} \boldsymbol{\beta} + b\sigma_e^2 + \sum_{i \in s} (w_i - 1) \hat{u}_i / \sum_{i \in s} (w_i - 1). \quad (5.11)$$

It follows from (3.3) and (5.11) that for given parameters  $\{\boldsymbol{\beta}, b, \sigma_u^2, \sigma_e^2\}$ , the mean  $\bar{Y}_i$  of area  $i$  not in the sample can be predicted as,

$$\hat{E}_p(\bar{Y}_i | D_s, \mathbf{I}_i = 0) = \bar{\mathbf{X}}'_i \boldsymbol{\beta} + b\sigma_e^2 + [\sum_{i \in s} (w_i - 1) \hat{u}_i / \sum_{i \in s} (w_i - 1)]. \quad (5.12)$$

The last term of (5.12) corrects for the fact that the mean of the random effects for areas outside the sample is different from zero under informative selection of areas.

## 6. MSE ESTIMATION

Estimating  $MSE(\hat{Y}_i | D_s, \mathbf{I}_i) = E_p[(\hat{Y}_i - \bar{Y}_i)^2 | D_s, \mathbf{I}_i]$  for the predictors considered in section 5 requires strict model assumptions that could be hard to validate. In order to deal with this problem, we estimate instead,  $MSE(\hat{Y}_i | \mathbf{X}, \mathbf{I}) = E_p[MSE(\hat{Y}_i | D_s, \mathbf{I}_i) | \mathbf{X}, \mathbf{I}]$ , where  $\mathbf{X} = \{\mathbf{x}_{ij}, (i, j) \in U\}$  and  $\mathbf{I} = \{\mathbf{I}_i, \mathbf{I}_{ij}, (i, j) \in U\}$  is the set of first and second stage sample indicators. Note that  $MSE(\hat{Y}_i | D_s, \mathbf{I}_i)$  can be viewed as random, such that  $MSE(\hat{Y}_i | \mathbf{X}, \mathbf{I})$  defines its ‘best predictor’ under the mean square loss function over the distribution  $f_{D_s | \mathbf{X}, \mathbf{I}}$ .

Denote by  $\hat{Y}_i$  the predictor defined by (5.8) if  $i \in s$  or by (5.12) if  $i \notin s$ . For what follows in this section only, we make the following additional mild assumptions:

$$Ass.2 \quad Cov_p[y_{ij}, y_{mk} | \mathbf{X}, \mathbf{I}_i = 1, \mathbf{I}_m = 0] = 0; \quad Cov_p[y_{ij}, y_{ik} | \mathbf{X}, u_i, \mathbf{I}_i = 1, \mathbf{I}_{ij} = \mathbf{I}_{ik} = 0] = 0;$$

implying that outcomes in sampled areas are uncorrelated with outcomes in nonsampled areas, and that the unobserved outcomes in a sampled area are

uncorrelated when conditioning on the random effect. The first assumption generally holds if the random effects are independent between the areas. The second assumption is also not restrictive because the population outcomes are conditionally independent given the random effect and by extending Remark 2 of Sverchkov and Pfeffermann (2004) to the case of a joint distribution for a pair of units, it follows that for small sampling fractions the joint sample-complement distribution and the population distribution are approximately the same.

*Ass.3*  $\text{Cov}_p[y_{ij}, y_{ik} | \mathbf{X}, u_i, \mathbf{I}_i = 0] = 0$ ; implying that the outcomes in a nonsampled area are uncorrelated conditionally on the random effect. This assumption holds if the area selection probability only depends on the area mean. (See the discussion below Ass.1)

*Ass.4* The predictor  $\hat{Y}_i$ ,  $i \notin s$  is approximately unbiased for  $E_p(\bar{Y}_i | \mathbf{X}, \mathbf{I}_i = 0)$  in the sense that,  $E_p(\hat{Y}_i | \mathbf{X}, \mathbf{I}) \cong E_p(\bar{Y}_i | \mathbf{X}, \mathbf{I}_i = 0)$  (follows from Section 5).

Consider first *sampled* areas. Denote  $Y_{Ri} = Y_i - \sum_{j \in s_i} y_{ij}$  where  $Y_i = N_i \bar{Y}_i$ , such that

$\hat{Y}_{Ri} = N_i \hat{Y}_i - \sum_{j \in s_i} y_{ij}$ . It is shown in Appendix A that under *Ass.1* and *Ass.2*,

$$E_p[(\hat{Y}_i - Y_i)^2 | \mathbf{X}, \mathbf{I}] = E_p[G(u_i, D_s) | \mathbf{X}, \mathbf{I}]; \quad (6.1)$$

$$G(u_i, D_s) = [\hat{Y}_{Ri} - \sum_{l \notin s_i} (u_{il} + \frac{\lambda_{il}}{\lambda_{il} - 1} b \sigma_\varepsilon^2)]^2 + \sum_{l \notin s_i} (\sigma_\varepsilon^2 - \frac{\lambda_{il} b^2 \sigma_\varepsilon^4}{(\lambda_{il} - 1)^2}), \quad (6.2)$$

where, as before,  $u_{il} = \mathbf{x}'_{il} \boldsymbol{\beta} + u_i$  and  $\lambda_{il} = k_i \exp[(b^2 \sigma_\varepsilon^2 / 2) + \mathbf{x}'_{il} \mathbf{a} + b u_{il}]$ . Note that the expectation in the right hand side of (6.1) is with respect to the *pdf*  $f(u_i, D_s | \mathbf{X}, \mathbf{I})$ . All the terms in (6.2) are either fixed values or functions of the data  $D_s$  and the random effect  $u_i$ . It follows therefore that for a sampled area,  $MSE(\hat{Y}_i | \mathbf{X}, \mathbf{I}, \mathbf{I}_i = 1)$

$= E_p[G(u_i, D_s) | \mathbf{X}, \mathbf{I}] / N_i^2$  can be estimated by the following parametric bootstrap procedure (see Remark 4 below):

1. Estimate  $\mathbf{a}, b, k_i, \boldsymbol{\beta}, \sigma_u^2, \sigma_e^2$  (see Section 9),
2. Generate  $B$  bootstrap samples  $\{u_i^b, y_{ij}^b\}$ ,  $i = 1, \dots, m$ ,  $j = 1, \dots, n_i$  from the sample model (5.1) with parameters  $\hat{\boldsymbol{\beta}}, \hat{\sigma}_u^2, \hat{\sigma}_e^2$ , using the covariates  $\mathbf{x}_{ij}$ ,  $(i, j) \in s$ . Compute,  $w_{ji}^b = \hat{k}_i \exp(\hat{\mathbf{x}}_{ij}' \mathbf{a} + \hat{b} y_{ij}^b)$ .
3. Re-compute the predictors  $\hat{Y}_{Ri}^b$   $i = 1 \dots m$ ,  $b = 1, \dots, B$  (with new parameter estimates) and compute  $G^b(u_i^b, D_s^b)$ ; the new parameter estimates are only used for computing  $\hat{Y}_{Ri}^b$ , the other terms of  $G^b(u_i^b, D_s^b)$  use the original parameter estimates.
4. Estimate,

$$M\hat{S}E(\hat{\bar{Y}}_i | \mathbf{X}, \mathbf{I}, I_i = 1) = \sum_{b=1}^B G^b(u_i^b, D_s^b) / N_i^2 / B \quad (6.3)$$

**Remark 4:** The estimator (6.3) ignores the contribution to the variance from estimating the hyperparameters  $\{\boldsymbol{\beta}, k_i, \mathbf{a}, b, \sigma_u^2, \sigma_e^2\}$ . Accounting for this extra source of variation requires a ‘double bootstrap’ procedure. See Hall and Maiti (2006) for bootstrap bias corrections in small area estimation that warrant MSE estimation of order  $O(1/m^2)$ .

