



Catalogue no. 11-522-XIE

**Statistics Canada International Symposium
Series - Proceedings**

**Symposium 2003: Challenges
in Survey Taking for the Next
Decade**

2003



Statistics
Canada

Statistique
Canada

Canada

Proceedings of Statistics Canada Symposium 2003
Challenges in Survey Taking for the Next Decade

PRACTICAL ISSUES IN MODEL-BASED SMALL AREA ESTIMATION

J.N.K. Rao¹

ABSTRACT

Traditional area-specific direct estimators may not provide acceptable precision for small areas because sample sizes in small areas are seldom large enough. This makes it necessary to “borrow strength” across related areas through indirect estimation based on linking models. Both area level and unit level models have been extensively studied in the literature to derive empirical best linear unbiased prediction (EBLUP), empirical Bayes (EB) and hierarchical Bayes (HB) estimators for small areas and associated measures of variability. This paper covers several practical issues related to model-based small area estimation, including model determination, choice of priors on model parameters in HB estimation, benchmarking to reliable direct large area estimators and use of survey weights in model-based estimation.

KEYWORDS: Area Level Model; Composite Estimate; Unit Level Model.

1. INTRODUCTION

Due to cost and other considerations, sample surveys are typically designed to provide area-specific (or direct) estimators with small sampling coefficient of variation (CV) for large areas (or domains). In fact, survey practitioners often stress that nonsampling errors, including measurement and coverage errors, contribute much more than sampling errors to total mean squared error (MSE) which is often used as a measure of quality of estimators. However, sampling errors play a dominant role in small area estimation because sample sizes in small areas are seldom large enough to provide direct estimators with acceptable quality in terms of sampling MSE (or CV). In fact, sample sizes can be zero in many small areas of interest. For example, data from the Current Population Survey (CPS) are used to estimate county (and school district) counts of poor school age children in the United States, but the CPS sample sizes are zero in many of the counties (National Research Council, 2000).

Due to difficulties with direct estimators, it is often necessary to employ indirect estimates that borrow information from related areas through explicit (or implicit) linking models, using census and administrative data associated with the small areas. Indirect estimators based on explicit linking models have received a lot of attention in recent years because of the following advantages over the traditional indirect estimators based on implicit models: (i) Explicit model-based methods make specific allowance for local variation through complex error structures in the model that link the small areas. (ii) Models can be validated from the sample data. (iii) Methods can handle complex cases such as cross-sectional and time series data, binary or count data, spatially-correlated data and multivariate data. (iv) Area-specific measures of variability associated with the estimates may be obtained, unlike overall measures commonly used with the traditional indirect estimators.

Basic area level and unit level models have been extensively studied in the literature to derive empirical best linear unbiased prediction (EBLUP), empirical Bayes (EB) and hierarchical Bayes (HB) small area estimators of totals (or means) and associated measures of variability. The EBLUP method is applicable for linear mixed models that cover the basic area level and unit level models. On the other hand, EB and HB methods are more generally applicable, covering generalized linear mixed models that are used to handle categorical (e.g., binary) and count data. MSE is used as a measure of variable under the EBLUP and EB approaches, while the HB approach uses the posterior variance as a measure of variability, assuming a prior distribution on the model parameters. We refer the reader to Rao (2003) for an extensive account of EBLUP, EB and HB methods for small area estimation. In this paper, I will cover several practical issues related to model-based small area estimation, including model determination, methods

¹ Carleton University, School of Mathematics and Statistics, Ottawa, Canada, e-mail: jrao@math.carleton.ca

for MSE estimation, choice of priors on model parameters in HB estimation, benchmarking to reliable direct large area estimators and use of survey weights in model-based estimation.

2. SMALL AREA MODELS

Two types of basic small area models have been studied in the literature. In the first type, called the basic area level model, only area-specific auxiliary data $z_i = (z_{i1}, \dots, z_{ip})^T$, related to some suitable function $\theta_i = g(Y_i)$ of the small area total Y_i ($i = 1, \dots, m$), are used to develop a linking model of the form $\theta_i = z_i^T \beta + v_i$ with $v_i \stackrel{iid}{\sim} N(0, \sigma_v^2)$, where σ_v^2 is the model variance. The linking model is combined with the matching sampling model $\hat{\theta}_i = \theta_i + e_i$, where $\hat{\theta}_i = g(\hat{Y}_i)$ is a direct estimator of θ_i and $e_i | \theta_i \stackrel{ind}{\sim} N(0, \psi_i)$ with known sampling variance ψ_i . The combined model, $\hat{\theta}_i = z_i^T \beta + v_i + e_i$, is a special case of the linear mixed model.

The basic area level model has at least two limitations. First, the assumption of known sampling variances, ψ_i , is restrictive, although methods based on generalized variance functions (GVF) have been proposed to produce smoothed estimates of the ψ_i 's. Secondly, the assumption $E(e_i | \theta_i) = 0$ may not be tenable if the small area sample size, n_i , is small and θ_i is a nonlinear function of the total Y_i , even if the director estimator \hat{Y}_i is design-unbiased for Y_i . It is more realistic to use the sampling model $\hat{Y}_i = Y_i + f_i$ with $E(f_i | Y_i) = 0$, which simply says that \hat{Y} is design-unbiased for Y_i . Further, we assume that $V(f_i | Y_i) = \sigma_i^2$, where the sampling variance may depend on Y_i ; for example, $\sigma_i^2 = Y_i^2 c_i^2$, where c_i is the known coefficient of variation of \hat{Y}_i ascertained from fitting GVF's. The sampling model is now unmatched with the linking model in the sense that they cannot be combined directly to produce a linear mixed model. Various extensions of the basic area level (also called Fay-Herriot model) have been proposed to handle correlated sampling errors, spatial dependence of the model errors v_i and time-series and cross-sectional data (see Rao, 2003, Chapter 8).

