

Article

Small area estimation under transformation to linearity

by Hukum Chandra and Ray Chambers



June 2011

Small area estimation under transformation to linearity

Hukum Chandra and Ray Chambers¹

Abstract

Small area estimation based on linear mixed models can be inefficient when the underlying relationships are non-linear. In this paper we introduce SAE techniques for variables that can be modelled linearly following a non-linear transformation. In particular, we extend the model-based direct estimator of Chandra and Chambers (2005, 2009) to data that are consistent with a linear mixed model in the logarithmic scale, using model calibration to define appropriate weights for use in this estimator. Our results show that the resulting transformation-based estimator is both efficient and robust with respect to the distribution of the random effects in the model. An application to business survey data demonstrates the satisfactory performance of the method.

Key Words: Sample survey; Survey estimation; Business surveys; Model calibration; Skewed data; Model-based direct estimation; Empirical best linear unbiased prediction.

1. Introduction

Commonly used methods for small area estimation (SAE) assume that a linear mixed model can be used to characterize the regression relationship between the survey variable Y and an auxiliary variable X in the small areas of interest. In particular, empirical best linear unbiased prediction (EBLUP), see Rao (2003, chapters 6 - 8) is typically based on a linear mixed model assumption. However, when the data are skewed, as is often the case in business surveys, the relationship between Y and X may not be linear in the original (raw) scale, but can be linear in a transformed scale, *e.g.*, the logarithmic (log) scale. In such cases we would expect estimation based on a linear mixed model for Y to be inefficient compared with one based on a similar model for a transformed version of Y . See Hidioglou and Smith (2005). The use of transformations in inference has a long history, see for example Carroll and Ruppert (1988, chapter 4). Recently, Chen and Chen (1996) and Karlberg (2000a) have investigated the use of a 'transform to linearity' approach for regression estimation of survey variables that behave non-linearly. However, to the best of our knowledge there has been no application of this idea in SAE, even though economic theory (and casual observation) suggests that regression relationships in business survey data are typically multiplicative, and hence linear in the log scale.

In this paper we extend the model-based direct (MBD) estimation ideas described in Chandra and Chambers (2005, 2009) to the situation where the linear mixed model underpinning SAE holds on the log scale, using weights derived via model calibration (Wu and Sitter 2001). In doing so, we note that our approach easily generalises to

other monotone (*i.e.*, invertible) transformations. In contrast, extension of the EBLUP approach to where the data follow a linear mixed model under transformation is complicated. We also relax the usual normality assumption for the area effects in order to examine robustness with respect to this assumption.

In the following section we summarise the MBD approach to SAE under a linear mixed model. In section 3 we describe an alternative to the linear mixed model for skewed data which reduces to the linear mixed model under log transformation, and in section 4 we use a model-based perspective to motivate model calibrated estimation of population quantities where the underlying variable is linear after suitable transformation. In section 5 we bring these two ideas together, introducing the concept of a fitted value model derived from a linear mixed model in the transformed scale. We then use this fitted value model to specify survey weights for use in an MBD estimator in SAE. In section 6 we present empirical results from a number of simulation studies that contrast the proposed transformation-based MBD estimator with both the EBLUP and the 'usual' MBD estimator defined by fitting a linear mixed model to the data as well as with an indirect empirical predictor based on the same transformed scale linear mixed model. Section 7 concludes the paper with a discussion of outstanding issues.

Note that the approach taken in this article is model-based. Consequently all moments are evaluated with respect to a model for the population data. Also, all sample data are assumed to have been obtained via a non-informative sampling method, *e.g.*, probability sampling with inclusion probabilities defined by known model covariates.

1. Hukum Chandra, Indian Agricultural Statistics Research Institute, Library Avenue, PUSA Campus, New Delhi-110012, India. E-mail: hchandra@iasri.res.in; Ray Chambers, Centre for Statistical and Survey Methodology, University of Wollongong, Wollongong, NSW, 2522, Australia. E-mail: ray@uow.edu.au.

2. Model-based direct estimation for small areas

To start, we fix our notation. Let U denote a population of size N and let \mathbf{y}_U denote the N -vector of population values of a characteristic Y of interest. Suppose that our primary aim is estimation of the total $t_{Uy} = \sum_U y_j$ of these population values (or their mean $m_{Uy} = N^{-1} \sum_U y_j$). Let \mathbf{X} denote a p -vector of auxiliary variables that are related, in some sense, to Y and let \mathbf{x}_U denote the corresponding $N \times p$ matrix of population values these variables. We assume that the individual sample values of \mathbf{X} are known. The non-sample values of \mathbf{X} may not be individually known, but are assumed known at some aggregate level. At a minimum, we know the vector of population totals \mathbf{t}_{Ux} of the columns of \mathbf{X} .

Suppose that it is reasonable to assume that the regression of Y on \mathbf{X} in the population is linear, *i.e.*,

$$E(\mathbf{y}_U | \mathbf{x}_U) = \mathbf{x}_U \boldsymbol{\beta} \text{ and } \text{Var}(\mathbf{y}_U | \mathbf{x}_U) = \mathbf{v}_U \quad (1)$$

where \mathbf{v}_U is known up to a multiplicative constant. Given a sample s of size n from this population, we can partition

$$\mathbf{x}_U = \begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_r \end{bmatrix}$$

and

$$\mathbf{v}_U = \begin{bmatrix} \mathbf{v}_{ss} & \mathbf{v}_{sr} \\ \mathbf{v}_{rs} & \mathbf{v}_{rr} \end{bmatrix}$$

into their sample and non-sample components. Here $r = U - s$ denotes the population units that are not in sample. The vector of weights that defines the Best Linear Unbiased Predictor (BLUP) of t_{Uy} is then (Royall 1976; Valliant, Dorfman and Royall 2000, section 2.4)

$$\mathbf{w}_s^{\text{BLUP}} = (w_j^{\text{BLUP}}; j \in s) \\ = \mathbf{1}_s + \mathbf{H}'_s (\mathbf{t}_{Ux} - \mathbf{t}_{sx}) + (\mathbf{I}_s - \mathbf{H}'_s \mathbf{x}'_s) \mathbf{v}_{ss}^{-1} \mathbf{v}_{sr} \mathbf{1}_r \quad (2)$$

where $\mathbf{H}_s = (\mathbf{x}'_s \mathbf{v}_{ss}^{-1} \mathbf{x}_s)^{-1} \mathbf{x}'_s \mathbf{v}_{ss}^{-1}$, \mathbf{I}_s is the identity matrix of order n , \mathbf{t}_{sx} is the vector of sample totals of \mathbf{X} and $\mathbf{1}_s$ ($\mathbf{1}_r$) denotes a vector of ones of size n ($N - n$).

We now assume that the target population U of size N can be partitioned into D non-overlapping small areas or domains, each of size N_i , $i = 1, \dots, D$, such that $N = \sum_{i=1}^D N_i$. Given a sample s of size n units is drawn from this population, we shall assume that a sub-sample s_i of size n_i units is drawn from area i , with $n = \sum_{i=1}^D n_i$. Note that we assume that all small areas are sampled and that there is at least one sample unit in each small area of interest.

As noted in section 1, linear mixed models are often used in SAE. Such models can be written in the form

$$\mathbf{y}_U = \mathbf{x}_U \boldsymbol{\beta} + \mathbf{g}_U \mathbf{u} + \mathbf{e}_U \quad (3)$$

where \mathbf{u} is a random vector of so-called area effects, \mathbf{e}_U is a population N -vector of random individual effects and \mathbf{g}_U is a known matrix. In general, area effects are vector-valued, so $\mathbf{u}' = (\mathbf{u}'_1 \mathbf{u}'_2 \dots \mathbf{u}'_D)$ and $\mathbf{g}_U = \text{diag}\{\mathbf{g}_i; i = 1, \dots, D\}$, where \mathbf{g}_i is of dimension $N_i \times q$. The area specific effects $\{\mathbf{u}_i; i = 1, \dots, D\}$ are assumed to be independent and identically distributed realisations of a random vector of dimension q with zero mean and covariance matrix Σ_u . Similarly, the scalar individual effects making up \mathbf{e}_U are assumed to be independent and identically distributed realisations of a random variable with zero mean and variance σ_e^2 , with area and individual effects mutually independent. The parameters $\theta = (\Sigma_u, \sigma_e^2)$ are typically referred to as the variance components of (3).