Next consider *nonsampled* areas. Under *Ass.2* and *Ass.4*,

$$\begin{aligned} E_p\{(\hat{\bar{Y}}_i - \bar{Y}_i)^2 | \mathbf{X}, \mathbf{I}\} &\cong E_p\{[\hat{\bar{Y}}_i - E_p(\hat{\bar{Y}}_i | \mathbf{X}, \mathbf{I})]^2 | \mathbf{X}, \mathbf{I}\} + E_p\{[\bar{Y}_i - E_p(\bar{Y}_i | \mathbf{X}, \mathbf{I})]^2 | \mathbf{X}, \mathbf{I}\} \\ &= Var_p(\hat{\bar{Y}}_i | \mathbf{X}, \mathbf{I}) + Var_p(\bar{Y}_i | \mathbf{X}, \mathbf{I}) \end{aligned} \quad (6.4)$$

The first variance in (6.4) can be estimated similarly to the estimation of  $MSE(\hat{\bar{Y}}_i | \mathbf{X}, \mathbf{I}, I_i = 1)$  in (6.3), that is, by applying the first 3 steps of the bootstrap procedure to obtain realizations  $\hat{Y}_i^b$ ,  $i \notin s$ , and then estimating,

$$\hat{V}ar_p(\hat{\bar{Y}}_i | \mathbf{X}, \mathbf{I}) = \frac{1}{B} \sum_{b=1}^B (\hat{Y}_i^b - \hat{Y}_{i,A})^2; \quad \hat{Y}_{i,A} = \frac{1}{B} \sum_{b=1}^B \hat{Y}_i^b. \quad (6.5)$$



Note that  $\hat{Y}_i$ , like  $G(u_i, D_s)$  in (6.2) only uses the sample data and hence it is only needed to generate data from the sample model.

In order to estimate the second variance in (6.4) we use the decomposition,

$$Var_p(\bar{Y}_i | \mathbf{X}, \mathbf{I}) = Var_p[E_p(\bar{Y}_i | u_i, \mathbf{X}, \mathbf{I}) | \mathbf{X}, \mathbf{I}] + E_p[Var_p(\bar{Y}_i | u_i, \mathbf{X}, \mathbf{I}) | \mathbf{X}, \mathbf{I}]. \quad (6.6)$$

By Ass.3, the second component in (6.6) is simply,

$$E_p[Var_p(\bar{Y}_i | u_i, \mathbf{X}, \mathbf{I}) | \mathbf{X}, \mathbf{I}] = \sigma_e^2 / N_i. \quad (6.7)$$

This result follows from a result in Pfeffermann *et al.* (1998), implying that under the sample model (5.1) and (5.2),  $E_p(e_{ij} | u_i, \mathbf{X}, I_i = 0) = const$  and  $Var_p(y_{ij} | u_i, \mathbf{x}_{ij}) = Var_{s_i}(y_{ij} | u_i, \mathbf{x}_{ij}) = \sigma_e^2$ .

Next consider the first term of (6.6). Again, under the sample model (5.1) - (5.2),

$$Var_p[E_p(\bar{Y}_i | u_i, \mathbf{X}, I_i = 0) | \mathbf{X}, \mathbf{I}] = Var_p(u_i | \mathbf{X}, I_i = 0). \quad (6.8)$$

It is shown in Appendix B that the latter variance can be estimated as,

$$\hat{V}ar_p(u_i | \mathbf{X}, I_i = 0) = \sum_{i \in s} \frac{w_i - 1}{\sum_{i \in s} (w_i - 1)} [\hat{r}_i - \sum_{i \in s} \frac{w_i - 1}{\sum_{i \in s} (w_i - 1)} \hat{r}_i]^2 - \frac{1}{m} \sum_{i \in s} \frac{\hat{\sigma}_e^2}{n_i}, \quad (6.9)$$

where  $\hat{\boldsymbol{\beta}}, \hat{\sigma}_e^2$  are sample estimates of  $\boldsymbol{\beta}, \sigma_e^2$  and  $\hat{r}_i = \frac{1}{n_i} \sum_{j \in s_i} (y_{ij} - \mathbf{x}'_{ij} \hat{\boldsymbol{\beta}})$ . Adding the estimator  $\hat{\sigma}_e^2 / N_i$  for (6.7) to (6.9) yields the estimator of the second variance in (6.4).

## 7. PREDICTION OF SMALL AREA MEANS UNDER A GENERAL SAMPLE MODEL AND SAMPLING SCHEME

In Section 5 we consider a particular two-level sample model (Eq. 5.1) and a particular relationship between the unit sampling weights and the outcome values within the selected areas (Eq. 5.2). Below we outline the basic steps in computing the predictors under a general two-level sample model fitted to the sample outcomes (Eqs. 2.1 and 2.3), with continuous or discrete outcomes and fixed and random effects, and a general relationship between the unit sampling weights and the outcomes. As in Section 5, we assume that the sample model and the conditional expectation of the unit

sampling weights have been identified and estimated based on the sample data. See, Rao (2003) for identification, estimation and diagnostic procedures in common use. The computations in Section 5 follow the same steps but take advantage of the normality assumptions in some of the derivations. We assume informative sampling of areas and within the areas and maintain *Ass.1* of Section 3.

First consider a sampled area. By (3.2), computation of the predictors of the area means requires estimating  $E_s[E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s]$ . By (2.7),

$$E_s[E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) | D_s] = E_s[H_{\mathbf{x}_{il}}(u_i) | D_s], \quad (7.1)$$

$$\text{where } H_{\mathbf{x}_{il}}(u_i) = \int y_{il} \frac{E_{si}(w_{li} | y_{il}, \mathbf{x}_{il}, u_i, I_i = 1) - 1}{E_{si}(w_{li} | \mathbf{x}_{il}, u_i, I_i = 1) - 1} f_{si}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1) dy_{il}.$$

The integral  $H_{\mathbf{x}_{il}}(u_i)$  depends on the sample model  $f_{si}(y_{il} | \mathbf{x}_{il}, u_i, I_i = 1)$  and the expectation  $E_{si}(w_{li} | y_{il}, \mathbf{x}_{il}, u_i, I_i = 1)$ , which as stated above are identified and estimated from the sample data. The integral can be computed either analytically or, if necessary, using numerical approximations. The expectation  $E_s[H_{\mathbf{x}_{il}}(u_i) | D_s]$  in (7.1) is with respect to the sample distribution  $f_s(u_i | D_s, I_i = 1)$ , which is obtained from the sample model defined by the Eqs. (2.1) and (2.3). This allows in principle to compute  $E_s[H_{\mathbf{x}_{il}}(u_i) | D_s]$  (with unknown parameters replaced by sample estimates). In practice it would often be sensible to assume  $E_{si}(w_{li} | y_{il}, \mathbf{x}_{il}, u_i, I_i = 1) = E_{si}(w_{li} | y_{il}, \mathbf{x}_{il}, I_i = 1)$ , as in (5.2).