In the second type, called basic unit level model, unit level auxiliary variables $x_{ij} = (x_{i1j}, \dots, x_{ipij})^T$ are related to the unit y-values y_{ij} through a nested error linear regression model $y_{ij} = x_{ij}^T \beta + v_i + e_{ij}$, where $v_i \stackrel{iid}{\sim} N(0, \sigma_v^2)$ and independent of $e_{ij} \stackrel{iid}{\sim} N(0, \sigma_e^2)$. Various extensions of the basic unit-level model have been proposed to handle binary responses, two-stage sampling within areas, multivariate responses and others (see Rao, 2003, Chapters 8, 9 and 10). For example, for binary responses y_{ij} , we may assume that $y_{ij} \stackrel{ind}{\sim} \text{Bernoulli}(p_{ij})$ and that the p_{ij} 's are linked by assuming a logistic regression model $\log\{p_{ij}/(1-p_{ij})\} = x_{ij}^T \beta + v_i$, where $v_i \stackrel{iid}{\sim} N(0, \sigma^2)$. This is a special case of generalized linear mixed models.

3. BASIC AREA LEVEL MODEL: EB

3.1 Estimation of θ_i

Under the basic area level model, the best estimator of θ_i in the sense of minimum MSE is given by $E(\theta_i | \hat{\theta}_i, \beta, \sigma_v^2)$ which depends on the model parameters β and σ_v^2 . Replacing (β, σ_v^2) by suitable estimators $(\hat{\beta}, \hat{\sigma}_v^2)$ obtained from the marginal distribution of $\hat{\theta}_i$'s, namely $\hat{\theta}_i \stackrel{ind}{\sim} N(z_i^T \hat{\beta}, \hat{\sigma}_v^2 + \psi_i)$, we obtain the empirical Bayes or empirical best (EB) estimator:

$$\hat{\theta}_i^{EB} = \hat{\gamma}_i \hat{\theta}_i + (1 - \hat{\gamma}_i) z_i^T \hat{\beta}, \quad (3.1)$$

where $\hat{\gamma}_i = \hat{\sigma}_v^2 / (\hat{\sigma}_v^2 + \psi_i)$. The form (3.1) shows that the EB estimator of θ_i is a weighted average of the direct estimator $\hat{\theta}_i$ and the regression synthetic estimator $z_i^T \hat{\beta}$ with weights $\hat{\gamma}_i$ and $1 - \hat{\gamma}_i$ respectively. The weight $\hat{\gamma}_i$ is a measure of between area variability relative to total variability associated with area i . The estimator $\hat{\theta}_i^{EB}$ is unbiased for θ_i under the combined model, i.e., $E(\hat{\theta}_i^{EB} - \theta_i) = 0$, and design-consistent as the sampling variance ψ_i goes to zero, provided the direct estimator is design-consistent. However, the resulting estimator $g^{-1}(\hat{\theta}_i^{EB})$ of Y_i is biased if $g(\cdot)$ is non linear. Note that $g^{-1}(\hat{\theta}_i^{EB})$ is not equal to the EB estimator \hat{Y}_i^{EB} obtained by evaluating $E[g^{-1}(\theta_i) | \hat{\theta}_i, \beta, \sigma_v^2]$ at $\hat{\beta}$ and $\hat{\sigma}_v^2$.

Under normality assumption, maximum likelihood (ML) or residual maximum likelihood (REML) method may be used to estimate β and σ_v^2 from the marginal distribution $\hat{\theta}_i \stackrel{iid}{\sim} N(z_i^T \beta, \sigma_v^2 + \psi_i)$. Alternatively, σ_v^2 may be estimated by a simple method of moments (Prasad and Rao, 1990) or by solving the following moment equation iteratively for σ_v^2 (Fay and Herriot, 1979):

$$a(\sigma_v^2) = \sum_{i=1}^m (\hat{\theta}_i - z_i^T \tilde{\beta}(\sigma_v^2))^2 / (\sigma_v^2 + \psi_i) = m - p, \quad (3.2)$$

where $\tilde{\beta}(\sigma_v^2)$ is the weighted least squares estimator of β for given σ_v^2 . The resulting estimators $\hat{\sigma}_v^2$ and $\hat{\beta} = \tilde{\beta}(\hat{\sigma}_v^2)$ lead to the EBLUP estimator of θ_i from (3.1). The EBLUP estimator does not depend on normality.

3.2 MSE Estimation

Methods of estimating $MSE(\hat{\theta}_i^{EB})$ that account for the variability of $\hat{\beta}$ and $\hat{\sigma}_v^2$ have been studied extensively in the literature, where $MSE(\hat{\theta}_i^{EB}) = E(\hat{\theta}_i^{EB} - \theta_i)^2$ and the expectation is with respect to the combined model (see Rao, 2003, Chapter 7). An accurate approximation to $MSE(\hat{\theta}_i^{EB})$ under normality is given by

$$MSE(\hat{\theta}_i^{EB}) \approx g_{1i}(\sigma_v^2) + g_{2i}(\sigma_v^2) + g_{3i}(\sigma_v^2) \quad (3.3)$$

where the leading term $g_{1i}(\sigma_v^2) = \gamma_i \psi_i$ with $\gamma_i = \sigma_v^2 / (\sigma_v^2 + \psi_i)$ is the contribution to MSE assuming β and σ_v^2 are known,

$$g_{2i}(\sigma_v^2) = (1 - \gamma_i)^2 z_i^T \left[\sum_{i=1}^m z_i z_i^T / (\sigma_v^2 + \psi_i) \right]^{-1} z_i \quad (3.4)$$

accounts for the variability of $\hat{\beta}$ and the term

$$g_{3i}(\sigma_v^2) = \left[\psi_i^2 / (\sigma_v^2 + \psi_i)^4 \right] E(\hat{\theta}_i - z_i^T \beta)^2 h(\sigma_v^2) \quad (3.5)$$

$$= \left[\psi_i^2 / (\sigma_v^2 + \psi_i)^3 \right] h(\sigma_v^2) \quad (3.6)$$

accounts for the variability of $\hat{\sigma}_v^2$, where $h(\sigma_v^2)$ is the asymptotic variance of $\hat{\sigma}_v^2$ for large m . Neglected terms in the approximation (3.3) are of lower order than m^{-1} , and the approximation is valid for the Prasad-Rao (PR), Fay-Herriot (FH), ML and REML methods of estimating σ_v^2 .

