

Building-block BLUPs for Aggregate Level Small Area Estimation for Survey Data

A.C. Singh¹ and P. Yuan²

¹ Center for Excellence in Survey Research, NORC at the University of Chicago, Chicago, IL 60603

² Human Resources and Skills Development Canada, Hull, QC K1A 0J9

Abstract

Often the choice of the level of aggregation in small area (SA) modeling is not governed by adequacy of modeling assumptions but by user needs which varies from user to user. This may have a serious impact on validity of the underlying exchangeability assumption of area-specific random effects in small area (SA) modeling for means. This problem is not likely to arise with unit (or individual) level models if all predictors or covariates for all units in every area are at the unit level and thus having similar predictive power across areas which is not the case with aggregate level models because areas vary considerably in size. However, unit level models are known to be difficult for taking the sampling design into account. As a compromise, we propose a building-block (B) model such that random effects at B-level sub-areas are more likely to be exchangeable, and model parameters are estimated at this level or at the group (G) level after grouping of B-level areas to avoid the problem of zero sample size or very unstable direct estimates. The target SA estimates are obtained as building-block best linear unbiased predictors (denoted by bBLUPs) using the target domain level model (derived from the B-level model) in the interest of providing automatically more weight to design consistent direct estimators because direct estimators at the target domain level are more precise than those at the building block level. Although bBLUPs are not optimal under the G-level model used for estimating fixed parameters, they provide a good compromise between estimation efficiency of SA estimators and reducing their absolute relative difference from direct estimators. An illustrative example is presented using Canadian Labour Force Survey data with provinces as target small areas and provincial economic regions classified by age and gender as building blocks.

Key Words: Building Blocks; Aggregate Level; BLUP at Target Level; Benchmarking

1. Introduction

The problem considered in this paper arose in the context of estimation of annual number of employed in three-digit occupation codes for each province in Canada using the monthly Labour Force survey (LFS); an example of a three digit occupation code is A39 for 'Managers in Manufacturing and Utilities'. The direct estimates of annual average at the province level are not very stable. For the year 2003, Table 1 shows the average annual estimates, standard error and coefficient of variation (CV) based on monthly LFS. In practice, it would be desirable to have low CVs (below 10%) but generally not more than 15%. There are several provinces for which CVs are near or higher than 15%

especially for the Atlantic provinces (NL, PEI, NS, and NB) and the province SK in the prairies. Treating provinces as small areas (SA), it would seem natural to apply the Fay-Herriot aggregate level SA model under certain assumptions to produce more efficient small area estimates (SAE) for each province for each three digit occupation code. The SAE will be a linear combination of the direct estimate and the indirect (or synthetic) estimate obtained under the model; the combining coefficient depends on the variance of the two estimates and assigns higher weight to the estimate with smaller variance.

In the Fay-Herriot aggregate level model, besides assuming that the covariate totals are known at the aggregate or area level, and the number of small areas is large enough for model fitting, a key assumption is that of the exchangeability of random effects used in defining the area-specific means; i.e., they are independent with common mean and variance. In our application, the number of small areas is only 10 which is clearly not large enough for precise estimation of model parameters consisting of first order parameters of regression coefficients and second order parameters of variances and covariances of random effects. This problem could have been alleviated if the user had defined SAs at lower levels such as economic regions (which are 73 in total) or even lower subprovincial domains (economic region by age by gender). In any such formulation with lower level areas, the exchangeability consideration is crucial in defining the aggregate model level. If SAs are defined such that they are generally more similar in their population size, and if all covariates are at the SA level, then the residuals (or random effects) in the model means are likely to be exchangeable. The reason for this is that the predictive power of covariates varies with the model level; it is better at lower levels as the covariates can better track the variation in the dependent variable.

The above observation raises a fundamental question of the choice of a suitable level of aggregation in the Fay-Herriot model (such as economic regions or other subprovincial regions in Canada, or MSA or State or some other group of counties or census tracts in the US) because the crucial assumption of exchangeability of random effects may not be valid for any ad hoc choice of the area level based on user needs which incidentally may change from user to user. In other words, if the chosen area level (such as province or state) varies quite a bit between areas, and since the predictive power of area level covariates varies with area sizes, it would be difficult to justify exchangeability of model errors over areas. On the other hand, if one were to model at the unit level (i.e., at the individual person level), the problem is much less severe. However, with unit level models it is in general quite difficult to properly take account of the sampling design because there is no corresponding finite population quantity that a unit level weighted estimate represents. For this reason, we restrict ourselves to aggregate level models.

It is reasonable to stipulate that the lower the level of aggregation, the more likely it would be for the exchangeability assumption to hold. We therefore propose the idea of building blocks for aggregate level modeling which stipulates working with approximately comparable building blocks (B) such as subprovincial domains (economic region by age by gender) in Canada or subcounties in US (county-level modeling might be adequate in practice although it might be better to create sub-counties which are more

comparable in terms of population size). In practice, creation of building blocks is subject to the restriction that subpopulation totals of covariates at the B-level are available. The underlying model is defined at the B-level but it is fitted by grouping building blocks (within the higher target SA level such as province in Canada or MSA or state in US) such that direct estimates can be constructed at the group (G) level, which are considerably more stable than the B-level estimates. Grouping is based on a priori considerations in partitioning the target population (such as geo-demographic stratification in some hierarchical manner) so that unbiased design-based estimates at the group level can be obtained via domain estimation techniques and that the number of groups is large enough to yield adequate degrees of freedom for model fitting. The choice of groups should not be data-driven to select building blocks with little or no sample for grouping. Once the model parameter estimates are obtained from G-level modeling, SAE at the G-level (or even at the B-level) can be obtained using best linear unbiased prediction (BLUP) theory to be denoted as bBLUPs for building block BLUPs. However, for the desired target small areas or domains (D) which are assumed to be at a level much higher than the B-level, it may be preferable not to use bBLUP estimates obtained from the G-level, but to adjust the direct estimate at the target domain (D) level using bBLUP; i.e., by estimating random effects at D-level but fixed parameters at G-level. This way, the resulting estimate is less likely to be dominated by synthetic or model based estimate because there would be much less shrinkage in general of random effects to zero using direct D-level estimates than using G-level direct estimates. The resulting bBLUP estimator, although not optimal under the G-level model, might be more desirable in practice than the bBLUP at G-level as it is expected to be closer to the direct estimator and more robust to model misspecification.

