# THE ROLE OF MODELS IN MODEL-ASSISTED AND MODEL-DEPENDENT ESTIMATION FOR DOMAINS AND SMALL AREAS

Risto Lehtonen[1]

[1] University of Helsinki, Finland
e-mail: risto.lehtonen@helsinki.fi

## Abstract

Estimation for population subgroups or domains is investigated for model-assisted generalized regression (GREG) and model-dependent EBLUP estimators, under different model choices and under unequal probability sampling. Two particular issues are addressed: (i) how to account for the domain differences in the model formulation, and (ii) how to account for the underlying unequal probability sampling design. Results on bias and accuracy of GREG and EBLUP are based on Monte Carlo experiments where PPS samples were drawn from an artificially generated population. The bias of GREG estimator remained negligible for all model formulations considered, and accuracy improved when including the PPS size variable in the assisting model. A "double-use" of the auxiliary data both in the sampling design and in the estimation design appeared favorable. In GREG, the mixed model formulation did not outperform the fixed-effects model formulation. For EBLUP, the model choice was critical and if not successful, large bias was introduced. For unweighted EBLUP, substantial bias reduction was attained with the inclusion of the PPS size variable in the model. We propose a new weighted EBLUP estimator for unequal probability sampling designs, as an alternative to the unweighted EBLUP. The results show that the weighted EBLUP behaves better that the unweighted EBLUP, but still the bias can be substantial and can dominate the MSE, which invalidates the construction of proper confidence intervals.

## 1 Introduction

Estimation of reliable statistics for population subgroups or domains constitutes an area of increasing importance in the production of official statistics. A good example is the estimation of the number of unemployed and employed, and the accompanying standard errors, for regional areas in a country by using sample survey data from a Labour Force Survey and auxiliary data taken from the available register and census data sources. Typically, a LFS is planned to produce reliable statistics for the entire population and large or major areas. Standard design-based direct estimators, such as the Horvitz-Thompson estimator, are often used for such cases. The task can become challenging when the number of sample elements in a number of domains remains small or minor. In this case, more advanced methods that effectively use the available auxiliary information are needed.

Methods available for the estimation of totals for domains and small areas include model-assisted design-based estimators, referring to the family of generalized regression (GREG) estimators (Särndal, Swensson and Wretman 1992, Estevao and Särndal 1999, 2004), and model-dependent techniques, such as the EBLUP estimator (Empirical Best Linear Unbiased Predictor) and synthetic estimators (Ghosh 2001, Rao 2003). Properties of these estimator types are discussed for example in Lehtonen and Veijanen (1998, 1999) and Lehtonen, Veijanen and Särndal (2003, 2005). The documentation of the EURAREA project includes use-

ful comparative materials on properties of model-dependent estimators (EURAREA Consortium 2004, Heady and Ralphs 2005).

Known design-based properties related to bias, precision and accuracy of model-assisted estimators and model-dependent estimators are summarized in Table 1. Model-assisted estimators are approximately design-unbiased by definition, but their variance can become large in domains where the sample size is small. Model-dependent estimators are design-biased: the bias can be large for domains where the model does not fit well. The variance of a model-dependent estimator can be small even for small domains, but the accuracy tends to be poor because the squared bias often dominates the mean squared error (MSE), as shown for example by Lehtonen, Veijanen and Särndal (2003 and 2005). The dominance of the bias component together with a small variance can cause poor coverage rates and invalid confidence intervals for a model-dependent estimator. For model-assisted design-based estimators, on the other hand, valid confidence intervals can be constructed. Typically, model-assisted estimators are used for major or not-so-small domains and model-dependent estimators are used for small domains where model-assisted estimators can fail.