Given the values of the variance components, it is straightforward to see that (3) is just a special case of the general linear model (1) that underpins the BLUP weights (2). In particular, under (3)

$$\mathbf{v}_{ss} = \text{diag}\{\mathbf{v}_{iss}; i = 1, \dots, D\} \\ = \text{diag}\{\mathbf{g}_{is} \Sigma_u \mathbf{g}'_{is} + \sigma_e^2 \mathbf{I}_{is}; i = 1, \dots, D\} \quad (4)$$

and

$$\mathbf{v}_{sr} = \text{diag}\{\mathbf{v}_{isr}; i = 1, \dots, D\} \\ = \text{diag}\{\mathbf{g}_{is} \Sigma_u \mathbf{g}'_{ir}; i = 1, \dots, D\}. \quad (5)$$

Here \mathbf{g}_{is} and \mathbf{g}_{ir} denote the restriction of \mathbf{g}_i to sampled and non-sampled units in area i respectively. Given estimated values $\hat{\theta} = (\hat{\Sigma}_u, \hat{\sigma}_e^2)$ of the variance components we can substitute these in (4) and (5) to obtain estimates $\hat{\mathbf{v}}_{ss}$ and $\hat{\mathbf{v}}_{sr}$ of \mathbf{v}_{ss} and \mathbf{v}_{sr} respectively, and therefore compute 'empirical' BLUP weights, or EBLUP weights for the population total of Y as

$$\mathbf{w}_s^{\text{EBLUP}} = (w_{ij}^{\text{EBLUP}}; j \in s_i; i = 1, \dots, D) \\ = \mathbf{1}_s + \hat{\mathbf{H}}'_s (\mathbf{t}_{Ux} - \mathbf{t}_{sx}) \\ + (\mathbf{I}_s - \hat{\mathbf{H}}'_s \mathbf{x}'_s) \hat{\mathbf{v}}_{ss}^{-1} \hat{\mathbf{v}}_{sr} \mathbf{1}_r \quad (6)$$

where $\hat{\mathbf{H}}_s = (\mathbf{x}'_s \hat{\mathbf{v}}_{ss}^{-1} \mathbf{x}_s)^{-1} \mathbf{x}'_s \hat{\mathbf{v}}_{ss}^{-1}$. Note that we now use a double index of ij to differentiate between population units in different areas.

The MBD estimator for the mean m_{iy} of Y in area i (Chandra and Chambers 2005, 2009) based on the EBLUP weights for the total (6) is simply the corresponding weighted average of the sample values of Y in area i ,

$$\hat{m}_{iy}^{HJ\text{-LinMBD}} = \left\{ \sum_{j \in s_i} w_{ij}^{EBLUP} \right\}^{-1} \sum_{j \in s_i} w_{ij}^{EBLUP} y_{ij}. \quad (7)$$

Note that (7) is *not* the EBLUP for m_{iy} under (3). This is (see Rao 2003, section 6.2.3)

$$\begin{aligned} \hat{m}_{iy}^{HT\text{-LinEBLUP}} &= \hat{E}\{m_{iy} | \mathbf{y}_{is}, \mathbf{x}_{is}, \mathbf{x}_{ir}\} \\ &= N_i^{-1} \left[\sum_{j \in s_i} y_j + \mathbf{1}'_{ir} \left\{ \mathbf{x}_{ir} \hat{\beta} + \hat{\mathbf{v}}_{irs}^{-1} (\mathbf{y}_{is} - \mathbf{x}_{is} \hat{\beta}) \right\} \right] \\ &= N_i^{-1} \left[n_i \bar{y}_{is} + (N_i - n_i) \right. \\ &\quad \left. \left\{ \bar{\mathbf{x}}'_{ir} \hat{\beta} + \bar{\mathbf{g}}'_{ir} \hat{\Sigma}_u \mathbf{g}'_{is} (\mathbf{g}_{is} \hat{\Sigma}_u \mathbf{g}'_{is} + \hat{\sigma}_e^2 \mathbf{I}_{is})^{-1} (\mathbf{y}_{is} - \mathbf{x}_{is} \hat{\beta}) \right\} \right]. \quad (8) \end{aligned}$$

Here \hat{E} denotes the expectation operator under (3) with unknown parameters replaced by estimates, \mathbf{x}_{is} and \mathbf{x}_{ir} are the matrices of sample and non-sample values of \mathbf{X} in area i , \mathbf{y}_{is} is the vector of sample values of Y in the same area, $\hat{\beta}$ is the ‘empirical’ BLUE of β , $\hat{\mathbf{v}}_{irs}$ is the transpose of the estimated value of \mathbf{v}_{irs} with $\hat{\mathbf{v}}_{iss}$ the corresponding estimate of \mathbf{v}_{iss} , see (4) and (5), and $\mathbf{1}_{ir}$ is a vector of ones of length $N_i - n_i$. Note that the last expression on the right hand side of (8) follows directly by substitution of (4) and (5), with $\bar{\mathbf{x}}_{ir}$ and $\bar{\mathbf{g}}_{ir}$ denoting the column vectors of order p and q defined by averaging the columns of \mathbf{x}_{ir} and \mathbf{g}_{ir} respectively. Like the EBLUP (8), the estimator (7) is a weighted function of all the sample values. Note that under random intercept specification of (3), (8) reduces to the expression (7.2.39) in Rao (2003, section 7.2).

Mean squared error (MSE) estimation for (8) is usually carried out using the theory described in Prasad and Rao (1990). Although this MSE estimator is somewhat complicated, it works well under (3). However, when (3) fails it can be misleading. It is also inadequate as an estimator of the repeated sampling MSE of (8), as has been pointed out by Longford (2007). In contrast, MSE estimation for (7) is quite straightforward. This is because if one treats the weights defining this estimator as fixed, then it is a linear estimator of a domain mean, and so its prediction variance V_i under (1) can be estimated using well-known methods (see Royall and Cumberland 1978). Since in general the EBLUP weights for the total (6) are not ‘locally calibrated’ (*i.e.*, they do not reproduce the area i mean $\bar{\mathbf{x}}_i$ of \mathbf{X}), (7) has a bias B_i under (1). A simple plug-in estimate of this bias is the difference between (7) and $\bar{\mathbf{x}}'_i \hat{\beta}$. The final MSE estimator used with (7) is therefore defined by summing the estimate of V_i and the square of this estimate of B_i . This method of MSE estimation has been empirically demonstrated to have good model-based as well as repeated sampling properties. See Chandra and Chambers (2005, 2009), Chambers and Tzavidis (2006), Chandra, Salvati and

Chambers (2007) and Tzavidis, Salvati, Pratesi and Chambers (2008).

3. Small area estimation under transformation

In this section we extend the MBD approach to SAE when the underlying regression relationships are non-linear. In doing so, we shall focus on the important case where the population values of Y follow a non-linear model in their original (raw) scale, but their logarithms can be modelled linearly. The extension to other ‘transform to linear’ models is straightforward.

Without loss of generality, suppose that both Y and X are scalar and strictly positive, with skewed population marginal distributions and clear evidence of non-linearity in their relationship, *e.g.*, as in many business surveys applications. Furthermore, a linear mixed model is appropriate for characterising how the regression of $\log(Y)$ on $\log(X)$ varies between the small areas. That is, for $i = 1, \dots, D$; $j = 1, \dots, N_i$ we have

$$l_{ij} = \log(y_{ij}) = \beta_0 + \beta_1 \log(x_{ij}) + \mathbf{g}'_{ij} \mathbf{u}_i + e_{ij} \quad (9)$$

where y_{ij} and x_{ij} are the values of Y and X respectively for population unit j in small area i , \mathbf{g}_{ij} denotes a ‘contextual’ covariate of dimension q , \mathbf{u}_i denotes a random effect for area i also of dimension q and e_{ij} is a scalar individual random effect. As usual with this type of model, we assume that all random effects are normally distributed and mutually uncorrelated, with zero expected values, $\text{Var}(\mathbf{u}_i) = \Sigma_u$ and $\text{Var}(e_{ij}) = \sigma_e^2$. Here Σ_u is the $q \times q$ matrix of covariances for the random effects. Note that $\text{Var}(l_{ij} | x_{ij}) = v_{ij} = \mathbf{g}'_{ij} \Sigma_u \mathbf{g}_{ij} + \sigma_e^2$ and $\text{Cov}(l_{ij}, l_{ik} | x_{ij}, x_{ik}, \mathbf{g}_{ij}, \mathbf{g}_{ik}) = v_{ijk} = \mathbf{g}'_{ij} \Sigma_u \mathbf{g}_{ik}$ under (9).

Given sample values of y_{ij} , x_{ij} and \mathbf{g}_{ij} , standard methods of estimation (*e.g.*, ML or REML, see Harville 1977) can be used to estimate the parameters of (9). Let $\hat{\Sigma}_u$ and $\hat{\sigma}_e^2$ denote the resulting estimates of the variance components of this linear mixed model. The estimate of $\beta = (\beta_0 \beta_1)'$ is then

$$\hat{\beta} = \left(\sum_i \mathbf{d}'_i \hat{\mathbf{v}}_{iss}^{-1} \mathbf{d}_{is} \right)^{-1} \left(\sum_i \mathbf{d}'_i \hat{\mathbf{v}}_{iss}^{-1} \mathbf{l}_{is} \right) \quad (10)$$

where $\hat{\mathbf{v}}_{iss}$, \mathbf{d}_{is} and \mathbf{l}_{is} are the sample components of $\hat{\mathbf{v}}_i = [\hat{v}_{ijk}] = \mathbf{g}_i \hat{\Sigma}_u \mathbf{g}'_i + \hat{\sigma}_e^2 \mathbf{I}_i$, $\mathbf{d}_i = [d_{ijk}] = [\mathbf{1}_i \log(\mathbf{x}_i)]$ and $\mathbf{l}_i = (l_{ij}; j = 1, \dots, N_i)$ respectively. Here \mathbf{g}_i is the $N_i \times q$ matrix defined by the covariates \mathbf{g}_{ij} in area i , \mathbf{I}_i is the identity matrix of order N_i , $\mathbf{1}_i$ denotes a vector of ones of dimension N_i and $\log(\mathbf{x}_i)$ denotes the vector of N_i values of $\log(X)$ in area i .