Comparing the leading term $\gamma_i \psi_i$ of (3.3) with ψ_i , the MSE of the direct estimator $\hat{\theta}_i$, it is clear that the EB estimator $\hat{\theta}_i^{EB}$ leads to large gain in efficiency when γ_i is small, i.e., when σ_v^2 , the variability of the model errors v_i , is small relative to the total variability, $\sigma_v^2 + \psi_i$. Note that ψ_i is also the design variance of $\hat{\theta}_i$.

Turning to MSE estimation, an estimator correct to the same order approximation as (3.3) is given by

$$mse(\hat{\theta}_i^{EB}) \approx g_{1i}(\hat{\sigma}_v^2) + g_{2i}(\hat{\sigma}_v^2) + 2g_{3i}(\hat{\sigma}_v^2). \quad (3.7)$$

The estimator (3.7) is approximately unbiased for $MSE(\hat{\theta}_i^{EB})$ in the sense that its bias is of lower order than m^{-1} , provided $\hat{\sigma}_v^2$ is based on REML or PR. For ML and FH methods of estimating σ_v^2 , an extra term $g_{0i}(\hat{\sigma}_v^2)$ is added to (3.8). This extra term for ML is positive (Datta and Lahiri, 2000). Therefore, ignoring this term and using (3.8) with ML estimator $\hat{\sigma}_v^2$ would lead to underestimation of MSE. On the other hand, the extra term for FH is negative (Datta, Rao and Smith, 2002). Therefore, ignoring this term and using (3.8) with FH estimator $\hat{\sigma}_v^2$ would lead to overestimation of MSE.

Lahiri and Rao (1995) showed that the MSE estimator (3.8) using the PR estimator of σ_v^2 is robust to nonnormality of the random effects in the sense that approximate unbiasedness remains valid, provided the normality of the sampling errors, e_i , holds. The latter assumption is less restrictive than the normality of the v_i 's because of the central limit theorem effect on the direct estimators $\hat{\theta}_i$. It is not known if the robustness property is also valid under REML, ML and FH methods.

A criticism of the MSE estimator (3.8) and its modification for ML and FH is that it is not area-specific in the sense that it does not explicitly depend on $\hat{\theta}_i$ although the area-specific auxiliary data z_i is involved in the $g_{2i}(\hat{\sigma}_v^2)$ -term. Rao (2000) used the expression (3.5) for $g_{3i}(\sigma_v^2)$ to get an alternative area-specific estimator of $g_{3i}(\sigma_v^2)$:

$$\tilde{g}_{3i}(\hat{\sigma}_v^2, \hat{\theta}_i) = \left[\psi_i^2 / (\sigma_v^2 + \psi_i)^4 \right] (\hat{\theta}_i - z_i^T \hat{\beta})^2 h(\hat{\sigma}_v^2). \quad (3.8)$$

Using (3.8), we get two different area-specific MSE estimators for REML or PR:

$$mse_1(\hat{\theta}_i^{EB}) \approx g_{1i}(\hat{\sigma}_v^2) + g_{2i}(\hat{\sigma}_v^2) + g_{3i}(\hat{\sigma}_v^2) + \tilde{g}_{3i}(\hat{\sigma}_v^2, \hat{\theta}_i) \quad (3.9)$$

and

$$mse_2(\hat{\theta}_i^{EB}) \approx g_{1i}(\hat{\sigma}_v^2) + g_{2i}(\hat{\sigma}_v^2) + 2\tilde{g}_{3i}(\hat{\sigma}_v^2, \hat{\theta}_i). \quad (3.10)$$

The term $\tilde{g}_{3i}(\hat{\sigma}_v^2, \hat{\theta}_i)$ is less stable than $g_{3i}(\hat{\sigma}_v^2)$ but it is of lower order than the leading term $g_{1i}(\hat{\sigma}_v^2)$ in (3.9) and (3.10). As a result, the coefficient of variation (CV) of $mse_1(\hat{\theta}_i^{EB})$ should be comparable to the CV of $mse(\hat{\theta}_i^{EB})$, at least for moderate to large m . Fuller (1989) estimated the conditional MSE of $\hat{\theta}_i^{EB}$ given the i -th area direct

estimator $\hat{\theta}_i$. His area-specific MSE estimator is closely related to the unconditional MSE estimator (3.9). Butar and Lahiri (1997) obtained an area-specific MSE estimator by correcting the bias of Laird and Louis' (1987) estimator of variability based on the parametric bootstrap method. This bias-corrected MSE estimator is identical to (3.9) which is obtained in a straightforward manner from the formula (3.5) for $g_{3i}(\sigma_v^2)$.

It is more appealing to survey practitioners to consider the estimation of sampling MSE of $\hat{\theta}_i^{EB}$, i.e., $MSE_p(\hat{\theta}_i^{EB}) = E_p(\hat{\theta}_i^{EB} - \theta_i)^2$, where the expectation E_p is with respect to the sampling design $p(\cdot)$, i.e., the distribution of sampling errors given the θ_i 's. Rivest and Belmonte (2000) derived a design-unbiased estimator of $MSE(\hat{\theta}_i^{EB})$ using the PR estimator of σ_v^2 . The leading term of this MSE estimator is area-specific, i.e., depends on $\hat{\theta}_i$, unlike the leading term $g_{li}(\hat{\sigma}_v^2)$ of the model-based MSE estimator, $mse(\hat{\theta}_i^{EB})$. However, it is highly unstable relative to $mse(\hat{\theta}_i^{EB})$ unless more weight is attached to the direct estimator $\hat{\theta}_i$, i.e., $1 - \hat{\gamma}_i$ is small.