**Table 1.** Design-based properties of model-assisted and model-dependent estimators for domains and small areas.

|  | **Design-based model-assisted methods - GREG family** | **Model-dependent methods SYN and EBLUP** |
|---|---|---|
| **Design bias** | Design unbiased (approximately) by the construction principle | Design biased Bias does not necessarily approach zero with increasing domain sample size |
| **Precision (Variance)** | Variance may be large for small domains Variance tends to decrease with increasing domain sample size | Variance can be small even for small domains Variance tends to decrease with increasing domain sample size |
| **Accuracy (Mean Squared Error, MSE)** | MSE = Variance (or nearly so) | MSE = Variance + squared Bias Accuracy can be poor if the bias is substantial |
| **Confidence intervals** | Valid intervals can be constructed | Valid intervals not necessarily obtained |

Survey statistician often faces challenging methodological choices when aiming at reliable estimation of population totals for domains and small areas. These choices include, for example, the inferential framework, model type (mathematical form, specification, parametrization, estimation of model parameters), and estimator type (point estimator, estimator of variance or MSE) for the unknown domain totals. Related to the problem of model choice, or the role of the model in model-assisted estimators and in model-dependent estimators, the two questions of special interest in this study are:

(i) How to account for the domain differences in the model formulation (relevant for model-assisted estimators in particular)?

(ii) How to account for the underlying unequal probability sampling design (relevant for model-dependent estimators in particular)?

We discuss points (i) and (ii) to some extent from a design-based perspective, under the fixed finite population approach. More specifically, we compare the relative performance (bias and accuracy) of the two estimator types of domain totals, GREG, and EBLUP, under different model choices. A continuous response variable is assumed. In the construction of models we

use both linear fixed-effects models and linear mixed models, where random effects are included in addition to the fixed effects. We fit the linear models with different parametrizations. In the estimation of the model parameters, we use both weighted and unweighted estimation procedures.

An underlying unequal probability sampling design is assumed. The case of unequal probability sampling is of importance for practical purposes in official statistics and many fields of empirical research. Without-replacement type fixed-size Probability Proportional to Size sampling (systematic PPS) was selected to represent an example of an unequal probability sampling design. This study extends the case of equal probability sampling investigated in Lehtonen, Särndal and Veijanen (2003, 2005) to unequal probability sampling designs.

The working paper is organized as follows. Chapter 2 introduces our notation and models and estimators used. Results for GREG and EBLUP estimators are given in Chapter 3. Conclusions are in Chapter 4.

## 2  Methods

### 2.1 Models and estimators of domain totals

We are interested in the estimation of totals of a continuous response variable *y* for the domains of interest. Availability of powerful auxiliary information is essential for the estimators of domain totals considered. We assume that we have access to unit-level data, which include domain membership indicators and vectors of auxiliary x-variables, for all units in the population. The auxiliary data vector also contains the size variable used in the PPS sampling procedure. The auxiliary data are incorporated in the estimation procedure by an appropriate model. Thus, the choice of the model that underlies the GREG, SYN and EBLUP estimators of domain totals is considered important.

Our question (i) was "How to account for the domain differences in the model formulation?". The domain differences can be accounted for by a proper model formulation. Basically, there are two options to facilitate the domain differences: (1) introduction of domain-specific fixed effects in the model, and (2) accounting for the domain differences by domain-specific random effects, such as random intercepts. It is obvious that these options are relevant for model-assisted estimators in particular. The reason is that in a standard GREG setting, a fixed-effects linear model is routinely used as the assisting model (Estevao and Särndal 1999, 2004), and a GREG estimator that uses a mixed model, the MGREG estimator, has been introduced only recently (Lehtonen and Veijanen 1999, Lehtonen et al. 2003, see also Goldstein 2003, p. 165). On the other hand, a mixed model formulation has a long tradition in the context of EBLUP estimation of small area totals (Fay and Herriot 1979, Rao 2003). The problem of model choice is discussed in a more general spirit in Firth and Bennett (1998).

To throw some light on question (ii) "How to account for the underlying unequal probability sampling design?", we study the different options to incorporate the information of the sampling design into the estimation procedure. In the modelling phase, there are two main options to account for the sampling design: (a) the incorporation of sampling weights in the estimation of model parameters, and (b) the inclusion of sampling design variables as additional covariates in the model. By default, sampling weights are incorporated in the estimation procedures for all assisting models of GREG estimators. As a rule, sampling weights are ignored in the estimation procedures for SYN estimators.