Note that when the variance components Σ_u and σ_e^2 are known, (10) is the BLUE for β . Consequently, $E(\hat{\beta}) \approx \beta$

and $\text{Var}(\hat{\beta}) \approx (\sum_i \mathbf{d}'_i \hat{\mathbf{V}}_{iss}^{-1} \mathbf{d}_i)^{-1}$. Put $\hat{\phi}_i = (\hat{\phi}_{ij}) = \mathbf{d}_i \hat{\beta}$. Then $E(\hat{\phi}_i) \approx \mathbf{d}_i \beta$ and $\text{Var}(\hat{\phi}_i) = \mathbf{A}_i = [a_{ijk}] \approx \mathbf{d}_i (\sum_g \mathbf{d}'_g \hat{\mathbf{V}}_{gss}^{-1} \mathbf{d}_g)^{-1} \mathbf{d}'_i$, where $a_{ijk} = \mathbf{d}'_{ij} \text{Var}(\hat{\beta}) \mathbf{d}_{ik} \rightarrow 0$ as $n \rightarrow \infty$.

Our aim is to use the log scale linear mixed model (9) for estimation of the small area means m_{iy} . In particular, we use model calibration (Wu and Sitter 2001) based on this model to develop sample weights for use in the MBD estimator (7) of this quantity.

4. Model calibrated weighting

Model calibration was introduced by Wu and Sitter (2001) as a model-assisted method of calibrated weighting when the underlying regression relationship is non-linear. Here we provide a model-based perspective on the method, as a precursor to using it for constructing weights for use in an MBD estimator in a similar situation.

Suppose that the underlying population model is non-linear, with the relationship between Y and \mathbf{X} in the population of form

$$E(y_j | \mathbf{x}_j) = h(\mathbf{x}_j; \eta) \text{ and } \text{Var}(y_j | \mathbf{x}_j) = \sigma_j^2. \quad (11)$$

Here $j = 1, \dots, N$, η (typically vector-valued) and σ_j^2 are unknown model parameters and the mean function $h(\mathbf{x}_j; \eta)$ is a known function of \mathbf{x}_j and η . We also assume that population units are mutually uncorrelated given their respective values of \mathbf{X} . Note that (11) is quite general, and includes linear, non-linear, and generalized linear models as special cases. In this situation, Wu and Sitter (2001) define the model-calibrated estimator of the population total t_{Uy} as $\hat{t}_{Uy}^{mc} = \sum_{j \in s} w_j^{mc} y_j$, where the vector of weights $\mathbf{w}_s^{mc} = (w_j^{mc})$ is chosen to minimise an appropriately chosen measure of the distance from \mathbf{w}_s^{mc} to the vector of Horvitz-Thompson weights $\mathbf{w}_s^\pi = (\pi_j^{-1})$, subject to the model calibration constraints

$$\sum_{j \in s} w_j^{mc} = N$$

and (12)

$$\sum_{j \in s} w_j^{mc} h(\mathbf{x}_j; \hat{\eta}_\pi) = \sum_{j \in U} h(\mathbf{x}_j; \hat{\eta}_\pi)$$

with $\hat{\eta}_\pi$ a design consistent estimator of η . Note that unlike standard calibration, the constraints (12) require that we know the individual population values of \mathbf{X} . The key idea behind this approach is that provided (11) fits reasonably, then y_j is (at least approximately) a linear function of its fitted value $h(\mathbf{x}_j; \hat{\eta}_\pi)$ under this model and so we can carry out linear estimation using these fitted values as auxiliary information.

A model-based perspective on model calibration can be developed as follows. Let $\hat{\eta}$ denote a ‘model-efficient’ estimator of η in (11), e.g., its maximum likelihood (ML) estimator, with associated fitted values $h(\mathbf{x}_j; \hat{\eta})$. In general, these fitted values will not be unbiased. They will also be correlated. However, there will still be a systematic relationship between the actual values of Y and their corresponding fitted values that we can approximate. Although there is nothing to stop us looking at more complex approximations, a linear model for the relationship between the population values y_j and the fitted values $\hat{y}_j = h(\mathbf{x}_j; \hat{\eta})$ seems a reasonable starting point. We therefore replace the non-linear model (11) by the linear model

$$E(y_j | \hat{y}_j) = \alpha_0 + \alpha_1 \hat{y}_j \quad (13)$$

and

$$\text{Cov}(y_j, y_k | \hat{y}_j, \hat{y}_k) = \omega_{jk}.$$

We refer to (13) as the ‘fitted value’ model corresponding to (11). Let \mathbf{J}_U denote the population ‘design matrix’ under (13), i.e., $\mathbf{J}_U = [\mathbf{1}_U \hat{\mathbf{y}}_U]$, where $\mathbf{1}_U$ denotes the unit vector of size N and $\hat{\mathbf{y}}_U = (\hat{y}_j; j = 1, \dots, N)$, and put $\Omega_U = [\omega_{jk}; j = 1, \dots, N; k = 1, \dots, N]$. We can then partition \mathbf{J}_U and Ω_U according to sample (s) and non-sample (r) units as

$$\mathbf{J}_U = \begin{bmatrix} \mathbf{J}_s \\ \mathbf{J}_r \end{bmatrix}$$

and

$$\Omega_U = \begin{bmatrix} \Omega_{ss} & \Omega_{sr} \\ \Omega_{rs} & \Omega_{rr} \end{bmatrix},$$

and hence write down the weights that define the BLUP of t_{Uy} under (13). These are the model-based model-calibrated weights

$$\begin{aligned} \mathbf{w}^{mbmc} &= (w_j^{mbmc}; j \in s) \\ &= \mathbf{1}_s + \mathbf{H}'_{cm} (\mathbf{J}'_U \mathbf{1}_U - \mathbf{J}'_s \mathbf{1}_s) + (\mathbf{I}_s - \mathbf{H}'_{cm} \mathbf{J}'_s) \Omega_{ss}^{-1} \Omega_{sr} \mathbf{1}_r \end{aligned} \quad (14)$$

where $\mathbf{H}_{mc} = (\mathbf{J}'_s \Omega_{ss}^{-1} \mathbf{J}_s)^{-1} \mathbf{J}'_s \Omega_{ss}^{-1}$. Clearly, these weights are model-calibrated since $\sum_{j \in s} w_j^{mbmc} = N$ and $\sum_{j \in s} w_j^{mbmc} \hat{y}_j = \sum_{j \in U} \hat{y}_j$. However, unlike the linear model EBLUP weights (2), they are *not* calibrated on \mathbf{X} . In practice, the components of Ω_U will not be known and will need to be estimated. When these estimates are substituted in (14), we obtain the empirical version \mathbf{w}^{embmc} of these model-calibrated weights.

5. Model calibrated weighting for small area estimation

We now use model calibration based on the log scale linear mixed model (9) to obtain sample weights for use in the MBD estimator (7). From the development in the previous section it can be seen that this requires us to first specify a fitted value model (13) for Y based on (9), *i.e.*, we need to calculate appropriate fitted values \hat{y}_{ij} as well as estimates $\hat{\omega}_{ijk}$ of $\omega_{ijk} = \text{Cov}(y_{ij}, y_{ik} | x_{ij}, x_{ik}, \mathbf{g}_{ij}, \mathbf{g}_{ik})$ under (9). The sample weights to use in the MBD estimator (7) are then given by (14).