In any application, it may not be feasible to choose the best level of building blocks (in the sense of comparable population size) due to practical constraints such as the availability of covariate totals at the B-level and the need for building blocks not to cut across SAs. Regardless, any serious consideration in choosing suitable and practically feasible building blocks is bound to lead to a more valid model than the one defined at an arbitrary SA level simply to satisfy user needs. The SA modeling formulation at the B-level is formally defined in Section 2 as well as the grouping of building blocks for model fitting which may be needed if the variance-covariance matrix at the B-level of direct estimates is very unstable due to very small or zero sample size. Incidentally, the number of parameters in the G-level model is still governed by the B-level model even though a reduced model via grouping is used for parameter estimation. In Section 3, estimation of B-level fixed model parameters using the G-level model is considered, and that of random effect parameters using either the G-level or the D-level model. Once model parameters are estimated, bBLUP estimation of SA parameters, and their mean squared error (MSE) estimation are considered in Section 4. Although the MSE for bBLUP is nonstandard due to difference in modeling levels, the second order adjustment to MSE could still be obtained along the lines of Prasad and Rao (1990). Next, modifications of bBLUP estimates of SAs under benchmark constraints and corresponding MSE are presented in Section 5. Empirical results from an application to LFS are presented in

Section 6 which, as expected, show that benchmarked SAE-D (the extension D indicates that random effects are estimated at D-level) tend to have slightly higher CVs (defined as relative root MSE) than the benchmarked SAE-G (when random effects are estimated at the G-level), but SAE-Ds are, in general, appreciably closer to direct estimates in view of the possibility of more over-shrinkage of SAE-Gs. Finally, concluding remarks and a discussion of related applications are presented in Section 7.

2. The Building Block and Related Models

The building block (B) model is a Fay-Herriot aggregate level model but at a very low level—as close as possible to the unit level except for the restriction that the covariate totals at B-level should remain available and building blocks be nested within target SAs. It can be expressed as a two stage model—sampling model for the estimated total and the linking model which links building block domain level parameters at the B-level (although the target SA parameters are at a higher domain or D-level such as the province level) to model parameters as follows.

$$\text{B-level Sampling Model: } \mathbf{y}_B = \boldsymbol{\theta}_B + \mathbf{e}_B$$

$$\text{B-level Linking Model: } \boldsymbol{\theta}_B = \mathbf{X}_B \boldsymbol{\beta}_B + \mathbf{Z}_B \boldsymbol{\eta}_B$$

where \mathbf{y}_B is a $k(B)$ -vector of sampling design-weighted direct total estimates for $k(B)$ -building block domains (e.g., 73 economic regions by 4 age groups by 2 genders for the example on Canadian LFS) for the study variable y which in our example is the indicator of employment for each individual in a given three digit occupation code, $\boldsymbol{\theta}_B$ is a $k(B)$ -vector of true subpopulation y -totals for building block domains, \mathbf{e}_B is the $k(B)$ -vector of sampling errors with known design-based variance-covariance matrix \mathbf{V}_B --this is only conceptual in that it may not be available due to zero sample size in some or all building block domains or otherwise very unstable for any practical utility due to very small sample size in some building blocks, \mathbf{X}_B is a $k(B) \times q$ matrix of known covariate totals at the building block domain level (covariates, for example, may consist of small area subpopulation counts for various geo-demographic groups available from population census-based projections, total count of persons with taxable employment income, and total count of employment beneficiary claims at the building block level available from administrative data), $\boldsymbol{\beta}_B$ is a $q(B)$ -vector of fixed regression parameters --the number $q(B)$ of parameters is assumed to be much less than $k(B)$ to ensure adequate degrees of freedom for estimating model parameters, $\mathbf{Z}_B \boldsymbol{\eta}_B$ is a $k(B)$ -vector of model errors where $\boldsymbol{\eta}_B$ is the $k(B)$ -vector of independent random effects with mean 0 and common variance $\sigma_{\eta(B)}^2$ and \mathbf{Z}_B is a $k(B) \times k(B)$ matrix of known coefficients that take out the variability in model errors over building blocks so that all random effects can be assumed to have a common variance. In our application, \mathbf{Z}_B will be a diagonal matrix, $\text{diag}(N_b)_{1 \leq b \leq k(B)}$, of building block subpopulation counts N_b corresponding to building blocks b .