3.3 Unknown Sampling Variances ψ_i

In sections 4.1 and 4.2 we assumed that the sampling variances, ψ_i , are known, but this is a restrictive assumption. Wang and Fuller (2003) and Rivest and Vandal (2002) studied the effect of estimating ψ_i on the MSE of the EB estimator (4.1) with $\hat{\gamma}_i$ replaced by $\hat{\sigma}_v^2 / (\hat{\sigma}_v^2 + \hat{\psi}_i)$, where $\hat{\psi}_i$ is an estimator of ψ_i . For example, suppose that we have a random sample $y_{ij} \stackrel{iid}{\sim} N(\theta_i, \sigma^2)$, $j = 1, \dots, n_i$ (≥ 2) from the i -th area and $\hat{\theta}_i = \bar{y}_i$, the sample mean. In this case $\hat{\psi}_i = s_i^2 / n_i$ is design-unbiased for ψ_i , where s_i^2 is the sample variance. Further, \bar{y}_i and $\hat{\psi}_i$ are independently distributed with $\hat{\psi}_i \approx N[\psi_i, \delta_i = 2\psi_i^2 / (n_i - 1)]$. Under this set-up, Rivest and Vandal (2002) obtained an appropriate MSE estimator by adding the term $2\hat{\delta}_i \hat{\sigma}_v^4 / (\hat{\psi}_i + \hat{\sigma}_v^2)^3$ to (3.8) to account for the estimation of ψ_i , where $\hat{\delta}_i = 2\hat{\psi}_i^2 / (n_i - 1)$. If the sample sizes n_i are small, then (3.8) can underestimate the MSE quite severely, unlike the Rivest-Vandal MSE estimator. If $\hat{\psi}_i$ is a smoothed estimator of ψ_i based on GVF model fitting, the contribution from the extra term is of the same order, $O(m^{-1})$, as the g_{3i} -term.

4. JACKKNIFE ESTIMATION OF MSE

Jiang, Lahiri and Wan (2002) proposed a jackknife method of estimating MSE of EB estimators that is applicable to generalized linear mixed models with block diagonal covariance structures, where the blocks correspond to small areas. This method also leads to approximately unbiased estimators of MSE of EB estimators. For example, consider the case of binary responses $y_{ij} \stackrel{iid}{\sim} \text{Bernoulli}(p_i)$, $j = 1, \dots, n_i$ and $\log\{p_i / (1 - p_i)\} = z_i^T \beta + v_i$, $i = 1, \dots, m$, where z_i is the vector of area-specific covariates, $v_i \stackrel{iid}{\sim} N(0, \sigma_v^2)$ and p_i is the i -th area proportion. The minimum MSE (or Bayes) estimator of p_i is given by $\hat{p}_i^B = E(p_i | y_i, \beta, \sigma_v^2) =: k(y_i, \beta, \sigma_v^2)$, where $y_i = \sum_j y_{ij}$. The EB estimator of p_i is $\hat{p}_i^{EB} = k(y_i, \hat{\beta}, \hat{\sigma}_v^2)$, where $\hat{\beta}$ and $\hat{\sigma}_v^2$ are suitable estimators of β and σ_v^2 obtained from the marginal distribution of y_i 's.

The jackknife method makes use of the following orthogonal decomposition of $MSE(\hat{p}_i^{EB})$:

$$MSE(\hat{p}_i^{EB}) = E(\hat{p}_i^B - p_i)^2 + E(\hat{p}_i^{EB} - \hat{p}_i^B)^2. \quad (4.1)$$

Based on the decomposition (4.1), Jiang et al. (2002) proposed the following jackknife steps to estimate $MSE(\hat{\theta}_i^{EB})$:

(1) Calculate $\hat{\beta}(\ell)$ and $\hat{\sigma}^2(\ell)$ deleting the R th area data (y_i, x_i) . Let $\hat{p}_i^{EB}(\ell) = k[y_i, \hat{\beta}(\ell), \hat{\sigma}^2(\ell)]$ be the EB estimator of p_i based on $\hat{\beta}(\ell)$ and $\hat{\sigma}^2(\ell)$. Note that y_i remains unchanged.

(2) Calculate the jackknife estimator of the last term in (4.1) as

$$\hat{M}_{2i} = \frac{m-1}{m} \sum_{\ell=1}^m [\hat{p}_i^{EB}(\ell) - \hat{p}_i^{EB}]^2. \quad (4.2)$$

(3) The first term $E(\hat{p}_i^B - p_i)^2$ may be written as $E[\tilde{g}_{li}(y_i, \beta, \sigma_v^2)] = g_{li}(\beta, \sigma_v^2)$, where $\tilde{g}_{li}(y_i, \beta, \sigma_v^2) = V(p_i | y_i, \beta, \sigma_v^2)$ is the posterior variance of p_i given y_i and (β, σ_v^2) . Adjust the bias of $g_{li}(\hat{\beta}, \hat{\sigma}_v^2)$ (as an estimator of $g_{li}(\beta, \sigma_v^2)$) using the jackknife bias reduction method. The bias-adjusted estimator is given by

$$\hat{M}_{1i} = g_{li}(\hat{\beta}, \hat{\sigma}_v^2) - \frac{m-1}{m} \sum_{\ell=1}^m [g_{li}(\hat{\beta}(\ell), \hat{\sigma}_v^2(\ell)) - g_{li}(\hat{\beta}, \hat{\sigma}_v^2)]. \quad (4.3)$$

Note that the leading term $g_{li}(\hat{\beta}, \hat{\sigma}_v^2)$ is not area-specific in the sense of not depending on y_i .