Typically, the underlying mixed model of a standard EBLUP estimator is fitted in an unweighted manner. Rao (2003) introduced a pseudo EBUP estimator, where sampling weights are included in the construction of the EBLUP estimator, but the parameters of the mixed model are estimated by unweighted techniques. As an alternative to the unweighted EBLUP

and pseudo EBLUP, we will introduce a new EBLUP estimator, where sampling weights are incorporated in the estimation of parameters of the underlying mixed model. We will also compare options (a) and (b) in their successfulness in accounting for the sampling design. It is obvious that these options are relevant for EBLUP estimators in particular.

We study the bias and accuracy properties of the estimators of domain totals by empirical methods. Our Monte Carlo simulation experiments consisted of repeated draws of systematic PPS samples from an artificially constructed fixed finite population.

Table 2 shows the model-dependent and model-assisted estimators to be discussed, in a two-way arrangement by estimator type and by model choice. Each of the rows corresponds to a different model choice. CC model (common intercepts, common slopes) is one whose only parameters are fixed effects defined at the population level; it contains no domain specific parameters. We obtain SYN-CC and GREG-CC estimators. SC model (separate intercepts, common slopes) is one having at least some of its parameters or effects defined at the domain level. These are fixed effects for SYN-SC and GREG-SC and random effects for EBLUP-SC, EBLUPW-SC and MGREG-SC. Table 2 also shows the estimation methods that are used in the estimation of model parameters.

To address points (i) and (ii) of Chapter 1, we discuss in more detail GREG-SC and MGREG-SC for GREG family estimators and EBLUP-SC and EBLUPW-SC for EBLUP family estimators.

**Table 2.** Schematic presentation of the model-dependent and model-assisted estimators of domain totals for a continuous response variable by model choice and estimator type, under unequal probability sampling.

| Model choice | | | | Estimator type | |
|---|---|---|---|---|---|
| **Model abbreviation** | **Model specification** | **Effect type** | **Estimation of model parameters** | *Model-dependent estimators* | *Model-assisted estimators* |
| CC | Common intercepts Common slopes | Fixed effects | OLS | SYN-CC | Not applicable(**) |
| | | | WLS | Not applicable(*) | GREG-CC |
| SC | Separate intercepts Common slopes | Fixed effects | OLS | SYN-SC | Not applicable (**) |
| | | | WLS | Not applicable(*) | GREG-SC |
| | | Fixed and random | REML GLS | EBLUP-SC | Not applicable (**) |
| | | | Weighted REML GWLS | EBLUPW-SC | MGREG-SC |

OLS Ordinary least squares
WLS Weighted least squares (sampling weights)
GLS Generalized least squares
GWLS Generalized weighted least squares (sampling weights)
REML Restricted (residual) maximum likelihood
Weighted REML Restricted pseudo maximum likelihood (sampling weights)

(*) In SYN, weights are ignored in the estimation procedure by default.
(**) In GREG, weights are incorporated in the estimation procedure by default.

We next introduce the notation used in this study.

**Population and sampling design**

$U = \{1, 2, ..., k, ..., N\}$   Population (fixed, finite)

$U_1, ..., U_d, ..., U_D$       Domains of interest (non-overlapping)

$Y_d = \sum_{U_d} y_k$, $d = 1, ..., D$   Target parameters (domain totals)

$\mathbf{x}_k = (x_{1k}, ..., x_{pk})'$   Auxiliary variable vector

$I_{dk} = 1$ if $k \in U_d$   Domain membership indicators,

$I_{dk} = 0$ otherwise   $d = 1, ..., D$

Note that we assume the vector value $\mathbf{x}_k$ and domain membership to be known for every population unit $k \in U$.

Sampling design: Systematic PPS with sample size $n$

$s$                 Sample from $U$

$s_d = s \cap U_d$   Random part of $s$ falling in domain $d$

$\pi_k = n \dfrac{x_{1k}}{\sum_{k \in U} x_{1k}}$   Inclusion probability for $k \in U$

$a_k = 1/\pi_k$   Sampling weight for $k \in s$

We observe $y_k$ for $k \in s$. Note that for estimation purposes, sample data and auxiliary data are merged at the micro level by using unique ID keys that are available in both data sources.