For estimation of model parameters, building blocks could be grouped in a pre-specified hierarchical manner (based on geo-demographic categories; e.g., first group gender within age and region, next age within regions, and finally regions if necessary) to obtain

a G-level model so that the variance-covariance matrix \mathbf{V}_G is reasonably stable and can be treated as known while the number $k(G)$ of building block groups is sufficiently large for a reasonably precise estimation of model parameters $\boldsymbol{\beta}_B$ and $\sigma_{\eta(B)}^2$. We obtain the G-level model as a reduced version of the B-level model as follows.

$$\begin{aligned} \text{G-level Sampling Model: } & \mathbf{y}_G = \boldsymbol{\theta}_G + \mathbf{e}_G \\ \text{G-level Linking Model: } & \boldsymbol{\theta}_G = \mathbf{X}_G \boldsymbol{\beta}_B + \mathbf{Z}_G \boldsymbol{\eta}_B \end{aligned}$$

where $\mathbf{y}_G = \mathbf{C}_G \mathbf{y}_B$, $\boldsymbol{\theta}_G = \mathbf{C}_G \boldsymbol{\theta}_B$, $\mathbf{e}_G = \mathbf{C}_G \mathbf{e}_B$, $\mathbf{X}_G = \mathbf{C}_G \mathbf{X}_B$, $\mathbf{Z}_G = \mathbf{C}_G \mathbf{Z}_B$, and \mathbf{C}_G is a $k(G) \times k(B)$ transformation or collapsing matrix of 1's and 0's to effect the appropriate grouping. Note that the B-level model parameters of fixed and random effects carry over to G-level but the dimension of the observation vector \mathbf{y}_G is reduced from $k(B)$ to $k(G)$.

Once the model parameters $\boldsymbol{\beta}_B$ and $\sigma_{\eta(B)}^2$ (and hence the random effects, $\boldsymbol{\eta}_B$) are estimated from the G-level model, estimates of the SA parameters $\boldsymbol{\theta}_D$ for the target domains (D) can be obtained from the relation $\boldsymbol{\theta}_D = \mathbf{C}_D \boldsymbol{\theta}_B$, where \mathbf{C}_D is the appropriate $k(D) \times k(B)$ transformation matrix to go from B-level to D-level analogous to the matrix \mathbf{C}_G . Given $\boldsymbol{\beta}_B$ and $\sigma_{\eta(B)}^2$, alternative estimates of the random effects $\boldsymbol{\eta}_B$, as discussed in the next section, can also be obtained from the following D-level model (i.e., provincial level in our example).

$$\begin{aligned} \text{D-level Sampling Model: } & \mathbf{y}_D = \boldsymbol{\theta}_D + \mathbf{e}_D \\ \text{D-level Linking Model: } & \boldsymbol{\theta}_D = \mathbf{X}_D \boldsymbol{\beta}_B + \mathbf{Z}_D \boldsymbol{\eta}_B \end{aligned}$$

where the sampling error vector \mathbf{e}_D has the variance-covariance matrix \mathbf{V}_D of dimension $k(D) \times k(D)$, and other quantities are defined in a manner analogous to G-level. It may be noted that the building blocks are defined to be nested within groups as well as within SAs or target domains, and the groups are defined to be nested within SAs.

3. Estimation of Model Parameters

Under the G-level model, the variance component $\sigma_{\eta(B)}^2$ can be estimated using the restricted maximum likelihood method; see e.g., Rao (2003, p.100). Given $\sigma_{\eta(B)}^2$, the BLUP estimators of $\boldsymbol{\beta}_B$ and $\boldsymbol{\eta}_B$ using the G-level model are given by

$$\begin{aligned} \hat{\boldsymbol{\beta}}_B^{(G)} &= (\mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{X}_G)^{-1} \mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{y}_G \\ \hat{\boldsymbol{\eta}}_B^{(G)} &= \Gamma_\eta \mathbf{Z}'_G \mathbf{W}_G^{-1} (\mathbf{y}_G - \mathbf{X}_G \hat{\boldsymbol{\beta}}_B^{(G)}) \end{aligned}$$

where $\mathbf{W}_G = \mathbf{V}_G + \mathbf{Z}_G \Gamma_\eta \mathbf{Z}'_G$, and $\Gamma_\eta = \sigma_{\eta(B)}^2 \mathbf{I}_{k(B) \times k(B)}$. Instead, if we use D-level direct estimator \mathbf{y}_D and the corresponding model to estimate $\boldsymbol{\eta}_B$, we obtain $\hat{\boldsymbol{\eta}}_B^{(D)}$ as

$$\hat{\boldsymbol{\eta}}_B^{(D)} = \Gamma_\eta \mathbf{Z}'_D \mathbf{W}_D^{-1} (\mathbf{y}_D - \mathbf{X}_D \hat{\boldsymbol{\beta}}_B^{(G)})$$

where $\mathbf{W}_D = \mathbf{V}_D + \mathbf{Z}_D \Gamma_\eta \mathbf{Z}'_D$. The estimator $\hat{\boldsymbol{\eta}}_B^{(D)}$ is expected in general to exhibit less shrinkage to the zero vector than the estimator $\hat{\boldsymbol{\eta}}_B^{(G)}$ because the covariance matrix \mathbf{W}_D is expected to be more stable than \mathbf{W}_G .

The MSE matrix of estimators $(\hat{\boldsymbol{\beta}}_B^{(G)}, \hat{\boldsymbol{\eta}}_B^{(G)})$ from G-level is obtained as

$$MSE \begin{pmatrix} \hat{\boldsymbol{\beta}}_B^{(G)} - \boldsymbol{\beta}_B \\ \hat{\boldsymbol{\eta}}_B^{(G)} - \boldsymbol{\eta}_B \end{pmatrix} = \begin{pmatrix} \boldsymbol{\Sigma}_\beta^{(G)} & \boldsymbol{\Sigma}_{\beta\eta}^{(G)} \\ (\boldsymbol{\Sigma}_{\beta\eta}^{(G)})' & \boldsymbol{\Sigma}_\eta^{(G)} \end{pmatrix},$$

where $\boldsymbol{\Sigma}_\beta^{(G)} = (\mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{X}_G)^{-1}$, $\boldsymbol{\Sigma}_{\beta\eta}^{(G)} = -\boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{Z}_G \Gamma_\eta$

$\boldsymbol{\Sigma}_\eta^{(G)} = \Gamma_\eta - \Gamma_\eta \mathbf{Z}'_G \mathbf{W}_G^{-1} (\mathbf{I} - \mathbf{P}_X^{(G)}) \mathbf{Z}_G \Gamma_\eta$, and $\mathbf{P}_X^{(G)} = \mathbf{X}_G \boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_G \mathbf{W}_G^{-1}$.