(4) Calculate the jackknife estimator of MSE as

$$mse_J(\hat{p}_i^{EB}) = \hat{M}_{1i} + \hat{M}_{2i}. \quad (4.4)$$

Booth and Hobert (1998) argued that for non-normal responses the MSE estimator should be area-specific because the posterior variance of p_i given (β, σ_v^2) depends on y_i , unlike the case of linear mixed models. They proposed the conditional MSE, given the i -th area data (y_i, z_i) , as the relevant measure of variability, and estimated the conditional MSE. Rao (2003, Chapter 9) addressed this criticism by simply modifying the bias-adjusted estimator \hat{M}_{1i} . Instead of evaluating the expectation of $\tilde{g}_{li}(y_i, \beta, \sigma_v^2)$ with respect to the marginal distribution of y_i (using numerical integration), he proposed to adjust the bias of $\tilde{g}_{li}(y_i, \hat{\beta}, \hat{\sigma}_v^2)$ as an estimator of $g_{li}(y_i, \beta, \sigma_v^2)$. This leads to

$$\tilde{M}_{1i}(y_i) = \tilde{g}_{li}(y_i, \hat{\beta}, \hat{\sigma}_v^2) - \frac{m-1}{m} \sum_{\ell=1}^m [\tilde{g}_{li}(y_i, \hat{\beta}(\ell), \hat{\sigma}_v^2(\ell)) - \tilde{g}_{li}(y_i, \hat{\beta}, \hat{\sigma}_v^2)] \quad (4.5)$$

which is area-specific including the leading term $\tilde{g}_{li}(y_i, \hat{\beta}, \hat{\sigma}_v^2)$. The modified jackknife estimator of MSE is given by

$$mse_J^*(\hat{p}_i^{EB}) = \tilde{M}_{1i}(y_i) + \hat{M}_{2i}. \quad (4.6)$$

Note that (4.6) is not only area-specific but also computationally simpler than (4.4) because it avoids the evaluation of the expectation of $\tilde{g}_{li}(y_i, \beta, \sigma_v^2)$ with respect to the marginal distribution of y_i . For the case of a linear mixed model, $\tilde{M}_{1i}(y_i) = \hat{M}_{1i}$ and hence (4.6) is identical to (4.4).

5. BASIC UNIT LEVEL MODEL: PSEUDO-EB

We now turn to the basic unit level model, $y_{ij} = x_{ij}^T \beta + v_i + e_{ij}$, and assume that the model holds for the sample, i.e., no sample selection bias within areas. If the areas are also sampled, then we assume the absence of selection bias for sampled areas as well. The mean for the i -th area, \bar{Y}_i , may be approximated as $\mu_i = \bar{X}_i^T \beta + v_i$, assuming that the number of population units in the i -th area, N_i , is large, where \bar{X}_i is the population mean of x for the i -th area. Assuming $v_i \stackrel{iid}{\sim} N(0, \sigma_v^2)$ and independent of $e_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$ and estimating the model parameters β and σ_v^2 from the marginal distribution of the sampled y_{ij} 's, we get the EB estimator of μ_i as

$$\hat{\mu}_i^{EB} = \hat{\gamma}_i \left[\bar{y}_i + (\bar{X}_i - \bar{x}_i)^T \hat{\beta} \right] + (1 - \hat{\gamma}_i) \bar{X}_i^T \hat{\beta}, \quad (5.1)$$

where $\hat{\gamma}_i = \hat{\sigma}_v^2 / (\hat{\sigma}_v^2 + \hat{\sigma}_e^2 / n_i)$, (\bar{y}_i, \bar{x}_i) are the i -th area sample means and $(\hat{\beta}, \hat{\sigma}_v^2)$ are the estimators of (β, σ_v^2) ; see Battese, Harter and Fuller (1988). This estimator is also the EBLUP estimator without normality assumption, provided β and σ_v^2 are estimated by a moments method such as the method of fitting-of-constants. As $n_i \rightarrow \infty$, the EB estimator converges to the "survey regression" estimator $\bar{y}_i + (\bar{X}_i - \bar{x}_i)^T \hat{\beta}$, while it converges to the "synthetic regression" estimator $\bar{X}_i^T \hat{\beta}$ as n_i decreases. We refer to Rao (2003), Chapter 7 for a detailed discussion of EBLUP estimation and MSE estimation. A drawback of (5.1) is that it is purely model-based and does not take account of the survey weights w_{ij} . As a result, it is not design-consistent as n_i increases, unless the sampling design is self-weighting within areas, i.e., $w_{ij} = w_i$. On the other hand, the EBLUP estimator under the area level model is design-consistent. It is desirable to ensure design-consistency for the unit level models as well because n_i could be moderately large for some of the areas under consideration. Also, it is desirable to ensure that the estimators of area totals automatically add up to the direct survey regression estimator of the large area total. You and Rao (2002a) developed a pseudo-EBLUP estimator of μ_i that satisfies both the desirable properties.

We have assumed that the random small area effects v_i are normally distributed in the basic unit level model $y_{ij} = x_{ij}^T \beta + v_i + e_{ij}$. The MSE estimator of the purely model-based estimator $\hat{\mu}_i^{EB}$ derived under normality of the v_i 's is not robust to deviations from normality, unlike the case of basic area level model studied in Section 3.2. It would be worthwhile to study EB inference under semi-nonparametric (SNP) representations of the density of v_i . Zhang and Davidian (2001) approximated the density of v_i by a SNP representation which includes normality as a special case, and it provides flexibility in capturing nonnormality through a user-chosen tuning parameter. Maiti (2001) used a finite mixture of normal distributions for the distribution of v_i , and developed hierarchical Bayes (HB) estimates of small area means, assuming a prior distribution on the model parameters. EB estimation of small area means and associated MSE estimation under broad classes of densities of v_i , such as the above, would be practically useful.

6. HIERARCHICAL BAYES (HB) APPROACH

We illustrate the HB approach using the basic area level model of Section 2. Under this approach, a prior distribution on the model parameters $\delta = (\beta^T, \sigma_v^2)$ is specified, and inferences are then based on the posterior distribution, $f(\theta_i | \hat{\theta})$, of θ_i given the data $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_m)^T$. In particular, θ_i is estimated by its posterior mean $E(\theta_i | \hat{\theta})$, called the HB estimator $\hat{\theta}_i^{HB}$. The variability of $\hat{\theta}_i^{HB}$ is measured by the posterior variance $V(\theta_i | \hat{\theta})$. The calculation of $\hat{\theta}_i^{HB}$ and $V(\theta_i | \hat{\theta})$ involves integrations with respect to the posterior distribution of β, σ_v^2 ,