A simple method of defining fitted values \hat{y}_{ij} under (9) is one where parameter estimates derived under this model are used to obtain predicted values on the log scale which are then back-transformed. Unfortunately, as is well known, this approach is biased. We therefore develop the first and second order moments of an appropriate bias-corrected fitted value model based on (9). Let \mathbf{x}_s and \mathbf{g}_s denote the sample values of x_{ij} and \mathbf{g}_{ij} respectively. Under (9),

$$E(y_{ij} | x_{ij}, \mathbf{g}_{ij}) = E\{e^{y_{ij}} | x_{ij}, \mathbf{g}_{ij}\} = e^{\phi_{ij} + v_{ij}/2} \neq E(e^{\hat{\phi}_{ij} + \hat{v}_{ij}/2} | \mathbf{x}_s, \mathbf{g}_s) = E(\hat{y}_{ij} | x_{ij}, \mathbf{g}_{ij})$$

so the usual bias correction that makes use of the fact that the conditional distribution of y_{ij} is lognormal is inadequate. Let $\hat{\eta}_{ij} = (\hat{\beta}, \hat{v}_{ij})'$ be an estimate of $\eta_{ij} = (\beta, v_{ij})'$ such that $E(\hat{\eta}_{ij} - \eta_{ij}) \approx 0$ for large n . Put $z(\eta_{ij}) = e^{\phi_{ij} + v_{ij}/2}$. Using a second order Taylor series approximation we can write

$$z(\hat{\eta}_{ij}) \approx z(\eta_{ij}) + (\hat{\eta}_{ij} - \eta_{ij})' z^{(1)}(\eta_{ij}) + \frac{1}{2} (\hat{\eta}_{ij} - \eta_{ij})' z^{(2)}(\eta_{ij}) (\hat{\eta}_{ij} - \eta_{ij})$$

and so

$$E\{z(\hat{\eta}_{ij})\} \approx z(\eta_{ij}) + \frac{1}{2} \text{tr}[E\{z^{(2)}(\eta_{ij}) (\hat{\eta}_{ij} - \eta_{ij}) (\hat{\eta}_{ij} - \eta_{ij})'\}]$$

Here

$$z^{(1)}(\eta_{ij}) = \left(\mathbf{d}'_{ij} e^{\phi_{ij} + v_{ij}/2} \quad \frac{1}{2} e^{\phi_{ij} + v_{ij}/2} \right)'$$

and

$$z^{(2)}(\eta_{ij}) = \begin{pmatrix} \mathbf{d}_{ij} \mathbf{d}'_{ij} e^{\phi_{ij} + v_{ij}/2} & \frac{1}{2} \mathbf{d}_{ij} e^{\phi_{ij} + v_{ij}/2} \\ \frac{1}{2} \mathbf{d}'_{ij} e^{\phi_{ij} + v_{ij}/2} & \frac{1}{4} e^{\phi_{ij} + v_{ij}/2} \end{pmatrix}$$

are the vector and matrix respectively containing the first and second order derivatives of $z(\eta_{ij})$ with respect to η_{ij} . Since the asymptotic covariance between ML (or REML) estimators of the fixed and variance components of a linear mixed model is zero (McCulloch and Searle 2001, chapter 2, pages 40 - 45), the covariance between $\hat{\beta}$ and \hat{v}_{ij} will be negligible. It follows that

$$\begin{aligned} & \text{tr}[E\{z^{(2)}(\eta_{ij})(\hat{\eta}_{ij} - \eta_{ij})(\hat{\eta}_{ij} - \eta_{ij})'\}] \\ &= \text{tr}[z^{(2)}(\eta_{ij})E\{(\hat{\eta}_{ij} - \eta_{ij})(\hat{\eta}_{ij} - \eta_{ij})'\}] \\ &\approx e^{\phi_{ij} + v_{ij}/2} \left[\mathbf{d}'_{ij} \left(\sum_g \mathbf{d}'_{gs} \hat{\mathbf{v}}_{gss}^{-1} \mathbf{d}_{gs} \right)^{-1} \mathbf{d}_{ij} + \frac{1}{4} \text{Var}(\hat{v}_{ij}) \right] \\ &= E(y_{ij} | x_{ij}, \mathbf{g}_{ij}) \left[\hat{a}_{ij} + \frac{1}{4} \text{Var}(\hat{v}_{ij}) \right] \end{aligned}$$

where $\hat{a}_{ij} = \mathbf{d}'_{ij} \hat{\mathbf{V}}(\hat{\beta}) \mathbf{d}_{ij}$ and $\hat{\mathbf{V}}(\hat{\beta}) = (\sum_i \mathbf{d}'_{is} \hat{\mathbf{v}}_{iss}^{-1} \mathbf{d}_{is})^{-1}$ is the usual estimator of $\text{Var}(\hat{\beta})$. Our fitted values are therefore defined by the second order bias corrected estimator of $E(y_{ij} | x_{ij}, \mathbf{g}_{ij})$,

$$\hat{y}_{ij} = h(\mathbf{d}_{ij}; \hat{\eta}_{ij}) = \hat{k}_{ij}^{-1} e^{\hat{\phi}_{ij} + \hat{v}_{ij}/2} \tag{15}$$

where

$$\hat{k}_{ij} = 1 + \frac{1}{2} \left\{ \hat{a}_{ij} + \frac{1}{4} \hat{\mathbf{V}}(\hat{v}_{ij}) \right\}$$

and $\hat{\mathbf{V}}(\hat{v}_{ij})$ is the estimated asymptotic variance of \hat{v}_{ij} . Under ML and REML estimation of the variance components of (9), this estimated asymptotic variance is obtained from the inverse of the relevant information matrix. Note that the bias adjustment of Karlberg (2000a) is a special case of (15).

In order to use (14) to define model-based model-calibrated sample weights, we also need estimates of the second order moments of the population values of Y given these fitted values. The conditional moments ω_{ijk} are a first order approximation to these moments. In particular, given normal random effects

$$\omega_{ijk} = e^{(\phi_{ij} + \phi_{ik}) + (v_{ij} + v_{ik})/2} (e^{v_{ijk}} - 1) \tag{16}$$

Our estimate $\hat{\omega}_{ijk}$ of ω_{ijk} is obtained by substituting $\hat{\phi}_{ij}$ and \hat{v}_{ijk} for ϕ_{ij} and v_{ijk} in (16).

The empirical model-based model-calibrated weights (14) corresponding to the fitted value model defined by (15) and (16) are

$$\begin{aligned} \mathbf{w}^{embmc} &= (w_{ij}^{embmc}; j \in s_i; i = 1, \dots, D) \\ &= \mathbf{1}_s + \hat{\mathbf{H}}'_{mc} (\mathbf{J}'_U \mathbf{1}_U - \mathbf{J}'_s \mathbf{1}_s) \\ &\quad + (\mathbf{I}_s - \hat{\mathbf{H}}'_{mc} \mathbf{J}'_s) \hat{\Omega}_{ss}^{-1} \hat{\Omega}_{sr} \mathbf{1}_r. \end{aligned} \tag{17}$$

Here $\mathbf{J}_U = [\mathbf{1}_U \hat{\mathbf{y}}_U]$, so

$$\mathbf{J}'_U \mathbf{1}_U - \mathbf{J}'_s \mathbf{1}_s = \begin{pmatrix} N - n \\ \sum_i \sum_{j \in r_i} \hat{y}_{ij} \end{pmatrix},$$

and $\hat{\mathbf{H}}_{mc} = (\mathbf{J}'_s \hat{\Omega}_{ss}^{-1} \mathbf{J}_s)^{-1} \mathbf{J}'_s \hat{\Omega}_{ss}^{-1}$. Also $\hat{\Omega}_{ss} = \text{diag}\{\hat{\Omega}_{iss}; i = 1, \dots, D\}$ and $\hat{\Omega}_{sr} = \text{diag}\{\hat{\Omega}_{isr}; i = 1, \dots, D\}$, where $\hat{\Omega}_{iss}$ and $\hat{\Omega}_{isr}$ are defined by the sample/non-sample decomposition of $\hat{\Omega}_i$. For example, when (9) corresponds to a random intercepts specification, $\hat{v}_{ijk} = \hat{\sigma}_u^2 + \hat{\sigma}_e^2 I(j = k)$ and so the components of $\hat{\Omega}_i$ are

$$\hat{\omega}_{ijk} = e^{\hat{\phi}_{ij} + \hat{\phi}_{ik} + \hat{\sigma}_u^2 + \hat{\sigma}_e^2} [e^{\hat{\sigma}_e^2} \{1 + I(j = k)(e^{\hat{\sigma}_e^2} - 1)\} - 1].$$

The development so far has assumed normality of log-scale random effects. However, there is no good reason (beyond convenience) to assume that with skewed data these random area effects should be normal. One alternative, given a scalar area effect in (9), is to assume that the random effects in this model are drawn from the *gamma* family of distributions. From the properties of this distribution and using binomial and exponential expansions (ignoring higher order terms) we can show that $E(y_{ij} | x_{ij}, \mathbf{g}_{ij}) \approx e^{\hat{\phi}_{ij} + v_{ij}/2} = z(\eta_{ij})$ as in the normal case. This indicates that an MBD estimator based on the model-based model-calibrated weights (17) should be robust with respect to the distribution of the random effects in (9).