Similarly, the MSE of $(\hat{\boldsymbol{\beta}}_B^{(G)}, \hat{\boldsymbol{\eta}}_B^{(D)})$ (where the random effects are estimated from D-level) is given by

$$MSE \begin{pmatrix} \hat{\boldsymbol{\beta}}_B^{(G)} - \boldsymbol{\beta}_B \\ \hat{\boldsymbol{\eta}}_B^{(D)} - \boldsymbol{\eta}_B \end{pmatrix} = \begin{pmatrix} \boldsymbol{\Sigma}_\beta^{(G)} & \boldsymbol{\Sigma}_{\beta\eta}^{(GD)} \\ (\boldsymbol{\Sigma}_{\beta\eta}^{(GD)})' & \boldsymbol{\Sigma}_\eta^{(D)} \end{pmatrix},$$

where

$$\boldsymbol{\Sigma}_{\beta\eta}^{(GD)} = -\boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{Z}_G \Gamma_\eta + \boldsymbol{\Sigma}_\beta^{(G)} (\mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{W}_{GD} - \mathbf{X}'_D) = -\boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{Z}_G \Gamma_\eta$$

because $\mathbf{W}_{GD} = \mathbf{Z}_G \Gamma_\eta \mathbf{Z}'_D + \mathbf{C}_G \mathbf{V}_B \mathbf{C}'_D = \mathbf{W}_G \mathbf{C}'_D$ under the assumption that G-level domains are nested within D-level domains and where the transformation matrix \mathbf{C}_D^* is such that $\mathbf{C}_D = \mathbf{C}_D^* \mathbf{C}_G$.

$$\begin{aligned} \boldsymbol{\Sigma}_\eta^{(D)} &= \Gamma_\eta - \Gamma_\eta \mathbf{Z}'_D \mathbf{W}_D^{-1} \mathbf{A}_1 \Gamma_\eta - \Gamma_\eta \mathbf{A}'_1 \mathbf{W}_D^{-1} \mathbf{Z}_D \Gamma_\eta + \\ &\quad \Gamma_\eta (\mathbf{Z}'_D - \mathbf{A}_2 - \mathbf{A}_3 + \mathbf{A}_4) \mathbf{W}_D^{-1} \mathbf{Z}_D \Gamma_\eta \\ &= \Gamma_\eta - \Gamma_\eta \mathbf{Z}'_D \mathbf{W}_D^{-1} \mathbf{A}_1 \Gamma_\eta - \Gamma_\eta \mathbf{A}'_1 \mathbf{W}_D^{-1} \mathbf{Z}_D \Gamma_\eta + \Gamma_\eta (\mathbf{Z}'_D - \mathbf{A}_2) \mathbf{W}_D^{-1} \mathbf{Z}_D \Gamma_\eta \end{aligned}$$

where $\mathbf{A}_1 = \mathbf{Z}_D - \mathbf{X}_D \boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{Z}_G = \mathbf{C}_D^* (\mathbf{I} - \mathbf{P}_X^{(G)}) \mathbf{Z}_G$,

$\mathbf{A}_2 = \mathbf{Z}'_D \mathbf{W}_D^{-1} \mathbf{X}_D \boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_G \mathbf{W}_G^{-1} \mathbf{W}_{GD} = \mathbf{Z}'_D \mathbf{W}_D^{-1} \mathbf{X}_D \boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_D$

$\mathbf{A}_3 = \mathbf{Z}'_D \mathbf{W}_D^{-1} \mathbf{W}'_{GD} \mathbf{W}_G^{-1} \mathbf{X}_G \boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_D = \mathbf{A}_2$, $\mathbf{A}_4 = \mathbf{Z}'_D \mathbf{W}_D^{-1} \mathbf{X}_D \boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}'_D = \mathbf{A}_2$.

Therefore,

$$\begin{aligned} \boldsymbol{\Sigma}_\eta^{(D)} = & \Gamma_\eta - \Gamma_\eta \mathbf{Z}'_D \mathbf{W}_D^{-1} \mathbf{C}_D^* (\mathbf{I} - \mathbf{P}_X^{(G)}) \mathbf{Z}_G \Gamma_\eta - \Gamma_\eta \mathbf{Z}'_G (\mathbf{I} - \mathbf{P}_X^{(G)'}) \mathbf{C}_D^{*'} \mathbf{W}_D^{-1} \mathbf{Z}_D \Gamma_\eta + \\ & \Gamma_\eta \mathbf{Z}'_D (\mathbf{I} - \mathbf{P}_X^{(D)'}) \mathbf{W}_D^{-1} \mathbf{Z}_D \Gamma_\eta. \end{aligned}$$

It is easily seen that $\boldsymbol{\Sigma}_\eta^{(D)}$ reduces to $\boldsymbol{\Sigma}_\eta^{(G)}$ if $\mathbf{C}_D = \mathbf{C}_G$ as in the case of usual BLUP theory.