**Models for continuous response $y$**

Linear fixed-effects models

   CC models $y_k = \beta_0 + \beta_1 x_{1k} + ... + \beta_p x_{pk} + \varepsilon_k$

   SC models $y_k = \beta_{0d} + \beta_1 x_{1k} + ... + \beta_p x_{pk} + \varepsilon_k$, $d = 1, ..., D$

   Fitted values under fixed-effects models $\hat{y}_k = \mathbf{x}_k' \hat{\boldsymbol{\beta}}$

Linear mixed models

   SC models $y_k = \beta_0 + u_d + \beta_1 x_{1k} + ... + \beta_p x_{pk} + \varepsilon_k$, $d = 1, ..., D$

   where $u_d$ are domain-specific random intercepts

   Fitted values under mixed models $\hat{y}_k = \mathbf{x}_k' \hat{\boldsymbol{\beta}} + \hat{u}_d$, $d = 1, ..., D$

Note that fitted values $\hat{y}_k$ are calculated for every $k \in U$.

**Estimators of domain totals**

The predictions $\{\hat{y}_k ; k \in U\}$ differ from one model specification to another. For a given model specification, the estimator of the domain total $Y_d = \sum_{U_d} y_k$ has the following structure for the three estimator types (SYN, GREG, EBLUP):

Model-assisted GREG estimators

$$\hat{Y}_{dGREG} = \sum_{k \in U_d} \hat{y}_k + \sum_{k \in s_d} a_k(y_k - \hat{y}_k)$$

Model-dependent SYN estimators

$$\hat{Y}_{dSYN} = \sum_{k \in U_d} \hat{y}_k$$

Model-dependent EBLUP estimators

$$\hat{Y}_{dEBLUP} = \sum_{k \in s_d} y_k + \sum_{k \in U_d - s_d} \hat{y}_k$$

where $d = 1,...,D.$

Note that $\hat{Y}_{dSYN}$ and $\hat{Y}_{dEBLUP}$ rely heavily on the truth of the model, and can be biased if the model is misspecified. On the other hand, $\hat{Y}_{dGREG}$ has a second term that protects against model misspecification.

We adopt the following conventions (Table 2). In SYN-CC, SYN-SC, GREG-CC and GREG-SC, a fixed-effects model formulation is assumed. A mixed model is assigned for EBLUP-SC, EBLUPW-SC and MGREG-SC estimators.

## Measures used in Monte Carlo simulations

In Monte Carlo simulation experiments, by using estimates $\hat{Y}_d(s_v)$ from repeated samples $s_v; v = 1, 2,..., K$, we computed for each domain $d = 1,...,D$ the following Monte Carlo summary measures of bias and accuracy.

(i) Absolute relative bias (ARB), defined as the ratio of the absolute value of bias to the true value:

$$\left| \frac{1}{K} \sum_{v=1}^{K} \hat{Y}_d(s_v) - Y_d \right| / Y_d$$

(ii) Relative root mean squared error (RRMSE), defined as the ratio of the root MSE to the true value:

$$\sqrt{\frac{1}{K} \sum_{v=1}^{K} (\hat{Y}_d(s_v) - Y_d)^2} / Y_d$$

## Details of the simulations

There were 100 domains in the population. The size of domain $d$ was proportional to $\exp(q_d)$, where $q_d$ was simulated from U(0,2.9). Each observation was allocated to a domain by geometric probability: intervals of length $\exp(q_d)$ were concatenated and a random point was chosen in (0, $\sum_d \exp(q_d)$). The interval containing the point determined the domain of the observation.

There were 47 minor domains, 19 medium-sized domains and 34 major domains in the population. These three classes were defined on the basis of expected sample size $n(N_d / N)$: less than 70, 70-119 and 120 or more units, respectively. The smallest domain of the generated population had 1,711 units and the largest had 28,296.

The variable $x_1$ is the size variable used in PPS sampling. The variable was simulated from uniform distribution U(1,11). Another auxiliary variable $x_2$ was simulated from N(0,9). The random effects $u_d$ were simulated independently from N(0,0.25). The error term $\varepsilon$ followed N(0,1).