Finally, we consider definition of the MBD estimator itself. As noted in section 2, this estimator is just the weighted average of the sample Y -values in an area. However, use of such a weighted average pre-supposes that the weights are reasonably close to being ‘locally calibrated on N ’, *i.e.*, when summed over the sample units in small area i we obtain a value that is not too different from the actual small area population size N_i . This property usually holds if the weights are the EBLUP weights for the total (6) defined by a linear mixed model for Y . It does not necessarily hold for the model-based model-calibrated weights (17). Consequently, we consider two specifications for the MBD estimator given these weights. The first, which we refer to as a ‘Hájek specification’, is just the weighted average (7), with weights defined by (17). The second, which we refer to as a ‘Horvitz-Thompson specification’, replaces the denominator in (7) by the actual value of N_i . That is, the two types of MBD estimator under model-based model-calibrated weighting that we consider are

$$\hat{m}_{iy}^{\text{HJ-TrMBD}} = \left\{ \sum_{j \in s_i} w_{ij}^{\text{embmc}} \right\}^{-1} \sum_{j \in s_i} w_{ij}^{\text{embmc}} y_{ij} \quad (18)$$

and

$$\hat{m}_{iy}^{\text{HT-TrMBD}} = N_i^{-1} \sum_{j \in s_i} w_{ij}^{\text{embmc}} y_{ij}. \quad (19)$$

Alternatively we can adopt a prediction-based approach to obtain an alternative indirect predictor for the small area mean under the log-transformed model (9). Our approach extends that of Karlberg (2000a). In this case, assuming model (9) holds, we predict each nonsample Y in small area i and then sum these predictions. Note that we need to correct for bias following back-transformation to the raw scale when calculating these predicted values for the nonsample Y . Under model (9), the resulting empirical predictor for the mean m_{iy} of Y in area i (denoted TrEP) can be defined as

$$\hat{m}_{iy}^{\text{TrEP}} = N_i^{-1} \left\{ \sum_{j \in s_i} y_{ij} + \sum_{j \in r_i} \hat{y}_{ij} \right\}, \quad (20)$$

where \hat{y}_{ij} is given by (15).

Estimation of the MSE of (18) and (19) is carried out in the usual way for MBD estimators, *i.e.*, via the MSE estimation approach described in section 2. Estimation of the MSE of (20) is not straightforward since this predictor is a non-linear function of Y values. We do not pursue this issue in this paper.

6. An empirical evaluation

In this section we provide empirical results on the comparative performances of five different methods of SAE. These are the two ‘transformation-based’ MBD estimators (18) and (19), both based on the model-based model-calibrated weights (17) and denoted by HJ-TrMBD and HT-TrMBD respectively; the log-transformation based predictor (20) under model (9), denoted TrEP, the ‘standard’ MBD estimator (7) based on the linear mixed model (3) and the empirical EBLUP weights for the total (6), which we denote by HJ-LinMBD to emphasise that it is a Hájek-type weighted mean based on weights derived under a linear mixed model; and the EBLUP (8) derived under the same linear mixed model, which we denote HT-LinEBLUP. Note that the MSEs for all three MBD estimators were estimated using the method described in section 2, while the MSE of HT-LinEBLUP was estimated using the method described in Prasad and Rao (1990). Note that we have not considered estimation of the MSE of TrEP.

Our empirical results are based on two types of simulation studies. The first type used model-based simulation to generate artificial population and sample data. That is, at each simulation population data were first generated under the model and a single sample was then taken from this simulated population by stratified simple random sampling without replacement with small area as strata. These data were then used to compare the performances of the different estimators. In section 6.1 we present the results from these model-based simulations. We carried out two

sets of model-based simulations. In the first set of simulations (Set A), we investigated the performance of these estimators given population data generated using the log-scale linear mixed model (9). In second set of simulations (Set B), we examined the robustness of these estimators to misspecification of this model. The second type of simulation study was design-based. In section 6.2 we describe design-based simulations. Here we evaluated these estimators in the context of repeated sampling from a real population using realistic sampling methods. That is, real survey data were first used to simulate a population, and this fixed population was then repeatedly sampled according to a pre-specified design. In particular, the sample design used was stratified random sampling with strata corresponding to the small areas of interest and with stratum allocations set to the small area sample sizes in the original datasets.

Four measures of estimator performance were computed using the various estimates generated in these simulation studies. They were the relative bias (RB) and the relative root mean squared error (RRMSE) of these estimates, together with the coverage rate and average width of the nominal 95 per cent confidence intervals based on them. In Tables 2 to 4 these measures are presented as averages over the small areas of interest.

6.1 The model-based simulation study

Model-based simulations are a common way of illustrating the sensitivity of an estimation procedure to variation in assumptions about the structure of the population of interest. Here we fixed the population size at $N = 15,000$ and randomly generated the small area population sizes $N_i, i = 1, \dots, D = 30$ so that $\sum_i N_i = N$. We used an overall sample size of $n = 600$ with small area sample sizes set so that they were proportional to the corresponding small area population sizes. These area-specific population and sample sizes were kept fixed in all our simulations. The population and sample sizes are given in Table 1a.

Table 1a
Area specific population (N_i) and sample (n_i) sizes for model-based simulation

Area	1	2	3	4	5	6	7	8	9	10
N_i	525	538	510	468	526	484	516	458	529	518
n_i	21	22	20	19	21	19	21	19	21	21
Area	11	12	13	14	15	16	17	18	19	20
N_i	502	524	509	484	487	459	542	498	512	500
n_i	20	21	20	19	19	18	22	20	20	20
Area	21	22	23	24	25	26	27	28	29	30
N_i	497	492	443	506	513	536	506	495	463	460
n_i	20	20	18	20	21	21	20	20	19	18

In Set A of our model-based simulations the population values y_{ij} were generated using the multiplicative model $y_{ij} = 5.0x_{ij}^\beta u_i e_{ij} (j = 1, \dots, N_i; i = 1, \dots, 30)$, with random samples then taken from each small area. Here the values of x_{ij} were independently drawn from the log-normal distribution $\log(x_{ij}) \sim N(6, \sigma_x^2)$, with the individual effects and area effects independently drawn as $\log(e_{ij}) \sim N(0, \sigma_e^2)$ and $\log(u_i) \sim N(0, \sigma_u^2)$ respectively. The population values of x were re-generated in each simulation. In particular, in each simulation we first generated the values of x 's for a population of size N and then randomly assigned these values to different areas of sizes N_i . The values of σ_e and σ_u were chosen so that the intra-area correlation in the population varied between 0.20 and 0.25. Table 1b shows the six different sets of parameter values that were used in Set A. These ensured that the simulated populations contained a wide range of variation. For each generated population and for each area i we selected a simple random sample (without replacement) of size n_i , leading to an overall sample size of $n = 600$. The sample values of y and the population values of x obtained in each simulation were then used to estimate the small area means. That is, using the sample data in each case, parameter values were estimated using the *lme* function in R (Bates and Pinheiro 1998), and estimates for the small area means then calculated, along with appropriate nominal 95% confidence intervals. The process of generating population and sample data, estimation of parameters and calculation of small area estimates was independently replicated 1,000 times. The results from this part of the simulation study are shown in Table 2.

Table 1b
Population specifications for model-based simulation Set A

Parameter Set	β	σ_u	σ_e	σ_x
1	0.5	0.30	0.50	3.00
2	0.8	0.35	0.60	2.50
3	1.0	0.40	0.70	2.25
4	1.3	0.45	0.80	1.75
5	1.5	0.50	0.90	1.50
6	2.0	0.60	1.00	1.20

In Set B of the model-based simulations, population data were generated using the model $y_{ij} = 5.0x_{ij} [\exp(\log^2(x_{ij}))]^\gamma u_i e_{ij}$. Here the individual effects e_{ij} and the area effects u_i were independently drawn as $\log(e_{ij}) \sim N(0, 1)$ and $\log(u_i) \sim N(0, 0.25)$ respectively, while the covariate values x_{ij} were drawn as $\log(x_{ij}) \sim N(3, 0.04)$. Five different values for the parameter γ (-1.0, -0.5, 0.0, 0.5, 1.0) were investigated, thus generating population data with different degrees of curvature. All other aspects of these simulations, including the estimators considered, were the same as in Set A. Table 3 presents results from this component of the simulation study.