4. Estimation of SA Parameters and their MSE

There are two estimators (SAE-G and SAE-D) that could be obtained for target SA parameters by using G-level and D-level estimators of random effects respectively. Using the G-level, we get results similar to the standard BLUP theory as given below.

$$\begin{aligned} \hat{\boldsymbol{\theta}}_D^{(G)} &= \mathbf{X}_D \hat{\boldsymbol{\beta}}_B^{(G)} + \mathbf{Z}_D \hat{\boldsymbol{\eta}}_B^{(G)} \\ \text{MSE}(\hat{\boldsymbol{\theta}}_D^{(G)} - \boldsymbol{\theta}_D) &= (\mathbf{X}_D \quad \mathbf{Z}_D) \begin{pmatrix} \boldsymbol{\Sigma}_\beta^{(G)} & \boldsymbol{\Sigma}_{\beta\eta}^{(G)} \\ (\boldsymbol{\Sigma}_{\beta\eta}^{(G)})' & \boldsymbol{\Sigma}_\eta^{(D)} \end{pmatrix} (\mathbf{X}_D \quad \mathbf{Z}_D)' \end{aligned}$$

which simplifies to $\text{MSE}(\hat{\boldsymbol{\theta}}_D^{(G)}) = g_1^{(G)} + g_2^{(G)}$ where

$$g_1^{(G)} = \mathbf{C}_D^* (\mathbf{V}_G - \mathbf{V}_G \mathbf{W}_G^{-1} \mathbf{V}_G) \mathbf{C}_D^{*'}, \quad g_2^{(G)} = \mathbf{C}_D^* (\mathbf{V}_G \mathbf{W}_G^{-1} \mathbf{X}_G \boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}_G' \mathbf{W}_G^{-1} \mathbf{V}_G) \mathbf{C}_D^{*'}.$$

Substituting the REML estimator of $\sigma_{\eta(B)}^2$ in the above expression, we get a naïve estimator which can be adjusted along the lines of Prasad and Rao (1990) to get the second order adjusted MSE estimator; see also Rao (2003, p. 104). Using D-level for estimating random effects, we get alternate estimators SAE-D as follows.

$$\begin{aligned} \hat{\boldsymbol{\theta}}_D^{(D)} &= \mathbf{X}_D \hat{\boldsymbol{\beta}}_B^{(G)} + \mathbf{Z}_D \hat{\boldsymbol{\eta}}_B^{(D)} \\ \text{MSE}(\hat{\boldsymbol{\theta}}_D^{(D)}) &= (\mathbf{X}_D \quad \mathbf{Z}_D) \begin{pmatrix} \boldsymbol{\Sigma}_\beta^{(G)} & \boldsymbol{\Sigma}_{\beta\eta}^{(GD)} \\ (\boldsymbol{\Sigma}_{\beta\eta}^{(GD)})' & \boldsymbol{\Sigma}_\eta^{(D)} \end{pmatrix} (\mathbf{X}_D \quad \mathbf{Z}_D)' \end{aligned}$$

which can be simplified as in the previous case to obtain $\text{MSE}(\hat{\boldsymbol{\theta}}_D^{(D)}) = g_1^{(D)} + g_2^{(D)}$ where denoting $\mathbf{Z}_G \Gamma_\eta \mathbf{Z}_G'$ by Γ_G , $\mathbf{Z}_D \Gamma_\eta \mathbf{Z}_D'$ by Γ_D , and $\mathbf{Z}_D \mathbf{V}_G \mathbf{Z}_D'$ by \mathbf{V}_D so that $\mathbf{W}_D = \Gamma_D + \mathbf{V}_D$, we have after some algebra

$$g_1^{(D)} = \mathbf{V}_D - \mathbf{V}_D \mathbf{W}_D^{-1} \mathbf{V}_D = g_1^{(G)} + \mathbf{C}_D^* \mathbf{V}_G (\mathbf{W}_G^{-1} - \mathbf{C}_D^{*'} \mathbf{W}_D^{-1} \mathbf{C}_D^*) \mathbf{V}_G \mathbf{C}_D^{*'}, \text{ and}$$

$$g_2^{(D)} = g_2^{(G)} - \mathbf{C}_D^* \Gamma_G (\mathbf{W}_G^{-1} - \mathbf{C}_D^{*'} \mathbf{W}_D^{-1} \mathbf{C}_D^*) \mathbf{X}_G \boldsymbol{\Sigma}_\beta^{(G)} \mathbf{X}_G' (\mathbf{W}_G^{-1} - \mathbf{C}_D^{*'} \mathbf{W}_D^{-1} \mathbf{C}_D^*) \Gamma_G \mathbf{C}_D^{*'}.$$

It is easily seen that the leading term $g_1^{(D)}$ in the MSE of the SAE-D estimator is larger than the term $g_1^{(G)}$ of MSE of SAE-G implying that SAE-D is not optimal for the G-level model as expected. The lower order term $g_2^{(D)}$, however, is smaller than the

corresponding term $g_2^{(G)}$. A second order adjustment to the above naive MSE estimator of SAE-D can also be developed.