$f(\beta, \sigma_v^2 | \hat{\theta})$. However, Monte Carlo Markov chain (MCMC) methods can be used to generate J simulated samples $\{\theta_1^{(j)}, \dots, \theta_m^{(j)}; j = 1, \dots, J\}$ directly from the joint posterior $f(\theta | \hat{\theta})$, where $\theta = (\theta_1, \dots, \theta_m)^T$ and J is sufficiently large. Using the simulated samples, we approximate $\hat{\theta}_i^{HB}$ by the mean $\theta_i^{(\cdot)} = J^{-1} \sum_j \theta_i^{(j)}$ and $V(\theta_i | \hat{\theta})$ by the variance $J^{-1} \sum_j (\theta_i^{(j)} - \theta_i^{(\cdot)})^2$ of the simulated samples $\theta_i^{(j)}$. The HB estimator of the total Y_i is approximated by the mean $Y_i^{(\cdot)} = J^{-1} \sum_j Y_i^{(j)}$ and the posterior variance of Y_i , $V(Y_i | \hat{\theta})$ by the variance $J^{-1} \sum_j (Y_i^{(j)} - Y_i^{(\cdot)})^2$, where $Y_i^{(j)} = g^{-1}(\theta_i^{(j)})$.

An advantage of the HB approach is that it is straightforward, the inferences are “exact” unlike the EB (or EBLUP) approach, and it can handle complex small area models using MCMC methods, but it requires the specification of a prior $f(\beta, \sigma_v^2)$ on the model parameters. It would be desirable to select a “matching” prior $f(\beta, \sigma_v^2) \propto f(\sigma_v^2)$ that leads to well-calibrated inferences. In particular, the posterior variance should be approximately unbiased for $MSE(\hat{\theta}_i^{HB})$, i.e., $E[V(\theta_i | \hat{\theta})] - MSE(\theta_i^{HB}) = o(m^{-1})$; asymptotically, $\hat{\theta}_i^{HB} \approx \hat{\theta}_i^{EB}$. This will provide a frequentist justification for the posterior variance as a measure of variability. Data, Rao and Smith (2002) showed that the matching prior is given by

$$f_i(\sigma_v^2) \propto (\sigma_v^2 + \psi_i)^2 \sum_{\ell=1}^m (\sigma_v^2 + \psi_\ell)^{-2}. \quad (6.1)$$

This prior depends collectively on the sampling variances, ψ_ℓ , for all the areas R as well as on the area-specific sampling variance ψ_i . For the balances case, $\psi_i = \psi$, the matching prior reduces to the “flat” prior $f(\sigma_v^2) \propto 1$. Note that the prior (6.1) on the common parameter σ_v^2 is designed for inference on the i -th area so that its dependence on ψ_i may not be problematic.

A disadvantage of the EB estimator $\hat{\theta}_i^{EB}$ is that the weight $\hat{\gamma}_i$ attached to the direct estimator takes zero value when $\hat{\sigma}_v^2 = 0$, in which case it reduces to the regression synthetic estimators $z_i^T \hat{\beta}$. Thus, all the direct estimator $\hat{\theta}_i$ receive zero weight even when the sample sizes in some areas are not small. This difficulty was encountered in using a state model to produce EB state estimates of poor school-age children in the United States (National Research Council, 2000). The HB approach avoids this difficulty by producing positive weights in all cases. Bell (1999) applied the HB approach to the state model using the prior $f(\beta, \sigma_v^2) = f(\beta)f(\sigma_v^2)$ with $f(\beta) \propto 1$ and $f(\sigma_v^2) \propto 1$, and obtained positive weights in all cases. But it is not clear if his method leads to well-calibrated inferences since the matching prior (6.1) is different from the flat prior, especially when the ψ_i -values vary significantly, as in the case of Bell (1999) with $\max(\psi_i) / \min(\psi_i)$ as large as 20.

You and Rao (2002b) used the HB approach to handle the case of unmatched sampling and linking area level models (Section 2) and applied it to Canadian census undercount estimation. In this application, C_i = census count, Y_i = number missing and \hat{Y}_i is a post-census survey estimator of Y_i with known sampling variance σ_i^2 for the i -th province in Canada ($i = 1, \dots, m = 10$). The σ_i^2 's were estimated by fitting a GVF model of the form $V(\hat{Y}_i) \propto C_i^\gamma$ and then treating as if known in the sampling model $\hat{Y}_i | Y_i \stackrel{ind}{\sim} N(Y_i, \sigma_i^2)$. The linking model is given by $\theta_i = \log\{Y_i / (Y_i + C_i)\} = \beta_0 + \beta_1 \log C_i + v_i$ with $v_i \stackrel{iid}{\sim} N(0, \sigma_v^2)$. HB estimates of undercounts, Y_i , and undercoverage rates, $U_i = Y_i / (Y_i + C_i)$, and associated coefficients of variation (based on the posterior variance) were calculated, using MCMC methods.

Singh, Folsom and Vaish (2003) studied the basic unit level model, $y_{ij} = x_{ij}^T \beta + v_i + e_{ij}$ with $v_i \stackrel{iid}{\sim} N(0, \sigma_v^2)$ independent of $e_{ij} \stackrel{iid}{\sim} N(0, \sigma_e^2)$, for the population, and allowed sample selection bias within small areas. They used methods based on survey-weighted estimating functions (EF) to account for the sample selection bias. They also extended the method to generalized linear mixed model such as the logistic mixed model $y_{ij} | p_{ij} \stackrel{ind}{\sim} \text{Bernoulli}(p_{ij})$ and logit $(p_{ij}) = \log\{p_{ij}/(1-p_{ij})\} = x_{ij}^T \beta + v_i$; see Rao (2003, pp. 253-254) for a brief account of the EF method. In situations where a sample of areas is selected, the random effects v_i are assumed to be free of sample selection bias.

The HB approach is powerful and attractive, but caution should be exercised when using MCMC methods. For example, MCMC algorithms could lead to seemingly reasonable inferences about a non-existent posterior distribution. This happens when the posterior distribution is improper and yet all the Gibbs conditional distributions, used in generating the MCMC simulated samples, are proper (Hobert and Casella, 1996). Another difficulty with MCMC is that the convergence diagnostic tools can fail to detect the sort of convergence failures they were designed to identify (Cowles and Carlin, 1996). We refer the reader to Rao (2003), Section 10.2.4 for a discussion of practical issues associated with MCMC.