Responses were simulated as

$$y_k = 1 + 2x_{1k} + 1.5x_{2k} + u_d + \varepsilon_k \ (k \in d)$$

Correlations of the variables in the population were: $corr(y, x_1) = 0.779$, $corr(y, x_2) = 0.607$ and $corr(x_1, x_2) = -0.001$. Domain means of the response variable were approximately equal, but the totals differed considerably: The means of domain totals were 45,614 for minor domains, 117,308 for medium domains and 241,527 for major domains.

Our population size is $N = 1,000,000$ and sample size $n = 10,000$. In Monte Carlo experiments, $K = 1000$ independent systematic PPS samples were generated. The inclusion probabilities are $\pi_k = nx_{1k} / \sum_k x_{1k}$. The weights $a_k = 1/\pi_k$ varied between 54.6 and 596.5.

## 3 Results

### 3.1 GREG estimators

We first discuss results for GREG estimators. Our point (i) devoted to GREG was "How to account for the domain differences in the model formulation". This is demonstrated by the eight different model formulations in Table 3. In models A1, B1, C1 and D1, the domain differences are accounted for by domain-specific fixed effects $\beta_{0d}$. In models A2, B2, C2 and D2, we use random intercepts $\beta_0 + u_d$, where $\beta_0$ is the fixed intercept common for all domains, and the random term $u_d$ is domain-specific. In addition, we have two explanatory variables at our disposal: the variable $x_1$, which was used in the PPS sampling design, and $x_2$, which is an auxiliary variable uncorrelated to $x_1$. Note that both variables correlate quite strongly with the response variable $y$. For $x_1$ and $x_2$, slope parameters $\beta_1$ and $\beta_2$ are common fixed effects for all domains.

For GREG, we incorporate the sampling weights in the estimation procedure of model parameters, including the mixed model underlying the MGREG-SC estimator. This facilitates the condition of "internal bias calibration" (a proper combination of model formulation and estimation procedure under a given sampling design) proposed by Firth and Bennett (1998).

Table 3 also shows our model building strategy. We start with simple models A1 and A2 and proceed step by step towards the population generating model D2. In all models considered, GREG family estimators are essentially unbiased, and a fixed-effects model formulation and a mixed model formulation yield similar accuracy. An explanation for this observation is that in the setting of this exercise, the average levels of the response variable did not vary much over the domains. Best accuracy (excluding the true model) is for models where the PPS size variable $x_1$ is included. This demonstrates the accuracy gains attained from the "double-use" of $x_1$ both in the sampling design and in the estimation design; see also Särndal (1996). We also note that accuracy differences between the different GREG estimators are substantial especially in minor and medium domains, and accuracy improves with increasing the domain sample size.

**Table 3.** Average absolute relative bias ARB (%) and average relative root mean squared error RRMSE (%) of model-assisted GREG estimators of domain totals for minor, medium-sized and major domains of the generated population.

| Model and estimator | Average ARB (%) | | | Average RRMSE (%) | | |
|---|---|---|---|---|---|---|
| | Domain size class | | | Domain size class | | |
| | Minor (20-69) | Medium (70-119) | Major (120+) | Minor (20-69) | Medium (70-119) | Major (120+) |
| **Model A1** $y_k = \beta_{0d} + \varepsilon_k$ | | | | | | |
| GREG-SC | 1.4 | 0.5 | 0.3 | 13.7 | 8.1 | 5.7 |
| **Model A2** $y_k = \beta_0 + u_d + \varepsilon_k$ | | | | | | |
| MGREG-SC | 0.2 | 0.2 | 0.1 | 13.7 | 8.1 | 5.6 |
| **Model B1** $y_k = \beta_{0d} + \beta_1 x_{1k} + \varepsilon_k$ | | | | | | |
| GREG-SC | 0.2 | 0.1 | 0.0 | 7.8 | 4.6 | 3.2 |
| **Model B2** $y_k = \beta_0 + u_d + \beta_1 x_{1k} + \varepsilon_k$ | | | | | | |
| MGREG-SC | 0.2 | 0.1 | 0.0 | 7.8 | 4.6 | 3.3 |
| **Model C1** $y_k = \beta_{0d} + \beta_2 x_{2k} + \varepsilon_k$ | | | | | | |
| GREG-SC | 1.4 | 0.5 | 0.3 | 11.6 | 6.8 | 4.8 |
| **Model C2** $y_k = \beta_0 + u_d + \beta_2 x_{2k} + \varepsilon_k$ | | | | | | |
| MGREG-SC | 0.2 | 0.1 | 0.1 | 11.6 | 6.8 | 4.7 |
| **Model D1** $y_k = \beta_{0d} + \beta_1 x_{1k} + \beta_2 x_{2k} + \varepsilon_k$ | | | | | | |
| GREG-SC | 0.0 | 0.0 | 0.0 | 1.7 | 1.0 | 0.7 |
| **Model D2** $y_k = \beta_0 + u_d + \beta_1 x_{1k} + \beta_2 x_{2k} + \varepsilon_k$ (Population generating model) | | | | | | |
| MGREG-SC | 0.0 | 0.0 | 0.0 | 1.7 | 1.0 | 0.7 |
| Variables | $x_1$ Size variable in PPS sampling, $x_2$ Auxiliary variable | | | | | |