5. Benchmarking of SAEs

For benchmarking of SAE-Ds (which use D-level estimators of random effects) to the national or a large subnational direct estimator, we add an extra covariate $\mathbf{x}_{r(D)}^+$ (which depends on covariance and transformation matrices) in the G-level model for the r th benchmark; see Singh, 2006. It is defined as

$$\mathbf{x}_{r(D)}^+ = \mathbf{W}_G \mathbf{C}_D^{*'} (\mathbf{C}_D \mathbf{W}_G \mathbf{C}_D')^{-1} \mathbf{V}_D \boldsymbol{\delta}_{r(D)}$$

where $\boldsymbol{\delta}_{r(D)}$ is simply a vector of 1s and 0s indicating membership of D-level areas contributing to the benchmark total, and the matrix \mathbf{C}_D^* , as before, is the transformation matrix from \mathbf{C}_G to \mathbf{C}_D under the assumption that the G-level domains are nested within D-level domains. Introduction of such new covariates in effect extends the $\boldsymbol{\beta}_B$ vector to $\boldsymbol{\beta}_B^+$ and the covariate matrix \mathbf{X}_G to \mathbf{X}_G^+ such that the additional estimating function is defined as

$$\mathbf{x}_{r(D)}^{+'} \mathbf{W}_G^{-1} (\mathbf{y}_G - \mathbf{X}_G^+ \hat{\boldsymbol{\beta}}_B^{+(G)}) = \mathbf{0}$$

which implies, in particular, that the benchmark condition is automatically satisfied; i.e.,

$$\boldsymbol{\delta}_{r(D)}' (\mathbf{y}_D - \mathbf{X}_D^+ \hat{\boldsymbol{\beta}}_B^{+(G)} - \mathbf{Z}_D \hat{\boldsymbol{\eta}}_B^{(D)}) = \mathbf{0}$$

Similarly, benchmark constraints on SAE-Gs (which use G-level estimators of random effects) can be realized by extending the \mathbf{X}_G matrix using new covariates $\mathbf{x}_{r(D)}^+$ defined by $\mathbf{V}_G \boldsymbol{\delta}_{r(G)}$.

6. Empirical Results: Application to the Canadian LFS

For the 2003 monthly Canadian LFS survey example, Table 1 shows the average sample size per province over the 12 month period, direct average estimate of employed in A39, standard error (SE) and coefficient of variation (CV). Except for ON and QC, estimates for all other provinces have CV near 15 % or higher which suggests a need for SAE and to find out how much efficiency gain can be realized with a simple linear mixed model. We consider building block domains defined by economic region by age by gender as building blocks—with 73 economic regions, 4 age categories (15-24, 25-34, 35-54, 55+), and 2 gender categories, the total number of building blocks $k(B)$ is 584. It is assumed that although not all building blocks are observed in the sample, the sampling design is such that there is a positive probability of nonzero sample size for each building block domain. Let N_{rj} denote the subpopulation size for rj domain in the b th building block where b represents the cross-classification rj of r th economic region by j th age-gender category. For an illustration of bBLUP, we consider a simple subgroup common mean model at the B-level for the mean μ_{rj} of $N_{rj}^{-1} y_{rj}$ which is given by

$$\begin{aligned}
 y_{rj} &= N_{rj}\mu_{rj} + \varepsilon_{rj} \\
 &= N_{rj}\beta_j + N_{rj}\eta_{rj} + \varepsilon_{rj},
 \end{aligned}$$

where the mean μ_{rj} (conditional on the domain rj) has common unconditional mean β_j over all economic regions r with model error η_{rj} —the random effect., and ε_{rj} being the sampling error in aggregate level estimator y_{rj} about $N_{rj}\mu_{rj}$. The random effects η_{rj} are assumed to be exchangeable with mean 0 and variance $\sigma_{\eta(B)}^2$. The G-level model at the economic region level is obtained from the above B-level model as

$$y_r = \sum_{j=1}^J N_{rj}\beta_j + \sum_{j=1}^J N_{rj}\eta_{rj} + \varepsilon_r,$$

where $\varepsilon_r (= \sum_{j=1}^{k(B)} \varepsilon_{rj})$ is the sampling error in the economic region level aggregate estimator, and J equals 8—the number of age-gender categories which in the earlier notation is simply $q(B)$. There are 73 groups or economic regions for the G-level model which were further regrouped to 69 to avoid groups with zero or very small sample size so that $k(G) = 69$. The D-level model at the target domain level can be obtained from the above G-level model by aggregating over regions within each province; total number $k(D)$ being 10. The design-based variance-covariance matrix V_G of the sampling error vector $\{\varepsilon_r\}_{1 \leq r \leq k(G)}$ is block-diagonal with provinces as strata but within provinces it is nondiagonal because primary sampling units (PSUs) in the LFS design cut across economic regions. In this illustration, the above model was extended to include only one benchmark consisting of the direct estimator at the national level.

Table 2 presents CVs (defined as relative root MSE) and relative difference of bBLUPs SAE-D* (with estimated random effects at D-level, and * indicates benchmarked) and compares with bBLUPs SAE-G (unbenchmarkd with estimated random effects at G-level) while Table 3 presents results for the complementary case where bBLUPs SAE-G are benchmarked and not bBLUPs SAE-D. Note that the model level, (at which random effects are estimated), determines which type of SAEs get benchmarked automatically through estimating functions for fixed effects. The term modified (mod) CV for SAE is used to signify that the denominator is the direct estimator and not the SAE for ease in comparison between $CV(\text{direct})$, $CV(\text{SAE-D}^*)$ and $CV(\text{SAE-G}^*)$, for example, all having a common denominator. The relative difference measures departure of SAEs from the direct estimator relative to direct. It is observed that both SAE-D1* and SAE-G* perform pretty much at par with respect to mod CV except for the province NL where SAE-D* has a mod CV of 19% compared to SAE-G* with a mod CV of 16%. However, in terms of relative difference, SAE-D* has a clear superior performance over SAE-G* for all provinces and especially for MA which was brought down in magnitude from 25% for SAE-G* to 18% for SAE-D*. It is desirable in practice to have SAEs that have good efficiency gains over direct estimators but exhibit only modest departures from them in general because of possible model misspecification. The SAE-D* estimators seem to provide a good balance between CV efficiency and reducing relative difference as anticipated because direct estimators tend to get more weight in SAE-D*. Interestingly,

the unbenchmarked SAE-D and SAE-G perform reasonably well in terms of mod CV but poorly with respect to relative difference. It may be noted that all MSE estimators in the above example were second order adjusted and REML was used for estimating the variance component.