Alternatively,  $E_s[H_{\mathbf{x}_{il}}(u_i) | D_s]$  can be estimated as,

$$\hat{E}_s[H_{\mathbf{x}_{il}}(u_i) | D_s] = \sum_{i \in s} H_{\mathbf{x}_{il}}(\hat{u}_i) / n, \quad (7.2)$$

where  $\hat{u}_i$  is the sample estimator of  $u_i$  under the sample model. (See below (5.5) for the estimator  $\hat{u}_i$  under the sample model (5.1)). Having estimated  $E_s[H_{\mathbf{x}_{il}}(u_i) | D_s]$ , the prediction of the area means follows the same steps as in Section 5.

Next consider nonsampled areas. By (3.3), predicting the means in such areas requires estimating,  $E_c[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 1) | D_s]$ . By (5.9),

$$\begin{aligned} E_c[E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 1) | D_s] &= \\ &= E_s[(w_i - 1) \frac{\int E_{si}\{w_{kij} | y_{ik}, \mathbf{x}_{ik}, u_i, I_i = 1\} y_{ik} f_{si}(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 1) dy_{ik}}{\int E_{si}\{w_{kij} | y_{ik}, \mathbf{x}_{ik}, u_i, I_i = 1\} f_{si}(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 1) dy_{ik}} | D_s] / [E_s(w_i | D_s) - 1], \quad (7.3) \\ &= E_s[(w_i - 1) K_{\mathbf{x}_{ik}}(u_i) | D_s] / [E_s(w_i | D_s) - 1] = K_{\mathbf{x}_{ik}}(D_s) \end{aligned}$$

where  $K_{\mathbf{x}_{ik}}(u_i)$  is the ratio of the two integrals. This ratio again only depends on the sample models  $f_{si}(y_{ik} | \mathbf{x}_{ik}, u_i, I_i = 1)$  and  $E_{si}(w_{kij} | y_{ik}, \mathbf{x}_{ik}, u_i, I_i = 1)$ , and it can be computed numerically if necessary. The last expectation in (7.3) can be estimated as,

$$\hat{K}_{\mathbf{x}_{ik}}(D_s) = \sum_{r \in S} (w_r - 1) K_{\mathbf{x}_{ik}}(\hat{u}_r) / \sum_{r \in S} (w_r - 1). \quad (7.4)$$

Note that no model is assumed for the area sampling weights so that  $K_{\mathbf{x}_{ik}}(D_s)$  can not be computed analytically under the joint distribution of  $(w_i, u_i) | D_s$  (see Eq. 5.11).

## 8. TESTING FOR PREDICTION BIAS

Evidently, predicting the small area means under informative sampling is more complicated and possibly less stable than under noninformative sampling. Thus, it is important to test the informativeness of the sample selection and if found noninformative, use standard optimal procedures. In what follows we propose simple test statistics for testing whether ignoring the sample selection biases the predictors.

### 8.1 Testing whether ignoring the selection of areas biases the predictors.

By (4.3), the selection of areas does not bias the predictors used under noninformative selection if  $Cov_s[\sum_{k=1}^{N_i} E_p(y_{ik} | \mathbf{x}_{ik}, u_i, I_{ik} = 1), w_i] | D_s = 0$ . However, one only needs to test  $Corr_s(u_i, w_i | D_s) = 0$  because if the true area mean is a function also of area covariates  $\mathbf{X}_i$ , say, the mean  $\bar{\mathbf{X}}_i = \sum_{j=1}^{N_i} \mathbf{x}_{ij} / N_i$  as under the model (5.1), dependence of  $w_i$  on  $\mathbf{X}_i$  alone does not bias the predictions. To see this, note that the

sample *pdf* of the area mean,  $\theta_i$  is by (2.5) and (2.6),  
 $f_s(\theta_i | \mathbf{X}_i) = E_s(w_i | \mathbf{X}_i) f_p(\theta_i | \mathbf{X}_i) / E_s(w_i | \theta_i, \mathbf{X}_i)$ , and if  $E_s(w_i | \theta_i, \mathbf{X}_i) = E_s(w_i | \mathbf{X}_i)$ ,  
 $f_s(\theta_i | \mathbf{X}_i) = f_p(\theta_i | \mathbf{X}_i)$ . This is true for general population models.

For testing  $H_0 : Corr_s(w_i, u_i | D_s) = 0$  we would ideally regress  $w_i$  against  $u_i$  but the random effects are unobservable. Thus, we regress instead  $w_i$  against the estimates  $\hat{u}_i$  as computed under the sample model. For the model (5.1), the estimates are defined in Section 5 as,  $\hat{u}_i = \hat{\gamma}_i [\bar{y}_i - \bar{\mathbf{x}}_i' \boldsymbol{\beta}]$ . Writing  $\hat{u}_i = u_i + \eta_i$ , it can be safely assumed that  $Cov(w_i, \eta_i) = 0$ , such that testing  $H_0$  can be implemented by regressing  $w_i = \delta_0 + \delta \hat{u}_i + \zeta_i$  and testing  $H_0 : \delta = 0$ .

Assuming that the  $\zeta_i$  are *iid* normal deviates, one can test  $H_0$  using the *t*-statistic,

$$t^A = \hat{\delta}_{OLS} / \sqrt{\hat{Var}(\hat{\delta}_{OLS})}, \quad (8.1)$$

which has then a *t*-distribution with  $(m-2)$  degrees of freedom under  $H_0$  (holding the estimates  $\hat{u}_i$  fixed). The null hypothesis refers to the sample distribution, thus justifying estimating  $\delta$  by OLS. For large number of sampled areas (large  $m$ ), the statistic  $t^A$  retains approximately the *t* distribution (effectively the standard normal distribution) even without the normality assumption, since the estimator  $\hat{\delta}_{OLS}$  has asymptotically a normal distribution based on the central limit theorem. Indeed, the empirical distribution of the statistic  $t^A$  in the simulation study described in Section 9 is extremely close to the nominal *t*-distribution, even though the distribution of the disturbances  $\zeta_i$  is in this case nonnormal (see table 3). On the other hand, one is obviously not restricted to using a *t*-test and other, more robust test procedures can be used instead, see, for example, Salibián-Barrera (2005). We mention also that regressing  $w_i$  against  $\hat{u}_i$  and testing  $H_0 : \delta = 0$  may not be very powerful if  $Var(\eta_i)$  is

large. An alternative test can possibly be constructed by noting that  $\hat{u}_i = u_i + \eta_i$  and using ‘errors in variables techniques’.

### 8.2 Testing whether ignoring the sampling within the areas biases the predictors.

By (4.1), sampling within the areas does not bias the predictors used under noninformative sampling if  $Cov_{si}(y_{il}, w_{li} | \mathbf{x}_{il}, u_i, I_i = 1) = 0$ . Hence, the ignorability of the sample selection within the selected areas can be tested by regressing  $w_{li} = \gamma_{0i} + \mathbf{x}'_{il}\boldsymbol{\gamma}_{1i} + \gamma_{2i}y_{il} + \eta_{il}$  and testing  $H_0 : \gamma_{2i} = 0$ , separately for every  $i \in s$ . However, with a large number of sampled areas, testing  $H_0$  for every area is not practical, and with small sample sizes within the selected areas, the tests have low power. Assuming the same sampling design within the areas, a more practical and powerful test is therefore,

$$F_{max}^w = \max[F_i, i = 1 \dots m], \quad (8.2)$$

where  $F_i$  defines the test statistic in area  $i$ . For a given distribution of  $F_i$ , computation of the percentiles of  $F_{max}^w$  is straightforward. Here again, if the disturbances  $\eta_{ij}$  can be assumed to be *iid* normal deviates, one can use the test statistics  $F_i = [\hat{\gamma}_{2i} / \hat{SD}(\hat{\gamma}_{2i})]^2$ , where  $\hat{\gamma}_{2i}, \hat{SD}(\hat{\gamma}_{2i})$  are respectively the OLS estimator and its estimated standard deviation. Under the null hypothesis  $H_0 : \gamma_{2i} = 0$ ,  $F_i \sim F(1, n_i - 3)$ . On the other hand, if the *iid* normality assumption is not warranted, the  $F$  distribution cannot be justified by asymptotic arguments as in the case of the statistic  $t^A$  in Section 8.1 since the sample sizes within the areas are typically small, and one has to use in this case a more robust test procedure. As with the test statistic  $t^A$ , the use of the test statistic  $F_{max}^w$  in the simulation study of Section 9 with  $F_i$  computed as above matches very closely the

corresponding nominal distribution, despite the fact that the distribution of the disturbances  $\eta_{ji}$  in this study is far from normal. See table 3.