Table 2
Average relative bias (ARB), average relative RMSE (ARRMSE), average coverage rate (ACR) and average interval width (AW) for model-based simulation Set A

Criterion	Estimator	Parameter Set					
		1	2	3	4	5	6
ARB,%	HJ-TrMBD	-82.68	-95.02	-98.08	-98.50	-98.29	-99.00
	HT-TrMBD	0.09	0.10	-0.14	-0.25	-0.03	0.04
	TrEP	0.08	0.09	-0.18	-0.48	-0.05	0.01
	HJ-LinMBD	12.01	4.09	-1.35	-5.54	-6.60	-9.88
	HT-LinEBLUP	13.39	5.18	-0.67	-5.24	-6.41	-9.67
ARRMSE	HJ-TrMBD	4.80	1.39	1.25	1.44	1.42	1.62
	HT-TrMBD	0.15	0.26	0.45	0.64	0.66	0.91
	TrEP	0.30	0.41	0.58	0.80	0.81	1.09
	HJ-LinMBD	1.11	1.41	1.85	1.99	2.06	2.69
	HT-LinEBLUP	0.79	0.54	0.64	0.92	0.93	1.31
ACR	HJ-TrMBD	0.99	0.98	0.97	0.95	0.94	0.92
	HT-TrMBD	0.94	0.91	0.89	0.89	0.89	0.88
	HJ-LinMBD	0.87	0.85	0.85	0.88	0.88	0.87
	HT-LinEBLUP	0.85	0.85	0.86	0.87	0.88	0.87
AW	HJ-TrMBD	1,592	22,688	140,452	52×10^4	35×10^5	44×10^6
	HT-TrMBD	219	4,414	34,105	14×10^4	11×10^5	15×10^6
	HJ-LinMBD	1,005	19,232	139,420	57×10^4	41×10^5	56×10^6
	HT-LinEBLUP	382	7,099	57,039	26×10^4	21×10^5	32×10^6

Table 3
Average relative bias (ARB), average relative RMSE (ARRMSE), average coverage rate (ACR) and average interval width (AW) for model-based simulation Set B

Criterion	Estimator	$\gamma = -1.0$	$\gamma = -0.5$	$\gamma = 0.0$	$\gamma = 0.5$	$\gamma = 1.0$
		ARB,%	HT-TrMBD	4.92	0.66	0.14
HJ-LinMBD	-0.21		0.04	0.12	0.16	-0.85
HT-LinEBLUP	-0.19		0.04	0.13	0.17	-0.77
ARRMSE	HT-TrMBD	0.38	0.35	0.33	0.37	0.41
	HJ-LinMBD	0.56	0.36	0.34	0.53	1.20
	HT-LinEBLUP	0.38	0.30	0.29	0.36	0.56
ACR	HT-TrMBD	0.94	0.92	0.92	0.91	0.87
	HJ-LinMBD	0.91	0.92	0.92	0.92	0.90
	HT-LinEBLUP	0.93	0.94	0.94	0.93	0.92
AW	HT-TrMBD	0.04	2.50	211	29,070	5×10^6
	HJ-LinMBD	0.06	2.70	214	38,660	13×10^6
	HT-LinEBLUP	0.05	2.60	214	33,442	10×10^6

6.2 The design-based simulation study

This study used the same population and samples as the simulation studies described in Chandra and Chambers (2005) and Chambers and Tzavidis (2006), which was based on data obtained from a sample of 1,652 farms that participated in the Australian Agricultural and Grazing Industries Survey (AAGIS). A realistic population of 81,982 farms was defined by sampling with replacement from the original sample of 1,652 farms with probabilities proportional to their sample weights, all of which were strictly

greater than one. A total of 1,000 independent samples, each of size $n = 1,652$, were drawn from this fixed population by simple random sampling without replacement within strata defined by the 29 Australian agricultural regions represented in the AAGIS sample. These regions are the small areas of interest. Regional sample sizes were fixed to be the same as in this original sample, varying from a low of 6 to a high of 117, which allows an evaluation of the performance of the different estimation methods across a range of realistic small area sample sizes. Note that sampling fractions in these strata also varied disproportionately, ranging between 0.70

and 15.87 percent. The aim is to estimate average annual farm costs (TCC, measured in A\$) in each region using farm size (hectares) as the auxiliary variable. The same mixed model specification as in Chandra and Chambers (2005) is used. This includes an interaction term (zone by size) in the fixed effects and a random slope specification for the area effect. In its linear form the model does not fit the AAGIS sample data terribly well. This fit is improved (albeit marginally) when a log-scale linear specification is used. Our results are summarized in Table 4.

6.3 Discussion of simulation results

The most striking feature of Table 2 is the extremely large values of the averages relative bias of HJ-TrMBD under model-based model-calibrated weighting. The two best performers with respect to relative bias are HT-TrMBD, which is based on the same weights as HJ-TrMBD, and TrEP. An investigation of the reason for the poor performance of HJ-TrMBD revealed that summing the model-based model-calibrated weights (17) within small areas produced extremely variable estimates of the small area population sizes, implying that these weights cannot be considered as ‘multipurpose’ – they function well when used with variables that are reasonably correlated with the variable that defines the fitted value model, but can fail with other, less well correlated, variables (*e.g.*, the indicator variable for small area inclusion). We further note that this problem does not arise with the ‘standard’ empirical EBLUP weights for the total (6), as HJ-LinMBD performs consistently for all six of the scenarios explored in Set A of the simulation study. From now on we therefore focus our discussion on the four estimators, HT-TrMBD, TrEP, HJ-LinMBD and HT-LinEBLUP.

Table 2 shows that the average relative biases and the average relative RMSEs for HT-TrMBD are consistently lower than those generated by HJ-LinMBD and HT-LinEBLUP. The average relative biases of HT-TrMBD and TrEP are comparable. However, the average relative RMSEs of HT-TrMBD are consistently smaller than the TrEP. Furthermore, average coverage rates and interval widths for HT-TrMBD are better than those generated by HJ-LinMBD and HT-LinEBLUP. In comparison, for the same order of relative bias, the relative RMSEs of HT-LinEBLUP is smaller than that of HJ-LinMBD, and, although both estimators generate very similar coverage rates, confidence intervals generated via HT-LinEBLUP tend to have smaller average widths than those generated via HJ-LinMBD.

The plots in Figure 1 display the region-specific performance measures generated by these four estimators for the Set A simulations. These show that the relative bias and the relative RMSE values generated by HT-TrMBD are smaller than corresponding values for HJ-LinMBD and HT-LinEBLUP in all regions. With almost identical values of relative biases, the HT-TrMBD has smaller values of relative RMSEs than corresponding values for TrEP in all regions. Further, the relative bias and the relative RMSE of HJ-LinMBD and HT-LinEBLUP increase as the non-linearity in the data increases (*i.e.*, as we move from parameter set 1 to parameter set 6). We also see that HT-TrMBD generates better coverage rates across all regions compared with the coverage rates generated by HT-LinEBLUP and HJ-LinMBD.

Table 4
Average relative bias (ARB), average relative RMSE (ARRMSE) and average coverage rate (ACR) for design-based simulation using AAGIS data. Simulation standard errors of ARB and ARRMSE are shown in parentheses

Criterion	Estimator	Average of 29 regions	Average of 28 regions
ARB, %	HT-TrMBD	1.96 (0.20)	1.92 (0.11)
	HJ-LinMBD	-2.13 (0.15)	-2.21 (0.12)
	HT-LinEBLUP	2.98 (0.18)	3.36 (0.16)
	PseudoEBLUP	4.01 (0.22)	4.41 (0.20)
	JL	1.89 (0.19)	2.23 (0.17)
ARRMSE, %	HT-TrMBD	21.93 (4.47)	17.41 (1.18)
	HJ-LinMBD	20.15 (3.80)	16.91 (2.20)
	HT-LinEBLUP	19.87 (1.78)	19.30 (1.63)
	PseudoEBLUP	22.42 (2.52)	21.95 (2.46)
	JL	20.97 (1.48)	20.48 (1.31)
ACR	HT-TrMBD	0.89	0.92
	HJ-LinMBD	0.93	0.95
	HT-LinEBLUP	0.85	0.85