## 3.2 EBLUP estimators

For estimators of the EBLUP family, we asked "How to account for the underlying unequal probability sampling design?". We proposed two options for this purpose: (a) the incorporation of sampling weights in the estimation of model parameters, and (b) the inclusion of sampling design variables as additional covariates in the model.

We compare unweighted and weighted EBLUP estimators constructed with four mixed model formulations. Model A includes a random intercept, variable $x_1$ is included in Model B, variable $x_2$ is included in Model C and both variables appear in the population generating model D. Similarly as for GREG, domain differences are accounted for by random intercept terms, and slope parameters are common for all domains. For all models (except D), EBLUP estimators are calculated with unweighted and weighted estimation of model parameters.

For Models A and C, unweighted estimators EBLUP-SC are seriously biased. For these models, the PPS sampling design is not accounted for. The bias declines considerably when the sampling weights are incorporated in the estimation of the mixed model, as shown by the new EBLUPW-SC estimators for Models A and C. The unweighted estimator EBLUP-SC under Model B shows best bias behaviour, indicating that the inclusion of the PPS size variable in the model can offer a powerful tool for bias reduction for EBLUP family estimators. Use of both weighting and the inclusion of $x_1$ in the model appears to be less powerful.

Accuracy behaviour of all EBLUP estimators is infected by the dominance of the squared bias component in the MSE, as indicated by the RRMSE figures. This holds for all three do-

main size classes. Because of large bias and small variance, invalid confidence intervals can be obtained. This means that point estimates can be systematically far away from the true value, independently of the domain sample size. In addition, accuracy does not improve much with increasing the domain sample size.

**Table 4.** Average absolute relative bias ARB (%) and average relative root mean squared error RRMSE (%) of model-dependent EBLUP estimators of domain totals for minor, medium-sized and major domains of the generated population.

| Model and estimator | Average ARB (%) | | | Average RRMSE (%) | | |
|---|---|---|---|---|---|---|
| | Domain size class | | | Domain size class | | |
| | Minor (20-69) | Medium (70-119) | Major (120+) | Minor (20-69) | Medium (70-119) | Major (120+) |
| **Model A** $y_k = \beta_0 + u_d + \varepsilon_k$ | | | | | | |
| EBLUP-SC | 22.9 | 23.1 | 21.7 | 22.9 | 23.3 | 21.8 |
| EBLUPW-SC | 3.7 | 3.5 | 3.3 | 3.9 | 3.6 | 3.5 |
| **Model B** $y_k = \beta_0 + u_d + \beta_1 x_{1k} + \varepsilon_k$ | | | | | | |
| EBLUP-SC | 1.8 | 1.4 | 0.7 | 2.8 | 2.5 | 2.2 |
| EBLUPW-SC | 3.5 | 3.5 | 3.3 | 3.5 | 3.6 | 3.3 |
| **Model C** $y_k = \beta_0 + u_d + \beta_2 x_{2k} + \varepsilon_k$ | | | | | | |
| EBLUP-SC | 22.3 | 23.1 | 21.8 | 22.4 | 23.2 | 21.9 |
| EBLUPW-SC | 3.7 | 3.6 | 3.2 | 3.9 | 3.7 | 3.3 |
| **Model D** $y_k = \beta_0 + u_d + \beta_1 x_{1k} + \beta_2 x_{2k} + \varepsilon_k$ (Population generating model) | | | | | | |
| EBLUP-SC | 0.3 | 0.1 | 0.0 | 1.3 | 0.8 | 0.6 |
| Variables | $x_1$ Size variable in PPS sampling, $x_2$ Auxiliary variable | | | | | |