Remark 5: Instead of testing  $Corr_{.si}(w_{jli}, y_{ij} | \mathbf{x}_{ij}, u_i, I_i=1) = 0$  by fitting a linear model, one can test  $H_0 : E_{.si}(w_{jli} | y_{ij}, \mathbf{x}_{ij}, u_i, I_i=1) = E_{.si}(w_{jli} | \mathbf{x}_{ij}, u_i, I_i=1)$ , allowing for other more plausible relationships between the weights  $w_{jli}$  and  $(\mathbf{x}_{ij}, y_{ij})$ , such as in (5.2). Note from (2.5) and (2.6) that  $E_{.si}(w_{jli} | y_{ij}, \mathbf{x}_{ij}, u_i, I_i=1) = E_{.si}(w_{jli} | \mathbf{x}_{ij}, u_i, I_i=1)$  implies that the population and sample distributions within the selected areas are the same.

## 9. MONTE-CARLO SIMULATION STUDY

In order to illustrate the bias that can occur when ignoring an informative sampling scheme and to assess the performance of the procedures developed in this article, we designed a small simulation study. The study consists of the following steps:

1- Generate area sizes,  $N_i = Int\{1000 \times (0.5 + \xi_i)\}$ ;  $\xi_i \sim U[0,1]$  and covariates

$$\mathbf{x}_{ij} = (50, t_i + \zeta_{ij})', \quad t_i = 1 + 3 \times Int\left[i - \frac{50}{3} \times Int\left(\frac{3}{50} \times i\right)\right] / 10, \quad \zeta_{ij} \sim U[0,5] \quad i = 1, \dots, 50,$$

$j = 1, \dots, N_i$ . Stratify the areas into 3 strata; stratum  $U_1$  consists of areas  $1 \leq i \leq 16$ , stratum  $U_2$  of areas  $16 < i \leq 33$  and stratum  $U_3$  of areas  $33 < i \leq 50$ . The complicated formula for generating the area values  $t_i$  satisfies that they are the same for each of the strata, except that Stratum 1 has only 16 areas.

2- Generate population random area effects,  $u_i \sim N(0, \sigma_u^2)$ ,  $i = 1, \dots, 50$ ,  $\sigma_u^2 = 100$ .

3- Generate  $y$ -values using the model (5.1) with  $\boldsymbol{\beta} = (1,1)'$ ,  $\sigma_e^2 = 100$ .

In order to avoid extreme selection probabilities, the random effects were truncated at  $\pm 2.5\sigma_u$ , and similarly for the residuals  $e_{ij}$  in (5.1).

4- Select 10 areas from each stratum with probabilities  $\pi_i = 10z_i / \sum_{j \in U_h} z_j$  by systematic PPS sampling, where  $z_i = \text{int}[1000 \times \exp(-u_i / 8\sigma_u)]$ , thus making the area selection informative.

5- Sample  $n_i$  units from selected area  $i$  by systematic PPS sampling with probabilities

$$\pi_{j|i} = n_i z_{ij} / \sum_{k=1}^{N_i} z_{ik}, \quad \text{where } z_{ij} = \exp\{[-(y_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta}) / \sigma_\varepsilon + \delta_{ij} / 5] / 3\}, \quad \delta_{ij} \sim N(0,1).$$

Note that the sampling of units is informative and that the sampling probabilities satisfy the relationship (5.2). The area sample sizes are fixed in a given stratum;  $n_i = 5$  if  $i \in U_1$ ,  $n_i = 25$  if  $i \in U_2$  and  $n_i = 50$  if  $i \in U_3$ .

6- Repeat Steps 2-5 10,000 times.

For each sample we computed 3 predictors of the area means:

A - 'Ordinary' small area predictors,

$$\hat{Y}_i^O = N_i^{-1} [n_i \bar{y}_i + (N_i - n_i) \hat{u}_i + (N_i \bar{\mathbf{X}}_i - n_i \bar{\mathbf{x}}_i)' \hat{\boldsymbol{\beta}}_{GLS}], \quad i \in s; \quad \hat{Y}_i^O = \bar{\mathbf{X}}_i' \hat{\boldsymbol{\beta}}_{GLS}, \quad i \notin s, \quad (9.1)$$

where  $\bar{\mathbf{X}}_i = \sum_{j=1}^{N_i} \mathbf{x}_{ij} / N_i$ ,  $\hat{u}_i = \hat{\gamma}_i (\bar{y}_i - \bar{\mathbf{x}}_i' \hat{\boldsymbol{\beta}}_{GLS})$ ,  $\hat{\gamma}_i = \hat{\sigma}_u^2 / [\hat{\sigma}_u^2 + \hat{\sigma}_\varepsilon^2 / n_i]$ ;  $(\hat{\sigma}_u^2, \hat{\sigma}_\varepsilon^2)$  were computed by the method of moments (fitting of constants) and  $\boldsymbol{\beta}$  by Generalized Least Squares with the unknown variances replaced by their estimators; see Rao (2003) for details. The predictors  $\{\hat{Y}_i^O\}$  are the EBLUP predictors of  $\bar{Y}_i$  for this model under noninformative sampling.

B- 'Design-based' estimators,

$$\hat{Y}_i^D = \bar{y}_{i,w} + (\bar{\mathbf{X}}_i - \bar{\mathbf{x}}_{i,w})' \hat{\boldsymbol{\beta}}_{PW} \quad \text{if } i \in s, \quad \hat{Y}_i^D = \bar{\mathbf{X}}_i' \hat{\boldsymbol{\beta}}_{PW} \quad \text{for } i \notin s, \quad (9.2)$$

$$(\bar{y}_{i,w}, \bar{\mathbf{x}}_{i,w}) = \sum_{j \in s_i} w_{ij} (y_{ij}, \mathbf{x}_{ij}) / \sum_{j \in s_i} w_{ij},$$

$$\hat{\boldsymbol{\beta}}_{PW} = \left[ \sum_{i \in s, j \in s_i} w_i w_{j|i} \mathbf{x}_{ij} \mathbf{x}'_{ij} \right]^{-1} \sum_{i \in s, j \in s_i} w_i w_{j|i} \mathbf{x}_{ij} y_{ij}.$$

The predictor  $\hat{Y}_i^D$  for  $i \notin S$  is not really a ‘design based’ estimator and is similar to the estimator in (9.1), except that  $\hat{\boldsymbol{\beta}}_{GLS}$  is replaced by the probability weighted estimator  $\hat{\boldsymbol{\beta}}_{PW}$ . As discussed in the introduction, design based theory is not suited for the prediction of means in nonsampled areas.