7. Concluding Remarks and Discussion of Related Applications

In this paper, the method of bBLUP estimation was introduced as an alternative to BLUP in order to have a principled approach of having a very low level building-block model for deriving all other aggregate level SAE models. The model was fitted by grouping of building blocks which should be based on a priori considerations (or past experience) but should not be data dependent. The new estimators can also be benchmarked in the interest of robustification by adding suitable covariates at the group model level. This research was conducted under an R&D initiative at Statistics Canada with the goal of developing a client-oriented SAE product (Singh, 2006). The goal was to develop a system which retains the simplicity of linear mixed models and use of aggregate levels for taking the sampling design into account as in the Fay-Herriot model but has new features of benchmarking, smoothing of sampling covariance of direct estimators, and grouping of building blocks to avoid an ad hoc user-specified choice of the area level in real applications. Although the building block idea does not completely resolve the ad hocry in the choice of aggregation level, it can assist a great deal in satisfying the exchangeability assumption.

For the above SAE project, we also considered the problem of ensuring SAEs to satisfy range restrictions such as 0 to 1 in the case of estimating the proportion of employed in a three digit occupation code. The linear mixed model approach doesn't guarantee that range restrictions on parameter estimates (even nonnegativity for positive parameters) will be satisfied. In this context, we did not wish to use nonlinear models for the mean function in order that the simplicity of linear mixed models was not lost. We therefore proposed a new estimation method termed quasi-bBLUP (to be published in a separate paper) which, analogous to quasi-likelihood estimation for fixed effects models, employs working expressions of covariance matrices \mathbf{V}_D and $\mathbf{\Gamma}_\eta$ to ensure range restrictions. The resulting estimates remain unbiased but not optimal. The reason for this is that the BLUP estimator can also be obtained as a linear projection (or linear regression predictor) of the random parameter on the data space which requires only the specification of the covariance structure. If the covariance structure is working (e.g., assuming \mathbf{V}_D diagonal when it is not, because off-diagonal estimates are too unstable, or taking a pre-specified value of $\sigma_{\eta(B)}^2$ because the estimated value may be too close to zero in cases where sampling error is high), the resulting BLUP-type estimating functions remain unbiased unconditionally on the random effects. It has been our experience that the synthetic estimator (i.e., the linear mean predictor) typically satisfies range restrictions which are of course satisfied by direct estimators. Therefore, with a diagonal version of \mathbf{V}_D , the resulting composite estimates or SAE necessarily satisfy range restrictions. However due to their

nonoptimality, the MSE expression for BLUPs needs to be adjusted for using a working version of \mathbf{V}_D .

Similarly the idea of using working values or biased estimates of parameters in a general $\mathbf{\Gamma}_\eta$ (it is simply $\sigma_{\eta(B)}^2 \mathbf{I}$ in our case) can be useful in dealing with SAE for combined cross-sectional and time series data which involves correlation in random effects over time that are known to be difficult to estimate in a stable manner. In many such situations, the generalization of bBLUPs to quasi bBLUP might be quite useful which unlike Fay-Herriot, relies on suboptimal estimation by using the ideas of building blocks and a working covariance structure. However, in some cases especially when dealing with rare outcomes, it may be necessary to use nonlinear models because the synthetic linear predictor may not satisfy range restrictions. In these cases, Singh and Verret (2006) proposed nonlinear aggregate level models with additive random components in the sense that random effects are outside the nonlinear mean function which contains only fixed effects. By having additive random effects, it is possible to simplify considerably estimation of model parameters somewhat similar to the case of linear mixed models.

Although in this paper we considered the case when the target domains were at a level higher than the group level, the method of grouping can also be used to deal with situations where target SAs are at a low level resulting in SAs with zero or very small sample size. If the sampling design is such that there are no samples planned for certain SAs, then synthetic estimation is the only option for those areas. However, if there is no or a very small sample in a SA by chance (e.g., when SAs are counties in US), then bBLUPS at the G-level can be used to define suitable nonsynthetic estimates; here building blocks are at a level lower than the target SA level (e.g., subcounties defined by demographics if SAs are counties). For this purpose, first SAEs for all groups (i.e., using G-level direct estimators and random effect estimators at the G-level) are obtained while allowing for benchmarking to higher level direct estimators. To define SAE-D where D-level is now lower than the G-level, it may not be practical to use direct estimators which may not be available or may be very unstable. An alternative would be to first obtain synthetic estimators at D-level, and then adjust them so that all D-level estimates for domains nested within a group add up or are benchmarked to the SAE for that group. This is in fact satisfied by SAE-G or $\hat{\theta}_D^{(G)}$ defined in Section 4. However, these estimators are not designed to respect range restrictions such as that of nonnegative estimation for nonnegative study variables or that of being in the interval (0,1) for estimating proportions. A way out is to use range restricted calibration methods (commonly used in adjusting sampling weights) to adjust synthetic estimates within a group to satisfy the group level SAE while meeting range restrictions.

Acknowledgments

This research work was an outcome of a new initiative on developing a small area product at Statistics Canada. The authors would like to thank Jon Rao and Wayne Fuller

for several useful discussions during the early stages of this project. They would also like to thank Francois Verret for his contributions during the initial part of the project.

References

Prasad, N.G.N. and Rao, J.N.K. (1990), The estimation of mean squared error of small area estimators, *Jour. Amer. Statist. Assoc.*, 85, 163-171.

Rao, J.N.K. (2003). *Small Area Estimation*, Wiley-Interscience, N.J.