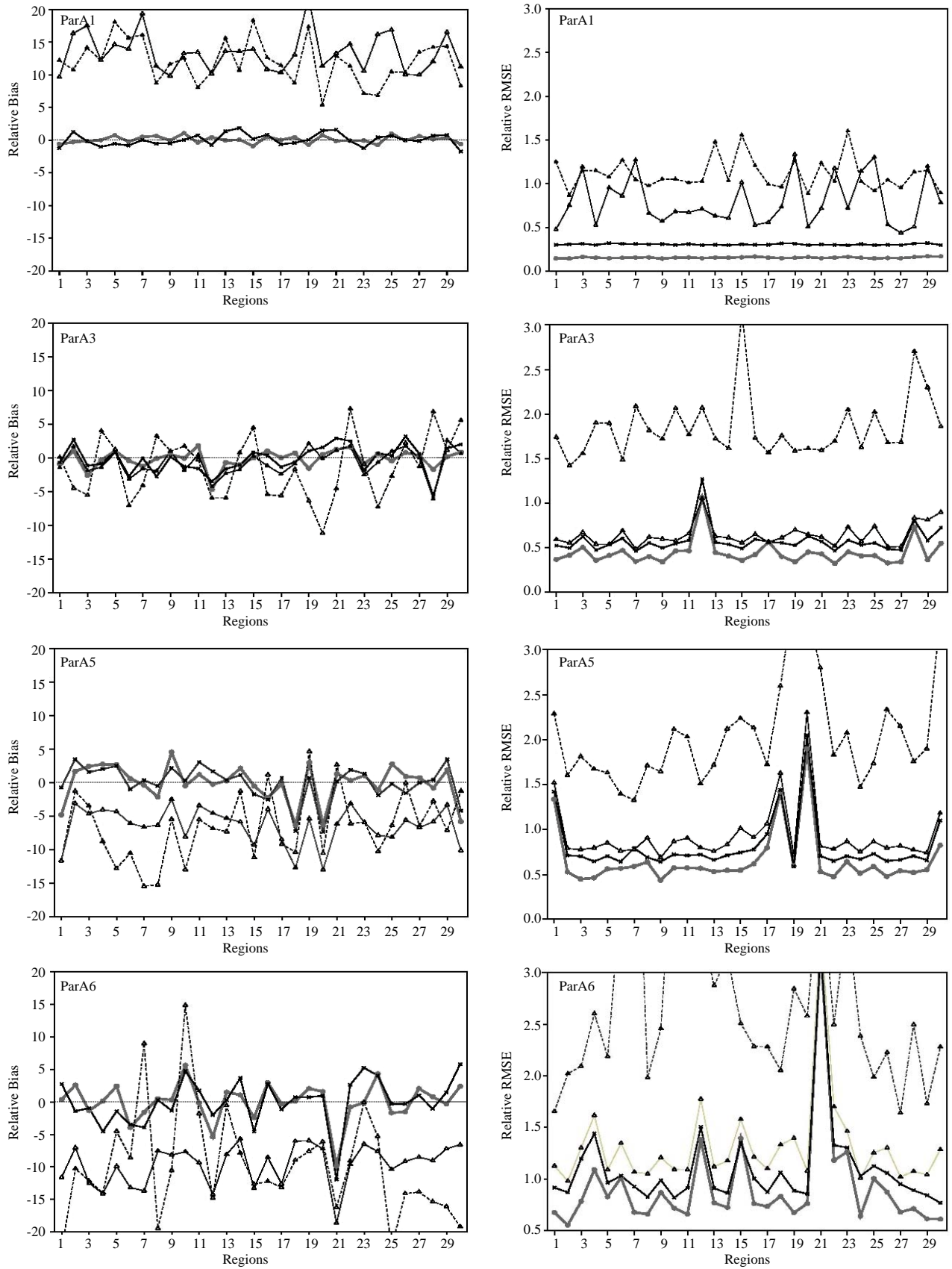


Figure 1 Area specific results for HT-TrMBD (solid line, ●), TrEP (thick line, ×), HT-LinEBLUP (thin line Δ) and HJ-LinMBD (dashed line, Δ) under parameter sets 1 (ParA1), 3 (ParA3), 5 (ParA5) and 6 (ParA6). Left column is Relative Bias (%) and right column is Relative RMSE

Overall, these results show that when the model for the underlying population is non-linear there can be significant gains from the use of HT-type MBD estimators for small area means based on the model-calibrated weights (17) compared with standard linear mixed model-based estimators like HJ-LinMBD and HT-LinEBLUP. They also show that the indirect estimator HT-LinEBLUP performs relatively better than the direct estimator HJ-LinMBD in these situations. The indirect predictor TrEP based on log-transformed model (9) performs well in terms of relative bias but is less efficient than the MBD estimator under the same model.

In Set B of the model-based simulations we investigated the robustness of model-based model-calibrated direct estimation to misspecification of the non-linear model. The results in Table 3 show that in this case the biases generated by HT-TrMBD increase as the actual non-linear model deviates more from the assumed non-linear model ($\gamma = 0.0$ in the table). However, these biases are offset by small variability, so in terms of average relative RMSE, HT-TrMBD still performs as well or better than HT-LinEBLUP and continues to dominate HJ-LinMBD. The biases generated by HJ-LinMBD and HT-LinEBLUP are of the same order, while the average relative RMSE of HT-LinEBLUP dominates that of HJ-LinMBD. Average coverage rates for HT-LinEBLUP are marginally better than those of HJ-LinMBD and HT-TrMBD, but the average widths of the confidence intervals underpinning these rates tended to be smallest for HT-TrMBD, followed by HT-LinEBLUP and then HJ-LinMBD. Overall, our model-based simulation results for Set B indicate that although MBD-based SAE with model-based model-calibrated weights is susceptible to model misspecification bias, the overall performance of this approach appears relatively unaffected by slight deviations from the assumed non-linear model.

In Table 4 and Figure 2 we present the average and region-specific performance measure generated by different SAE methods for AAGIS data respectively. These results show that the average relative bias of HT-TrMBD is smaller than that of both HT-LinEBLUP and HJ-LinMBD, while the average relative RMSE of HT-TrMBD is marginally larger than the corresponding values for HJ-LinMBD and HT-LinEBLUP. Inspection of Figure 2 shows that this result is essentially due to one region (21) in the original AAGIS sample that contained a massive outlier (TCC > A\$30,000,000). This outlier was included in the simulation population (twice) and then selected (in one case, twice) in 37 of the 1000 simulation samples, leading to completely unrealistic estimates for region 21 being generated by HT-TrMBD and HJ-LinMBD. The right-hand column in Table 4 therefore shows the average performances of the different

methods when this region is excluded. Here we see that now HT-TrMBD and HJ-LinMBD are essentially on a par, with both dominating HT-LinEBLUP. The fact that HT-TrMBD does not provide significant gains over HJ-LinMBD in this case reflects the fact that the raw-scale and log-scale linear mixed models used in these estimators both provide relatively poor fits to the AAGIS data.

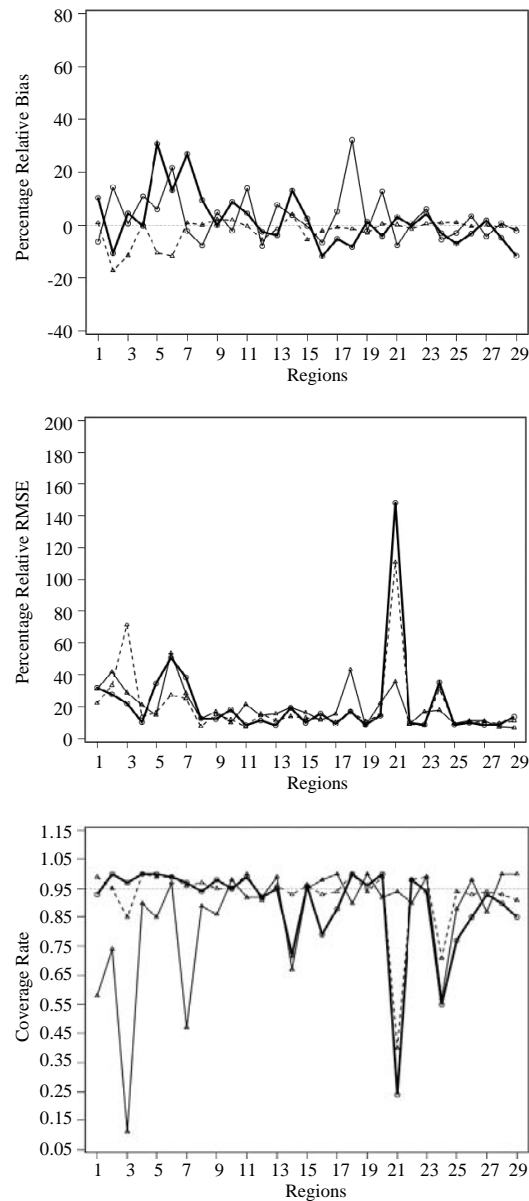


Figure 2 Region-specific simulation results for HT-TrMBD (thick line, \circ), HT-LinEBLUP (thin line Δ) and HJ-LinMBD (dashed line, Δ) in design-based simulations based on the AAGIS data. Plots show (in order from the top), RB (%), RRMSE (%) and CR. Regions are ordered in terms of increasing population size

7. Conclusions and further research

The simulation results discussed in the previous section show that combining model-based model-calibrated weights with direct estimation can bring significant gains in SAE efficiency if the population data are clearly non-linear. As one would expect, these gains are less when the assumed non-linear model is misspecified. Although we do not provide the details, our conclusions were essentially unaffected when we carried out similar simulations using gamma distributed random effects.

Our main caveat concerning the use of the model-based model-calibrated weights (17) for SAE is their specificity. These weights do not appear to have the same ‘multi-purpose’ characteristics as standard EBLUP weights for the total based on linear mixed models. Further research is therefore required on how to build model-calibrated weights for SAE that are more ‘general purpose’. It is to be expected that such weights would not be as efficient as the variable specific weights (17), but hopefully this will be more than offset by their increased utility. A further issue that is extremely important in practice is that positively skewed survey variables can also take zero (or even negative) values. For example, economic variables like debt and capital expenditure often take zero values, while variables defined as the difference of two non-negative quantities (*e.g.*, profit, which is the difference between income and expenditure) can be negative. Karlberg (2000b) uses a mixture model to characterise data that are a mix of zeros and strictly positive values. This type of model can be used in model-based model-calibrated weighting.