C- The new predictors  $\hat{Y}_i^N$ . The predictors are defined by (5.8) for sampled areas and by (5.12) for nonsampled areas. Note that since the population random effects are normal and because of the sampling scheme used to select the areas, the sample random effects also have a normal distribution but with different expectation, thus justifying the use of these predictors. The model parameters  $\sigma_u^2, \sigma_e^2, \boldsymbol{\beta}$  have been estimated in the same way as for the estimators in A. The coefficients  $\mathbf{a}, b, k_i$  indexing the relationship between the weights  $w_{ji}$  and the outcome and auxiliary variables were estimated by fitting the model (5.2), using the procedures REG and NLIN in SAS.

In addition to the three sets of predictors we computed also the test statistics developed in Section 8 and the variance estimators of the predictors  $\hat{Y}_i^N$  developed in Section 6, distinguishing between sampled and nonsampled areas. Since the computation of the variances requires generating bootstrap samples, we restricted this part of the simulation study to 300 samples and 300 bootstrap samples for each sample.

Table 1 shows the empirical prediction bias and root mean square error (RMSE) of the three predictors over the 10,000 simulations, separately for sampled and nonsampled areas. Denote by  $\bar{Y}_r$  the true mean of area  $t$  in simulation  $r$ ,  $r=1 \dots 10,000$ , and let  $\hat{Y}_r$  represent any of the predictors. Define  $D_r = 1$  if area  $t$  is sampled in simulation  $r$  and  $D_r = 0$  otherwise. For given area  $t$ , the prediction bias and RMSE when this area is sampled are computed as,



$$Bias_t = \sum_{r=1}^{10,000} D_{tr} (\hat{Y}_{tr} - \bar{Y}_{tr}) / \sum_{r=1}^{10,000} D_{tr} ; RMSE_t = \sqrt{\sum_{r=1}^{10,000} D_{tr} (\hat{Y}_{tr} - \bar{Y}_{tr})^2 / \sum_{r=1}^{10,000} D_{tr}} \quad (9.3)$$

The prediction bias and RMSE when area  $t$  is not sampled are obtained by replacing  $D_{tr}$  by  $(1 - D_{tr})$  in (9.3). The results in Table 1 are averages over the areas contained in the same stratum (having the same sample size). Table 1 shows also the means of the variance estimators developed in Section 6.

Table 1 about here

The conclusions from Table 1 are clear-cut:

- 1- Ignoring the informative sampling scheme induces large prediction bias for both sampled and nonsampled areas. The large biases induce large RMSEs.
- 2- The design based estimators are approximately unbiased in sampled areas when the sample sizes within the areas are sufficiently large ( $n_i = 25$  in our study), but are biased when estimating the means of nonsampled areas. Recall that no design unbiased predictor for a given nonsampled area exists in general.
- 3- The predictors  $\hat{Y}_i^N$  are literally unbiased for both sampled and nonsampled areas.
- 4- The RMSEs of all the predictors for sampled areas decrease as the sample sizes within the areas increase.
- 5- The RMSEs of the predictor  $\hat{Y}_i^N$  in nonsampled areas are lower than the RMSEs of the other predictors but they seem high, particularly when compared to the RMSEs obtained for the sampled areas. Note, however, that for nonsampled areas the standard deviation of the random effect is  $Std_c(u_i) \cong 9.75$ , which is only slightly smaller than the RMSEs of  $\hat{Y}_i^N$ .
- 6- The RMSE estimates are basically unbiased for both sampled and nonsampled areas. The magnitude of the bias and the precision of the RMSE estimators can be further assessed by the performance of confidence intervals for the area means that use them.

Table 2 shows the coverage rates of the conventional confidence intervals

$\hat{Y}_i^N \pm Z_{1-\frac{\alpha}{2}} \sqrt{\hat{Var}(\hat{Y}_i^N)}$  for  $(1-\alpha) = 0.90, 0.95, 0.99$ . The match with the nominal levels

is almost perfect.

Table 2 about here

Finally, when applying the tests  $t^A$  (Eq. 8.1) and  $F_{\max}^w$  (Eq. 8.2), the null hypothesis of noninformative selection of areas was always rejected with p-values smaller than 0.01, and the null hypothesis of noninformative sampling within the selected areas was always rejected with p-values smaller than 0.025.

Notwithstanding, statistical tests have their limitations, and it is important to assess the performance of the new predictors when the sample selection is in fact noninformative. To this end, we sampled the areas with probabilities proportional to a size variable  $z_i \sim U[1,2]$  that is independent of the random area effects, and sampled the units within the selected areas with probabilities proportional to  $z_{ij} = \exp(\varphi_{ij}/15)$ , where  $\varphi_{ij} \sim N(0,1)$ , independently of  $y_{ij}$ . The results of this exercise can be summarized as follows:

1. For noninformative sampling of areas and within the areas, the ordinary small area predictor  $\hat{Y}_i^O$  is in common use, but the new predictor, although much more entangled, performs equally well both in terms of bias and RMSE.
2. The RMSE estimators perform well, with a positive bias of up to 4% in sampled areas and a negative bias of up to 3% in nonsampled areas. Also, the match between the empirical coverage rates and the nominal levels is again almost perfect.
3. Table 3 shows the distributions of the test statistics  $t^A$  (Eq. 8.1) and  $F_{\max}^w$  (Eq. 8.2 with) as obtained for this case. The empirical percentiles for both tests are almost

identical to the nominal percentiles, despite the fact that the disturbances in the models for  $w_i$  and  $w_{ji}$  are not normal.

## 10. ESTIMATION OF MEAN BODY MASS INDEX IN USA COUNTIES

### 10.1 *The sample data*

In this section we apply the methodology developed in the previous sections for estimating the mean body mass index (BMI) for counties in the U.S. The BMI is defined as the ratio between the weight, measured in kilograms, and the square of the height, measured in meters. An index higher than 27.8 for men and higher than 27.3 for women defines overweight, which is known to be a major health risk factor. Estimating the mean BMI at the national and sub-national level is therefore of prime importance for health authorities dealing with this problem. The data used for this study were collected as part of the third national health and nutrition examination survey (NHANES III). The survey was conducted in two phases during the years 1988-1991 and 1991-1994, and it represents the U.S. total civilian non-institutional population.

NHANES III is a stratified four-stage clustered survey that collects health, dietary and background information through questionnaires and physical examinations. The primary sampling units (PSU) are in most cases individual counties. There are 81 PSUs in the sample, selected with probability proportional to a measure of size without replacement. The size measure was constructed in such a way that the survey oversampled PSUs with large populations of Mexican-Americans and Blacks. The second stage of the sample selection consisted of sampling of area segments, which were then stratified based on the percent of Mexican-Americans. Next, households were sampled within the strata, with higher rates for strata with high minority concentrations. In the last stage a sample of persons was sampled from classes of households defined by age, sex and race that were sampled at different rates. For more

details and the computation of the sampling weights, see <http://www.cdc.gov/nchs/about/major/nhanes/nh3data.htm>. The data set used for this study refers to the 81 sampled counties. There are 3138 counties in total in the U.S. The sample sizes within the sampled counties exceed 80 in almost all the counties, with a total sample size of 16,521, divided into 8767 women and 7754 men. Thus, the major small area estimation problem with this survey is that only a small fraction of the counties that define the areas is represented in the sample.

## 10.2 Analysis

In a previous article, Malec *et al.* (1999) used NHANES III data for estimating overweight prevalence for states in the U.S. by fitting logistic models with fixed age/race/gender effects and random race/gender effects. In order to account for sampling effects within the selected counties, the authors estimated the sampling probabilities utilizing the sampling weights, and then substituted the estimates in the likelihood. The state prevalence estimates were obtained by applying the Bayesian approach with the resulting empirical likelihood, using MCMC simulations.