# 4 Conclusions

Results indicate that under unequal probability sampling, model-assisted GREG family estimators are quite insensitive to the model choice, a property also shown in our previous research to hold under SRSWOR. Model formulation and the estimation strategy of the model are critical for model-dependent EBLUP family estimators. This is especially true when using EBLUP for unequal sampling designs.

Bias of GREG estimators remained negligible for all model choices. "Double-use" of the same auxiliary information, that is, the use of the size variable in the PPS sampling design and in the assisting model, appeared to be beneficial with respect to accuracy. The accuracy improved with increasing the domain sample size. In this case, the mixed model formulation did not outperform the fixed-effects model formulation.

For model-dependent EBLUP family estimators, the bias can be large for a misspecified model. The PPS sampling design could be accounted for with two options, by the inclusion of the PPS size variable in the mixed model, or by the use of the weighted version of the EBLUP estimator, where the sampling weights are incorporated in the estimation procedure of model parameters. Of these two options, the first one appeared to be more effective, producing an EBLUP estimator with small bias and good accuracy. However, for both options, the squared bias component can still dominate the MSE, even in minor domains, tending to invalidate the construction of proper confidence intervals. Dominance of the bias component also can cause that the accuracy does not show improvement, when increasing the domain sample size.

# References

Estevao, V.M. and Särndal, C.-E. (1999) The use of auxiliary information in design-based estimation for domains. *Survey Methodology*, **25**, 213–221.

Estevao, V.M. and Särndal, C.-E. (2004) Borrowing strength is not the best technique within a wide class of design-consistent domain estimators. *Journal of Official Statistics*, **20**, 645–669.

EURAREA Consortium (2004) Project Reference Volume, Parts 1, 2 and 3 (PDF).
Website: www.statistics.gov.uk/eurarea/

Fay, R.E. and Herriot , R.A. (1979) Estimation of income for small places: an application of James-Stein procedures to census data. *Journal of the American Statistical Association*, 74, 269–277.

Firth, D. and Bennett, K.E. (1998) Robust models in probability sampling. (With discussion) *Journal of the Royal Statistical Society*, B, **60**, 3–56.

Ghosh, M. (2001) Model-dependent small area estimation: theory and practice. In: Lehtonen and Djerf K. (eds) *Lecture Notes on Estimation for Population Domains and Small Areas*. Helsinki: Statistics Finland, Reviews 2001/5, 51–108.

Goldstein, H. (2003) *Multilevel Statistical Models*. Third Edition. London: Arnold.

Heady, P. and Ralphs, M. (2005) EURAREA: an overview of the project and its findings. *Statistics in Transition,* 7, 557–570.

Lehtonen, R., Särndal, C.-E. and Veijanen, A. (2003) The effect of model choice in estimation for domains, including small domains. *Survey Methodology*, **29**, 33–44.

Lehtonen, R., Särndal, C.-E. and Veijanen, A. (2005) Does the model matter? Comparing model-assisted and model-dependent estimators of class frequencies for domains. *Statistics in Transition,* **7**, 649–673.

Lehtonen, R. and Veijanen, A. (1998). Logistic generalized regression estimators. *Survey Methodology*, 24, 51–55.

Lehtonen, R. and Veijanen, A. (1999) Domain estimation with logistic generalized regression and related estimators. *Proceedings, IASS Satellite Conference on Small Area Estimation*, Riga, August 1999. Riga: Latvian Council of Science, 121–128.

Rao, J.N.K. (2003) *Small Area Estimation*. Hoboken: Wiley.

Särndal, C.-E. (1996)  Efficient estimators with simple variance in unequal probability sampling. *Journal of the American Statistical Association*, **91**, 1286–1300.

Särndal, C.-E., Swensson, B. and Wretman, J.H. (1992) *Model Assisted Survey Sampling*. New York: Springer-Verlag.