MCMC methods are also extensively used for model determination which plays a vital role in developing model-based small area estimates. In particular, methods based on Bayes factors, posterior predictive densities and cross-validation predictive densities are employed for model determination; see Rao (2003), Section 10.2.6. The criterion of posterior predictive probability is often used to check the overall fit of a proposed model. Sinharay and Stern (2003) conducted a simulation study to investigate the effectiveness of this criterion for model checking, using the basic area level model with no covariates. Their study indicates that it is difficult to detect nonnormality of the random effects v_i using this criterion, unless the extent of violation is huge. Therefore, caution should be exercised in using HB criteria for model determination.

7. SOME PRACTICAL CONSIDERATIONS

In this section, we provide brief remarks on practical issues related to small area estimation.

(i) Design issues

It is important to consider design issues that have an impact on small area estimation. A proper resolution of design issues could lead to enhancement in the reliability of direct as well as indirect estimates for both planned and unplanned domains (areas). The following measures at the design stage might be useful in minimizing the need for indirect estimators, at least for some planned domains: Use of list frames to replace clusters wherever possible, use of many small strata from which samples are drawn, compromise sample allocations to satisfy reliability requirements at a small area level as well as large area level, integration of surveys through harmonizing questions across surveys of the same population, use of multiple frame surveys, use of “rolling samples” as a method of cumulating data over time. We refer the reader to Rao (2003), Chapter 2, Singh et al. (1994) and Marker (2001) for further details.

Despite the above preventive measures at the design stage, indirect estimates will be needed in practice because it is not possible to anticipate and plan for all possible areas (or domains) and uses of survey data: “the client will always require more than is specified at the design stage” (Fuller, 1999, p. 344).

(ii) Model selection and validation

Methodological developments and applications of model-based estimation are impressive, but caution should be exercised because of the model assumptions. Good auxiliary information related to the variables of interest plays a vital role in determining suitable linking models. Therefore, more attention should be given to the compilation of auxiliary variables that are good predictors of study variables. Subject matter specialists or end users should have

influence on the choice of models, particularly on the choice of auxiliary variables. However, model diagnostics should be used to find suitable model(s) that fit the data well. Such model diagnostics include residual analysis to detect departures from assumed models, selection of auxiliary variables and case-deletion diagnostics to detect influential observations. We refer the reader to Rao (2003), Chapter 6 for some methods of model validation in the frequentist framework.

Hierarchical Bayes (HB) approach has become very popular in recent years due to its ability to handle complex models using MCMC methods. However, caution should be exercised in the choice of improper priors on model parameters, as noted in Section 6. Limitations of MCMC methods, such as shortcomings of available convergence diagnostics, should also be noted. Carlin and Louis (2000) made an important observation on the dangers of “plug and play” implementation via MCMC: “Worse, the sheer power of MCMC methods has led to the temptation to fit models larger than the data can readily support without a strongly informative prior structure - now something of a rarity in applied Bayesian work”.

HB methods for model validation via MCMC have been extensively developed, but the effectiveness of some criteria for model checking is questionable as noted in Section 6. Further work on effective HB methods for model checking is needed.

(iii) Area level vs. unit level models

Area level models have wider scope than unit level models because area level auxiliary information is more readily available than unit level auxiliary data. Also, design weights are incorporated by modelling design-weighted direct estimators, and the resulting EB or HB estimators are design consistent. But the assumption of known sampling variances, ψ_i , is quite restrictive. Smoothed estimates of ψ_i 's based on GVF model fitting can also cause difficulties in MSE estimation, as noted in Section 3.3. We need more work on obtaining good approximations to the sampling variances as well as methods that incorporate the variability associated with estimated sampling variances in MSE estimation. This task becomes more difficult when using multivariate or time series area level models because sampling covariances are also needed.

Recent work on incorporating survey weights into model-based estimation, via pseudo-EB or pseudo-HB, is promising; in particular, the self-benchmarking property noted in Section 5. But the assumption that the sample selection bias is absent may not be true for some applications. The estimating functions (EF) approach of Singh, Folsom and Vaish (2003), mentioned in Section 6, allows sample selection bias within sampled areas but it assumes that the random effects v_i are free from sample selection bias in situations where a sample of areas is selected. Methods for handling the latter case are needed. Moreover, their method assumes known sampling variances, as in the area level model, and this assumption may be restrictive (see Rao, 2003, Section 10.5.4).

(iv) “Triple-goal” estimation

We focussed on model-based estimation of small area totals or means, but such estimates may not be suitable if the main objective is to produce an ensemble of parameter estimates whose distribution is in some sense close enough to the distribution of area-specific parameters. For example, we may be interested in ranking areas or identifying areas that fall below or above some prespecified level. Shen and Louis (1998) proposed “triple-goal” estimators that can produce good ranks, a good histogram and good area-specific estimators, assuming simple linking models. It would be useful to extend their methods to handle more complex models that are suitable for small area estimation.

(v) Nonsampling errors

We have assumed the absence of measurement errors in the responses and/or the covariates as well as nonresponse. But nonsampling errors can have a substantial effect on small area estimation, and it would be useful to develop suitable designs as well as methods of estimation that can account for nonsampling errors. Nandram and Choi (2002) used HB nonresponse models for binary data, and applied the theory to data from the U.S. National Crime Survey to estimate small area proportions. Measurement errors in the responses, even under an additivity assumption, can lead to biased estimators of quantiles and histograms. In the context of direct estimation, Fuller (1995) proposed methods at the design stage that can lead to bias-adjusted estimators of quantiles and histograms.