In our application we fit the model (5.1), separately for men and women, with county random effects and seven covariates: A constant, 3 dummy race variables and 3 age variables. The race variables are:  $x_1 = 1$  if non Hispanic white,  $x_2 = 1$  if non Hispanic black and  $x_3 = 1$  if Hispanic. The age variables are:  $x_4 = age \times I_{20 \leq age < 50}$ ,  $x_5 = age \times I_{50 \leq age < 75}$ ,  $x_6 = age \times I_{75 \leq age}$ . The age variables are used as proxy for a quadratic relationship between the BMI and age. We could not include  $age^2$  in the model because the county means of this variable are unknown. There are a few other covariates with sample measurements that affect the BMI but could not be used for the same reason. One of these variables is education, measured by the number of years at school, which was found to have a negative effect on the BMI of women. The data

files that we could use only contain information on the county numbers of adults with college and higher education, but this information is unknown at the individual level.

Table 4 shows the estimated regression coefficients, their standard errors (S.E.) and the estimates of the variance of the random effects and the residual variance. All the coefficients except in the case of ‘White non Hispanic’ in the women’s model are significant at the 5% level based on the conventional  $t$ -statistic. We tested the assumption that the residual variance is constant across the counties by first fitting the model for each of the sampled counties separately, assuming fixed county effects and then testing the homogeneity of the estimated residuals. After dropping 7 outlying counties, the hypothesis of homogeneity is accepted using Bartlett’s test with p-values of 0.99 for women and 0.13 for men.

Table 4 about here

Next we applied the tests of sample ignorability considered in Section 8. We found that for both men and women the sampling within the counties does not introduce prediction bias (given the covariates included in the model), and that the sampling of counties is informative for women, but not for men. The p-values when testing the sample ignorability within the counties are 0.56 for women and 0.41 for men. The sample ignorability within the counties has been tested also by regressing  $\log(w_{ji})$  against  $(y_{ij}, \mathbf{x}_{ij})$  instead of  $w_{ji}$  (see Remark 5 in Section 8.2), and by fitting the two regression models in each of the sampled counties separately, confirming in all the cases that for the present model the sample selection within the counties can be ignored. On the other hand, when testing the ignorability of the county selection using (8.1), the p-values are 0.0164 for women and 0.31 for men, suggesting an informative sampling of counties for the women’s model but not for the men’s model.

As explained in Section 10.1, the sampling probabilities within the counties were determined by the race and age characteristics, and hence it is not surprising that the

sampling within the counties was found to be ignorable for the present model that includes race and age as explanatory variables. It is interesting to mention in this regard that Malec *et al.* (1999) found that the sampling within the counties is informative, despite the fact that their model likewise accounts for age and race/gender categories. The authors do not elaborate on the reasons for this finding but they show results illustrating different national and state estimates, depending on whether the sampling process is accounted for or not.

The result that the sampling of counties is informative for the women's model is likewise not surprising because the county selection probabilities were determined by the true county race totals, and these totals are not included in the model (see below). The model of Malec *et al.* (1999) contains fixed and random race parameters, which is probably why the authors concluded that the selection of counties is not informative for their model. The fact that the selection of counties can be ignored for the men's model in our application is probably related to the fact that the variance of the county random effects is small,  $\hat{\sigma}_u^2 = 0.76$ , which makes it harder to detect selection effects.

As mentioned in the introduction, a possible way of controlling sampling effects is by including in the model all the design variables used for the sample selection. In the present application we are in a fortunate (but uncommon) situation where the county design variables;  $x_{8i}$  = county total of non Hispanic White,  $x_{9i}$  = county total of non Hispanic Black and  $x_{10i}$  = county total of Hispanic, are known. Adding these variables (divided by  $10^5$ ) to the model yields the following coefficients and standard errors. Women:  $\beta_8 = -0.112(0.076)$ ,  $\beta_9 = 0.089(0.200)$ ,  $\beta_{10} = 0.141(0.141)$ . Men:  $\beta_8 = -0.017(0.043)$ ,  $\beta_9 = -0.064(0.115)$ ,  $\beta_{10} = 0.037(0.079)$ . The coefficients and standard errors of the other covariates change only slightly from their values in Table 4 when fitting the model with only the six covariates. Thus, all three design variables are

highly insignificant individually, and they are also jointly insignificant with p-values of 0.42 for women and 0.69 for men. With such high p-values, many analysts would tend to drop the design variables from the model and conclude that the sampling of counties is noninformative for the six covariates model, which in view of the previous analysis is not true for the women's model. Furthermore, when re-estimating the random effects using the extended model that includes the three design variables, and applying the informativeness test in (8.1), we find that the sampling of counties is not informative for this model, with p-values 0.17 for women and 0.63 for men. Thus, the selection of counties can only be ignored when including the design variables in the model.

What are the implications of the use of the model with six covariates or the model with 9 covariates (including the 3 design variables) on the prediction of the county means? In what follows we restrict to the models for women because the selection of counties was found earlier to be ignorable for the men's model. Starting with the sampled areas, both models yield very similar predictors when using the predictors defined by (9.1), which are the empirical best linear predictors (EBLUP) under noninformative sampling within the areas ( $b = 0$  in (5.8)). For the nonsampled areas, however, they yield somewhat different predictors. Figure 1 shows four different predictors of the means in nonsampled areas. The predictor  $\bar{\mathbf{X}}_i' \hat{\boldsymbol{\beta}}_{GLS}$  under the reduced model (6 covariates) as obtained when ignoring the county selection (Eq. (9.1)), the predictor  $\bar{\mathbf{X}}_i' \hat{\boldsymbol{\beta}}_{GLS}$  under the extended model with 9 regressors, (the vector  $\bar{\mathbf{X}}_i$  contains in this case both the proportions and the totals of the three races), the empirical predictor (5.12) under the reduced model ( $b \neq 0$ ), and the predictor (5.12) under the extended model. The horizontal line at 27.3 marks the threshold defining overweight. For the predictor (5.12) under the reduced model the area selection bias correction,

$$\sum_{i \in S} (w_i - 1) \hat{u}_i / \sum_{i \in S} (w_i - 1) \text{ is } 0.47 \text{ with Jackknife estimated standard deviation of } 0.16.$$

For the predictor (5.12) under the extended model the bias correction is 0.25 with similar estimated standard deviation. The use of the bias correction for the extended model is therefore questionable, which is consistent with the test result that the selection of counties is ignorable for this model. The use of Jackknife for variance estimation assumes that the random effects  $\hat{u}_i$  are approximately independent. It is used here only as a rough measure for assessing the stability of the bias correction.

The 4 plots in Figure 1 suggest that ignoring the county selection method and using the synthetic predictor based on the 6 regressors model under-predicts the true county means. This becomes evident by comparing the synthetic predictors under this model with the synthetic predictors under the extended model. The latter predictors are lower than the predictors obtained under the 6 covariates model with the bias correction, but interesting enough, once the bias correction is added also to the predictors under the extended model, both sets of predictors behave very similarly. However, as discussed above, the use of a bias correction for the extended model is questionable.

The magnitudes of the bias corrections seem very small, but they are not negligible. To see this, we computed the percentage of nonsampled areas for which the predicted means are higher than the threshold of 27.3, as obtained by use of the four predictors. The use of the two synthetic predictors yields a percentage of 2.84% for the six covariates model and 5.56% for the extended model. Adding the bias correction of 0.47 to the first synthetic predictor increases the percentage to 9.2%, whereas adding the bias correction of 0.25 to the second synthetic predictor further increases the percentage to 10.3%. Thus, if areas with means that exceed the threshold are to be given extra attention, the use of the bias correction can be very important.