Finally, we note that using a transformation-based MBD approach where the usual linear model assumptions are only approximately valid (the situation considered in this paper) is not the only approach that has been suggested for this problem. Two alternative approaches in the literature are the pseudo-EBLUP (Rao 2003, section 7.2.7) and the model-assisted EB-type estimator of Jiang and Lahiri (2006). Recollect from (8) that the EBLUP is defined by replacing the unknown area i mean m_{iy} by an estimate of its expected value given the observed sample values of Y in area i and the area i values of \mathbf{X} . Let π_{ij} denote the sample inclusion probability of population unit j in small area i . The pseudo-EBLUP is then defined by replacing m_{iy} by an estimate of its expected value given the value of its design-consistent estimate

$$\hat{m}_{iy}^{\pi} = \left(\sum_{j \in s_i} \pi_{ij}^{-1} \right)^{-1} \sum_{j \in s_i} \pi_{ij}^{-1} y_{ij} = \sum_{j \in s_i} \tilde{w}_{ij} y_{ij} \quad (21)$$

and the area i values of \mathbf{X} . That is, under (3) the pseudo-EBLUP of m_{iy} is

$$\begin{aligned} \hat{m}_{iy}^{\text{psuedoEBLUP}} &= \hat{E}\{m_{iy} | \hat{m}_{iy}^{\pi}, \mathbf{x}_{is}, \mathbf{x}_{ir}\} \\ &= \bar{\mathbf{x}}'_i \hat{\beta}_{\tilde{w}} + (\bar{\mathbf{g}}'_i \hat{\Sigma}_{u\tilde{w}} \bar{\mathbf{g}}_{i\tilde{w}}) \\ &\quad \left(\bar{\mathbf{g}}'_{i\tilde{w}} \hat{\Sigma}_{u\tilde{w}} \bar{\mathbf{g}}_{i\tilde{w}} + \hat{\sigma}_{e\tilde{w}}^2 \sum_{j \in s_i} \tilde{w}_{ij}^2 \right)^{-1} (\hat{m}_{iy}^{\pi} - \bar{\mathbf{x}}'_{i\tilde{w}} \hat{\beta}_{\tilde{w}}) \end{aligned} \quad (22)$$

where $\hat{\beta}_{\tilde{w}}$, $\hat{\Sigma}_{u\tilde{w}}$ and $\hat{\sigma}_{e\tilde{w}}^2$ are pseudo-maximum likelihood estimates based on the weights \tilde{w}_{ij} and $\bar{\mathbf{g}}_{i\tilde{w}}$ and $\bar{\mathbf{x}}'_{i\tilde{w}}$ are design-consistent estimates of $\bar{\mathbf{g}}_i$ and $\bar{\mathbf{x}}_i$ that are defined in exactly the same way as \hat{m}_{iy}^{π} above. Under the same model the Jiang and Lahiri (2006) model-assisted EB-type approach leads to an estimator that is also defined by conditioning on the value of \hat{m}_{iy}^{π} ,

$$\begin{aligned} \hat{m}_{iy}^{JL} &= \sum_{j \in s_i} \tilde{w}_{ij} \hat{E}\{\hat{E}(y_{ij} | \mathbf{x}_{ij}, \mathbf{u}_i) | \hat{m}_{iy}^{\pi}, \mathbf{x}_i\} \\ &= \bar{\mathbf{x}}'_{i\tilde{w}} \hat{\beta} + \{\tilde{\mathbf{w}}'_{is} (\mathbf{g}_{is} \hat{\Sigma}_u \mathbf{g}'_{is} + \hat{\sigma}_e^2 \mathbf{I}_{is}) \tilde{\mathbf{w}}_{is}\}^{-1} \\ &\quad \{\tilde{\mathbf{w}}'_{is} \mathbf{g}_{is} \hat{\Sigma}_u \mathbf{g}'_{is} \tilde{\mathbf{w}}_{is}\} (\hat{m}_{iy}^{\pi} - \bar{\mathbf{x}}'_{i\tilde{w}} \hat{\beta}) \end{aligned} \quad (23)$$

where $\tilde{\mathbf{w}}_{is}$ is the vector of standardised sample weights \tilde{w}_{ij} in area i . Note that in (23) we use optimal (*i.e.*, ML or REML) estimates for model parameters.

Both (22) and (23) are essentially motivated by the idea of estimating the area i mean by its conditional expectation under (3) given the value of the usual design-consistent estimator (21) for this quantity. As such, they are indirect estimators like the HT-LinEBLUP. Under (3), neither will be as efficient as the HT-LinEBLUP, while if (9) rather than (3) holds, then both estimators rely on the design consistency of \hat{m}_{iy}^{π} for robustness. Since relying on a large sample property of a small sample statistic seems rather optimistic, we prefer to tackle the model specification problem directly, replacing (3) by (9) and using the transformation-based MBD approach described in section 5. Values of average relative bias and average relative RMSE for the pseudo-EBLUP (22) and the Jiang and Lahiri estimator (23) are shown in Table 4. It is interesting to note that neither estimator appears to perform any better than the standard EBLUP in these design-based simulations, and all three are substantially out performed in terms of average relative RMSE by the two MBD-type estimators that were investigated in this study. Clearly the results of a single (but reasonably realistic) simulation study should not be considered as anything more than indicative. However, they do provide some evidence that asymptotic design-based properties are no guarantee of small area estimation performance.

The indirect predictor (20) of the small area mean is obtained by using well known prediction-based ideas. Under log transformed models, there are alternative approaches to obtain better indirect predictor for small area mean. For example, Slud and Maiti (2006) described an

indirect predictor for the small area mean under an area level version of the log transformed model (9). Berg (2009, private communication) follows the Slud-Maiti approach to obtain a predictor for small area mean under a random intercepts specification of the unit level log transformed model (9). However, like the Slud-Maiti predictor, Berg's predictor ignores the bias correction necessary after back-transformation to the raw scale. The empirical properties of this predictor have yet to be examined.

Acknowledgements

The first author gratefully acknowledges the financial support provided by a PhD scholarship from the U.K. Commonwealth Scholarship Commission. Constructive comments from Editor, Associate Editor and two referees are also gratefully acknowledged. They resulted in the revised version of the article representing a considerable improvement on the original.

References

- Bates, D.M., and Pinheiro, J.-C. (1998). Computational Methods for Multilevel Models. <http://franz.stat.wisc.edu/pub/NLME/>.
- Carroll, R., and Ruppert, D. (1988). *Transformation and Weighting in Regression*. New York: Chapman and Hall.
- Chambers, R., and Tzavidis, N. (2006). M-quantile models for small area estimation. *Biometrika*, 93, 255-268.
- Chandra, H., and Chambers, R.L. (2005). Comparing EBLUP and C-EBLUP for small area estimation. *Statistics in Transition*, 7, 637-648.
- Chandra, H., and Chambers, R. (2009). Multipurpose weighting for small area estimation. *Journal of Official Statistics*, 25(3), 379-395.
- Chandra, H., Salvati, N. and Chambers, R. (2007) Small area estimation for spatially correlated populations. A comparison of direct and indirect model-based methods. *Statistics in Transition*, 8, 887-906.
- Chen, G., and Chen, J. (1996). A transformation method for finite population sampling calibrated with empirical likelihood. *Survey Methodology*, 22, 139-146.
- Harville, D.A. (1977). Maximum likelihood approaches to variance component estimation and to related problems. *Journal of the American Statistical Association*, 72, 320-338.
- Hidiroglou, M.A., and Smith, P.A. (2005). Developing small area estimates for business surveys at the ONS. *Statistics in Transition*, 7, 527-539.
- Jiang, J., and Lahiri, P. (2006). Estimation of finite population domain means: A model-assisted empirical best prediction approach. *Journal of the American Statistical Association*, 101, 301-311.
- Karlberg, F. (2000a). Population total prediction under a lognormal superpopulation model. *Metron*, LVIII, 53-80.
- Karlberg, F. (2000b). Survey estimation for highly skewed populations in the presence of zeroes. *Journal of Official Statistics*, 16, 229-241.
- Longford, N.T. (2007). On standard errors of model-based small-area estimators. *Survey Methodology*, 33, 69-79.
- McCulloch, C.E., and Searle, S.R. (2001). *Generalized, Linear and Mixed Models*. New York: John Wiley & Sons, Inc.
- 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. (2003). *Small Area Estimation*. New York: John Wiley & Sons, Inc.
- Royall, R.M. (1976). The linear least-squares prediction approach to two-stage sampling. *Journal of the American Statistical Association*, 71, 657-664.
- Royall, R.M., and Cumberland, W.G. (1978). Variance estimation in finite population sampling. *Journal of the American Statistical Association*, 73, 351-358.
- Slud, E. V., and Maiti, T. (2006). Mean squared error estimation in transformed Fay-Herriot models. *Journal of the Royal Statistical Society, Series B*, 68(2), 239-257.
- Tzavidis, N., Salvati, N., Pratesi, M. and Chambers, R. (2008). M-quantile models with application to poverty mapping. *Statistical Methods And Applications*, 17, 393-411.
- Valliant, R., Dorfman, A.H. and Royall, R.M. (2000). *Finite Population Sampling and Inference*. New York: John Wiley & Sons, Inc.
- Wu, C., and Sitter, R.R. (2001). A model calibration approach to using complete auxiliary information from survey data. *Journal of the American Statistical Association*, 96, 185-193.