Singh, A.C. (2006), Some problems and proposed solutions in developing a small area estimation product for clients. *ASA Proc. Surv. Res. Meth. Sec.*, 3673-3683.

Singh, A.C., and Verret. F. (2006). Mixed linear nonlinear models for small area estimation with application to the Canadian community health Survey. *Proc. Statistics Canada Symposium on Statistical Issues in Measuring Population Health, Ottawa, ON*

Table 1: Monthly Total Employed (A39)—Annual Average for 2003 LFS

| Province | Pop. Size | Sample size | Direct | SE | CV |
|----------|------------|-------------|--------|-------|-------|
| NL | 429,298 | 3,577 | 670 | 177 | 0.264 |
| PEI | 109,886 | 2,769 | 233 | 55 | 0.235 |
| NS | 758,549 | 5,858 | 1,532 | 292 | 0.19 |
| NB | 607,565 | 5,624 | 1,275 | 218 | 0.171 |
| QC | 6,059,655 | 18,062 | 25,273 | 2,204 | 0.087 |
| ON | 9,766,566 | 30,373 | 42,447 | 3,178 | 0.075 |
| MA | 876,396 | 7,117 | 3,023 | 432 | 0.143 |
| SK | 744,431 | 7,241 | 1,963 | 339 | 0.173 |
| AB | 2,467,412 | 10,317 | 7,643 | 1,098 | 0.144 |
| BC | 3,346,181 | 9,110 | 8,676 | 1,228 | 0.142 |
| Canada | 25,165,939 | 100,048 | 92,734 | 4,260 | 0.046 |

Note: NL: Newfoundland and Labrador, PEI: Prince Edward Island, NS: Nova Scotia, NB: New Brunswick, QC: Quebec, ON: Ontario, MA: Manitoba, SK: Saskatchewan, AB: Alberta, BC: British Columbia.

Table 2: Direct Estimators, SAE-G, and SAE-D* for Monthly Total Employed (A39)

(* indicates benchmarking; SAE-D and SAE-G correspond to estimated random effects at D- and G-levels respectively)

| PROV | Direct Estimate | | SAE-G | | | SAE-D* | | |
|--------|-----------------|-------|----------|---------|------------|----------|---------|------------|
| | Estimate | CV | Estimate | Mod. CV | Rel. Diff. | Estimate | Mod. CV | Rel. Diff. |
| NL | 670 | 0.264 | 621 | 0.182 | -0.07 | 648 | 0.186 | -0.03 |
| PEI | 233 | 0.235 | 207 | 0.199 | -0.11 | 207 | 0.199 | -0.11 |
| NS | 1,532 | 0.19 | 1,323 | 0.136 | -0.14 | 1,384 | 0.137 | -0.1 |
| NB | 1,275 | 0.171 | 1,224 | 0.131 | -0.04 | 1,236 | 0.133 | -0.03 |
| QC | 25,273 | 0.087 | 24,514 | 0.067 | -0.03 | 25,162 | 0.067 | 0 |
| ON | 42,447 | 0.075 | 41,748 | 0.065 | -0.02 | 43,098 | 0.065 | 0.02 |
| MA | 3,023 | 0.143 | 2,368 | 0.099 | -0.22 | 2,468 | 0.102 | -0.18 |
| SK | 1,963 | 0.173 | 1,819 | 0.124 | -0.07 | 1,848 | 0.124 | -0.06 |
| AB | 7,643 | 0.144 | 7,707 | 0.098 | 0.01 | 7,837 | 0.099 | 0.03 |
| BC | 8,676 | 0.142 | 8,807 | 0.11 | 0.02 | 8,846 | 0.111 | 0.02 |
| Canada | 92,734 | 0.046 | 90,337 | 0.046 | -0.03 | 92,734 | 0.046 | 0 |

Table 3: Direct Estimators, SAE-G* and SAE-D for Monthly Total Employed (A39)

(* indicates benchmarking; SAE-D and SAE-G correspond to estimated random effects at D- and G-levels respectively)

| PROV | Direct Estimate | | SAE-G* | | | SAE-D | | |
|--------|-----------------|-------|----------|---------|------------|----------|---------|------------|
| | Estimate | CV | Estimate | Mod. CV | Rel. Diff. | Estimate | Mod. CV | Rel. Diff. |
| NL | 670 | 0.264 | 600 | 0.159 | -0.1 | 624 | 0.162 | -0.07 |
| PEI | 233 | 0.235 | 195 | 0.191 | -0.16 | 195 | 0.191 | -0.16 |
| NS | 1,532 | 0.19 | 1,299 | 0.12 | -0.15 | 1,346 | 0.12 | -0.12 |
| NB | 1,275 | 0.171 | 1,184 | 0.12 | -0.07 | 1,191 | 0.121 | -0.07 |
| QC | 25,273 | 0.087 | 25,218 | 0.062 | 0 | 25,582 | 0.063 | 0.01 |
| ON | 42,447 | 0.075 | 43,584 | 0.064 | 0.03 | 44,315 | 0.064 | 0.04 |
| MA | 3,023 | 0.143 | 2,273 | 0.091 | -0.25 | 2,330 | 0.092 | -0.23 |
| SK | 1,963 | 0.173 | 1,783 | 0.113 | -0.09 | 1,803 | 0.113 | -0.08 |
| AB | 7,643 | 0.144 | 7,730 | 0.087 | 0.01 | 7,815 | 0.088 | 0.02 |
| BC | 8,676 | 0.142 | 8,868 | 0.099 | 0.02 | 8,878 | 0.1 | 0.02 |
| Canada | 92,734 | 0.046 | 92,734 | 0.046 | 0 | 94,079 | 0.046 | 0.01 |