## 11. CONCLUDING REMARKS

This article presents a first attempt of predicting small area means under informative sampling of areas and within the areas. The proposed procedure assumes



knowledge of the models holding for the sample data and for the sampling weights within the selected areas, but otherwise is ‘model free’. Both models can be identified estimated and tested from the sample data. In the present application we consider the familiar nested error regression model but as outlined in Section 7, the procedure can be applied to other models with continuous or discrete outcomes using similar steps.

Much of the research in small area estimation concerns the use of Bayesian methods that allow considering heavily structured models and accounting for all sources of variation when assessing the prediction errors. In this article we restrict to the frequentist approach but it would seem that the proposed procedure can be applied in a Bayesian set up, except that it will require modelling the relationship between the area selection probabilities and the true area means, which as discussed in Section 5 is more complicated but not necessary under the present procedure. Developing a Bayesian solution that does not require this extra step is an intriguing problem.

#### APPENDIX A: DERIVATION OF EQUATIONS 6.1, 6.2

We note first that  $\mathbf{X} \subset D_s$  and that  $D_s$  defines also  $\mathbf{I}$  such that conditioning on  $D_s$  implies conditioning on  $\mathbf{X}$  and  $\mathbf{I}$  as well. Since conditional on  $(D_s, u_i, I_i = 1)$ ,

$\hat{Y}_{Ri} - E_{ci}[Y_{Ri} | D_s, u_i, I_i = 1]$  is constant, it follows from *Ass.2* that,

$$\begin{aligned} E_p[(\hat{Y}_i - Y_i)^2 | \mathbf{X}, \mathbf{I}] &= E_p\{E_p[(\hat{Y}_i - Y_i)^2 | D_s, u_i, I_i = 1] | \mathbf{X}, \mathbf{I}\} \\ &= E_p\{[\hat{Y}_{Ri} - E_{ci}(Y_{Ri} | D_s, u_i, I_i = 1)]^2 | \mathbf{X}, \mathbf{I}\} \\ &+ E_p\{E_p[(Y_{Ri} - E_{ci}(Y_{Ri} | D_s, u_i, I_i = 1))^2 | D_s, u_i, I_i = 1] | \mathbf{X}, \mathbf{I}\} \\ &= E_p\{[\hat{Y}_{Ri} - E_{ci}(Y_{Ri} | D_s, u_i, I_i = 1)]^2 + E_{ci}[(Y_{Ri} - E_{ci}(Y_{Ri} | \mathbf{X}, u_i, I_i = 1))^2 | \mathbf{X}, u_i, I_i = 1]\} | \mathbf{X}, \mathbf{I}\} \\ &= E_p[G(u_i, D_s) | \mathbf{X}, \mathbf{I}], \text{ where} \end{aligned}$$

$$G(u_i, D_s) = [\hat{Y}_{Ri} - E_{ci}(Y_{Ri} | D_s, u_i, I_i = 1)]^2 + E_{ci}[(Y_{Ri} - E_{ci}(Y_{Ri} | \mathbf{X}, u_i, I_i = 1))^2 | \mathbf{X}, u_i, I_i = 1],$$

thus establishing Eq. (6.1). By (5.3) and (5.4),

$$E_{ci}(y_{il}^2 | \mathbf{x}_{il}, u_i, \mathbf{I}_i = 1) = \frac{\lambda_{il}}{\lambda_{il} - 1} [\sigma_e^2 + (u_{il} + b\sigma_e^2)^2] - \frac{1}{\lambda_{il} - 1} [\sigma_e^2 + u_{il}^2],$$

$$E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, \mathbf{I}_i = 1) = u_{il} + \frac{\lambda_{il}}{\lambda_{il} - 1} b\sigma_e^2; \quad u_{il} = \mathbf{x}_{il}' \boldsymbol{\beta} + u_i. \quad \text{Hence,}$$

$$E_{ci}[(y_{il} - E_{ci}(y_{il} | \mathbf{x}_{il}, u_i, \mathbf{I}_i = 1))^2 | \mathbf{x}_{il}, u_i, \mathbf{I}_i = 1] = \sigma_e^2 - \frac{\lambda_{il} b^2 \sigma_e^4}{(\lambda_{il} - 1)^2}. \quad \text{It follows that under}$$

*Ass.1* and *Ass.2*,  $G(u_i, D_s)$  in (6.1) can be written as in (6.2).

#### APPENDIX B: ESTIMATION OF $Var_p(u_i | \mathbf{X}, \mathbf{I}_i = 0)$

Let  $\tilde{\xi} = \{\tilde{u}_i, \tilde{e}_{ij}, \tilde{\mathbf{I}}_i, \tilde{w}_i, (i, j) \in U\}$  be a generic random vector distributed identically but independently of  $\xi = \{u_i, e_{ij}, \mathbf{I}_{ij}, \mathbf{I}_i, w_i, (i, j) \in U\}$  given  $\mathbf{X}$  under the population distribution, that is,  $P_p(\tilde{\xi} \in A | \mathbf{X}) = P_p(\xi \in A | \mathbf{X})$ , for every set  $A$  belonging to the  $\sigma$ -algebra generated by  $\xi$ . Define,  $\tilde{y}_{ij} = \mathbf{x}'_{ij} \boldsymbol{\beta} + \tilde{u}_i + \tilde{e}_{ij}$ . Then by *Ass.3*,

$$\begin{aligned} Var_p[\sum_{i \in S} \frac{1}{n_i} \sum_{j \in S_i} (\tilde{y}_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta}) | \mathbf{X}, \mathbf{I}, \tilde{\mathbf{I}}_i = 0] &= Var_p[(\sum_{i \in S} \tilde{u}_i + \sum_{i \in S} \frac{1}{n_i} \sum_{j \in S_i} \tilde{e}_{ij}) | \mathbf{X}, \mathbf{I}, \tilde{\mathbf{I}}_i = 0] \\ &= m Var_p(u_i | \mathbf{X}, \mathbf{I}_i = 0) + \sum_{i \in S} \frac{\sigma_e^2}{n_i}, \quad \text{such that,} \end{aligned}$$

$$Var_p(u_i | \mathbf{X}, \mathbf{I}_i = 0) = \frac{1}{m} Var_p[\sum_{i \in S} \frac{1}{n_i} \sum_{j \in S_i} (\tilde{y}_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta}) | \mathbf{X}, \mathbf{I}, \tilde{\mathbf{I}}_i = 0] - \frac{1}{m} \sum_{i \in S} \frac{\sigma_e^2}{n_i}. \quad (\text{B1})$$

Let  $\tilde{r}_i = \frac{1}{n_i} \sum_{j \in S_i} (\tilde{y}_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta})$ ;  $r_i = \frac{1}{n_i} \sum_{j \in S_i} (y_{ij} - \mathbf{x}'_{ij} \boldsymbol{\beta})$ . Then, by (B1) and (2.8),

$$Var_p(u_i | \mathbf{X}, \mathbf{I}_i = 0) = \frac{1}{m} \sum_{i \in S} E_s \left\{ \frac{\tilde{w}_i - 1}{E_s(\tilde{w}_i | \mathbf{X}) - 1} [\tilde{r}_i - E_s \frac{\tilde{w}_i - 1}{E_s(\tilde{w}_i | \mathbf{X}) - 1} \tilde{r}_i]^2 \right\} - \frac{1}{m} \sum_{i \in S} \frac{\sigma_e^2}{n_i}. \quad (\text{B2})$$

The estimator (6.9) follows from (B2) by substituting  $\hat{\boldsymbol{\beta}}, \hat{\sigma}_e^2$  for  $\boldsymbol{\beta}, \sigma_e^2$  and replacing the expectations under the sample distribution by the corresponding sample means.