REFERENCES

- Battese, G.E., Harter, R.M., Fuller, W.A. (1988). An Error-Components Model for Prediction of Crop Areas Using Survey and Satellite Data, *Journal of the American Statistical Association*, **83**, 28-36.
- Bell, W.R. (1999). Accounting for Uncertainty About Variances in Small Area Estimation. *Bulletin of the International Statistical Institute*.
- Booth, J.G. and Hobert, J.P. (1998). Standard Errors of Predictors in Generalized Linear Mixed Models. *Journal of the American Statistical Association*, **93**, 262-272.
- Butar, F.B. and Lahiri, P. (2001). On Measures of Uncertainty of Empirical Bayes Small-Area Estimators. Technical Report, Department of Statistics, University of Nebraska-Lincoln.
- Carlin, B.P. and Louis, T.A. (2000). Empirical Bayes: Past, Present and Future. *Journal of the American Statistical Association*, **95**, 1286-1289.
- Cowles, M.R. and Carlin, B.P. (1996). Markov Chain Monte Carlo Convergence Diagnostics: A Comparative Review. *Journal of the American Statistical Association*, **91**, 883-904.
- Datta, G.S. and Lahiri, P. (2000). A Unified Measure of Uncertainty of Estimated Best Linear Unbiased Predictors in Small Area Estimation Problems. *Statistica Sinica*, **10**, 613-627.
- Datta, G.S., Rao, J.N.K. and Smith, D.D. (2002). On Measures of Uncertainty of Small Area Estimators in the Fay-Herriot Model, Technical Report, University of Georgia, Athens.
- Fay, R.E. and Herriot, R.A. (1979). Estimation of Income for Small Places: An Application of James-Stein Procedures to Census Data. *Journal of the American Statistical Association*, **74**, 269-277.
- Fuller, W.A. (1989). Prediction of True Values for the Measurement Error Model, in *Conference on Statistical Analysis of Measurement Error Models and Applications*, Humboldt State University.
- Fuller, W.A. (1995). Estimation in the Presence of Measurement Error. *International Statistical Review*, **63**, 121-147.
- Fuller, W.A. (1999). Environmental Surveys Over Time. *Journal of Agricultural, Biological and Environmental Statistics*, **4**, 331-345.
- Hobert, J.P. and Casella, G. (1996). The Effect of Improper Priors on Gibbs Sampling in Hierarchical Linear Mixed Models. *Journal of the American Statistical Association*, **91**, 1461-1473.
- Jiang, J., Lahiri, P. and Wan, S.-M. (2002). A Unified Jackknife Theory for Empirical Best Prediction with M-Estimation. *Annals of Statistics*, **30**, 1782-1810.
- Lahiri, P. and Rao, J.N.K. (1995). Robust Estimation of Mean Squared Error of Small Area Estimators. *Journal of the American Statistical Association*, **82**, 758-766.
- Laird, N.M. and Louis, T.A. (1987). Empirical Bayes Confidence Intervals Based on Bootstrap Samples. *Journal of the American Statistical Association*, **82**, 739-750.
- Maiti, T. (2001). Robust Generalized Linear Mixed Models for Small Area Estimation. *Journal of Statistical Planning and Inference*, **98**, 225-238.
- Marker, D.A. (2001). Producing Small Area Estimates from National Surveys: Methods for Minimizing Use of Indirect Estimators. *Survey Methodology*, **27**, 183-188.

- Nandram, B. and Choi, J.W. (2002). Hierarchical Bayesian Nonresponse Models for Binary Data from Small Areas with Uncertainty About Ignorability. *Journal of the American Statistical Association*, **97**, 381-388.
- National Research Council (2000). *Small-Area Estimates of School-Age Children in Poverty: Evaluation of Current Methodology*, C.F. Citro and G. Kalton (Eds.), Committee on National Statistics, Washington, D.C.: National Academy Press.
- Prasad, N.G.N. and Rao, J.N.K. (1990). The Estimation of the Mean Squared Error of Small-Area Estimators, *Journal of the American Statistical Association*, **85**, 163-171.
- Rao, J.N.K. (2001). EB and EBLUP in Small Area Estimation, in S.E. Ahmed and N. Reid (Eds.); *Empirical Bayes and Likelihood Inference*. Lecture Notes in Statistics 148, New York: Springer, pp. 33-43.
- Rao, J.N.K. (2003). *Small Area Estimation*. New York: Wiley.
- Rivest, L-P. and Belmonte, E. (2000). A Conditional Mean Squared Error of Small Area Estimators. *Survey Methodology*, **26**, 67-78.
- Rivest, L-P. and Vandal, N. (2003). Mean Squared Error Estimation for Small Areas When the Small Area Variances are Estimated, in *Proceedings of the International Conference on Recent Advances in Survey Sampling*, Laboratory for Research in Statistics and Probability, Carleton University, Ottawa, Canada (in press).
- Singh, M.P., Gambino, J., Mantel, H.J. (1994). Issues and Strategies for Small Area Data, *Survey Methodology*, **20**, 3-22.
- Singh, A.C., Folsom, Jr., R.E. and Vaish, A.K. (2003). Estimating Function Based Approach to Hierarchical Bayes Small Area Estimation for Survey Data, in *Proceedings of the International Conference on Recent Advances in Survey Sampling*, Laboratory for Research in Statistics and Probability, Carleton University, Ottawa, Canada (in press).
- Shen, W. and Louis, T.A. (1998). Triple-Goal Estimates in Two-Stage Hierarchical Models. *Journal of the Royal Statistical Society, Series B*, **60**, 455-471.
- Sinharay, S. and Stern, H.S. (2003). Posterior Predictive Model Checking in Hierarchical Models. *Journal of Statistical Planning and Inference*, **111**, 209-221.
- You, Y. and Rao, J.N.K. (2002a). A Pseudo-Empirical Best Linear Unbiased Prediction Approach to Small Area Estimation Using Survey Weights. *Canadian Journal of Statistics*, **30**, 431-439.
- You, Y. and Rao, J.N.K. (2002b). Small Area Estimation Using Unmatched Sampling and Linking Models. *Canadian Journal of Statistics*, **30**, 3-15.
- Wang, J. and Fuller, W.A. (2003). The Mean Squared Error of Small Area Predictors Constructed With Estimated Area Variances. *Journal of the American Statistical Association*, **98**, 716-723.
- Zhang, D. and Davidian, M. (2001). Linear Mixed Models With Flexible Distributions of Random Effects for Longitudinal Data. *Biometrics*, **57**, 795-802.