#### REFERENCES

Arora, V. and Lahiri, P. (1997), "On the superiority of the Bayesian method over the BLUP in small area estimation problems," *Statistica Sinica* **7**, 1053-1063.

- Battese, G.E., Harter, R. M. and Fuller, W.A. (1988), "An error component model for prediction of county crop areas using survey and satellite data," *Journal of the American Statistical Association* 83, 28-36.
- Hall, P. and Maiti, T. (2006), "On parametric bootstrap methods for small area predictions," *Journal of the Royal Statistical Society* 68, Series B, 221-238.
- Kim, D. H. (2002), "Bayesian and empirical Bayesian analysis under informative sampling," *Sankhya B*, 64, 267-288.
- Kott, P.S. (1990), "Robust small domain estimation using random effects modeling," *Survey Methodology* 15, 3-12.
- Malec, D., Davis, W. W., and Cao, X. (1999). "Model-based small area estimates of overweight prevalence using sample selection adjustment," *Statistics in Medicine* 18, 3189-3200.
- Pfeffermann, D., Krieger, A. M. and Rinott, Y. (1998), "Parametric distributions of complex survey data under informative probability sampling," *Statistica Sinica* 8, 1087-1114.
- Pfeffermann, D., and Sverchkov, M. (1999), "Parametric and semi-parametric estimation of regression models fitted to survey data," *Sankhya* 61, 166-186.
- Pfeffermann, D., and Sverchkov, M. (2003), "Fitting generalized linear models under informative probability sampling," In: *Analysis of Survey Data*, eds. C. J. Skinner and R. L. Chambers, New York: Wiley, 175-195.
- Pfeffermann, D., Moura, F. A. S. and Nascimento-Silva, P. L. (2006), "Multilevel modeling under informative sampling," *Biometrika*, 93, 943-959.
- Prasad, N. G. N., and Rao, J. N. K. (1999), "On robust small area estimation using a simple random effects model," *Survey Methodology* 25, 67-72.
- Rao, J. N. K. (2003), *Small Area Estimation*. New York: Wiley.

Salibián-Barrera, M. (2005). “Estimating the p-values of robust tests for the linear model,” *Journal of Statistical Planning and Inference* 128, 241-257.

Skinner, C. J. (1994), “Sample models and weights,” *1994 Proceedings of the American Statistical Association*, Survey Research Methods Section, 133-142.

Sverchkov, M., and Pfeffermann, D. (2004), “Prediction of finite population totals based on the sample distribution,” *Survey Methodology*, 30, 79-92.

*Table 1. Bias, Root Mean Square Error (RMSE) and mean of RMSE estimators (RMSE-E). Informative sampling of areas and within areas. Simulation results.*

	Sample size	Sampled Areas			Nonsampled Areas		
		Ordinary $\hat{Y}_i^O$	Design $\hat{Y}_i^D$	New $\hat{Y}_i^N$	Ordinary $\hat{Y}_i^O$	Design $\hat{Y}_i^D$	New $\hat{Y}_i^N$
Bias	$n_i = 5$	-3.25	-0.71	-0.02	-6.36	-2.00	-0.32
	$n_i = 25$	-3.27	-0.14	-0.09	-6.10	-1.73	-0.06
	$n_i = 50$	-3.27	-0.07	-0.15	-6.10	-1.73	-0.06
RMSE	$n_i = 5$	5.26	4.88	4.14	11.77	10.04	9.85
	$n_i = 25$	3.80	2.19	1.95	11.70	10.08	9.93
	$n_i = 50$	3.54	1.54	1.39	11.71	10.11	9.96
RMSE-E	$n_i = 5$	---	---	4.28	---	---	9.90
	$n_i = 25$	---	---	2.02	---	---	9.91
	$n_i = 50$	---	---	1.46	---	---	9.91

*Table 2. Coverage rates of confidence intervals for true area means. Informative sampling of areas and within areas. Simulation results.*

	Nominal levels	Sampled Areas			Nonsampled Areas		
		0.90	0.95	0.99	0.90	0.95	0.99
Sample size	$n_i = 5$	0.90	0.94	0.98	0.91	0.95	0.99
	$n_i = 25$	0.89	0.94	0.98	0.91	0.95	0.99
	$n_i = 50$	0.89	0.94	0.98	0.92	0.96	0.99

Table 3. Distribution of test statistics, noninformative sampling. Simulation results.

Percentiles	0.01	0.025	0.05	0.10	0.90	0.95	0.975	0.99
Sampling of areas	0.013	0.029	0.053	0.107	0.896	0.952	0.975	0.988
Sampling within areas	0.009	0.025	0.049	0.093	0.903	0.948	0.976	0.988

Table 4. Regression coefficients, S.E. (in parentheses) and variances for BMI models

Coeff.	Intercept	White Non Hispanic	Black Non Hispanic	Hispanic	Age<50	50≤Age<75	Age≥75
Men	22.960 (0.414)	0.739 (0.314)	0.740 (0.316)	1.161 (0.322)	0.083 (0.008)	0.056 (0.005)	0.020 (0.004)
Women	21.852 (0.526)	-0.670 (0.374)	2.355 (0.375)	1.602 (0.394)	0.133 (0.010)	0.095 (0.006)	0.049 (0.005)

Men:  $\sigma_u^2 = 0.760$ ,  $\sigma_e^2 = 23.040$  ; Women:  $\sigma_u^2 = 2.830$ ,  $\sigma_e^2 = 39.560$

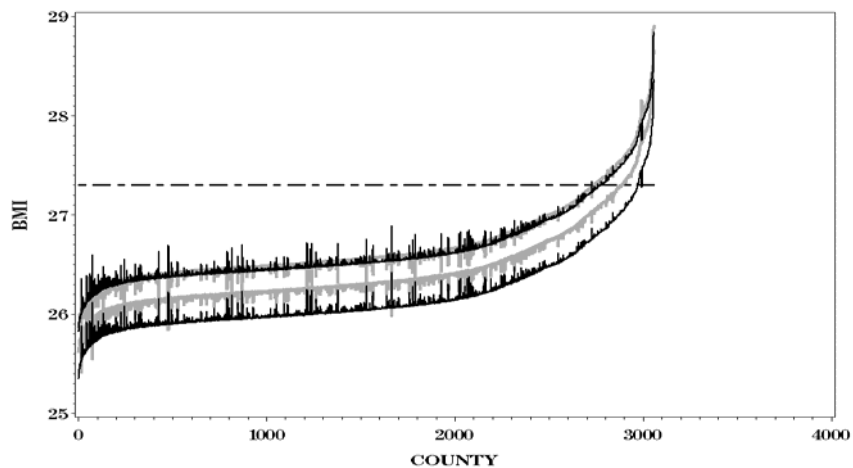


Figure 1. Prediction of mean body mass index of women in nonsampled counties of NHANES III. Values above the horizontal line at 27.3 define 'overweight'.

The lower dark and grey curves show the synthetic predictors under the six covariates model and the 9 covariates model respectively. The upper dark and grey curves show the corresponding predictors with bias corrections. The counties are ordered by the average values of the 4 predictors.