



# Multivariate Fay–Herriot models for small area estimation



Roberto Benavent, Domingo Morales\*

Centro de Investigación Operativa, Universidad Miguel Hernández de Elche, Avda. de la Universidad s/n, 03202 Elche, Spain

## ARTICLE INFO

### Article history:

Received 14 January 2015

Received in revised form 23 July 2015

Accepted 24 July 2015

Available online 1 August 2015

### Keywords:

Multivariate linear mixed models

Fay–Herriot model

REML method

EBLUP

MSE estimation

Bootstrap

Poverty

## ABSTRACT

Multivariate Fay–Herriot models for estimating small area indicators are introduced. Among the available procedures for fitting linear mixed models, the residual maximum likelihood (REML) is employed. The empirical best predictor (EBLUP) of the vector of area means is derived. An approximation to the matrix of mean squared crossed prediction errors (MSE) is given and four MSE estimators are proposed. The first MSE estimator is a plug-in version of the MSE approximation. The remaining MSE estimators combine parametric bootstrap with the analytic terms of the MSE approximation. Several simulation experiments are performed in order to assess the behavior of the multivariate EBLUP and for comparing the MSE estimators. The developed methodology and software are applied to data from the 2005 and 2006 Spanish living condition surveys. The target of the application is the estimation of poverty proportions and gaps at province level.

© 2015 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 1. Introduction

Surveys are designed for obtaining reliable estimates in the whole population or in some subpopulations called planned domains. However, it is quite common in practice to use survey data for estimating indicators of non-planned domains (small areas) with small samples sizes. Small area estimation deals with inference problems for this kind of domains. In these cases, direct estimators might have large sampling errors. Direct estimators can be improved by assuming regression models that link all the sample data by introducing a relation between the variable of interest and a set of explanatory variables.

Linear mixed models use random area effects for the extra between-area variation of the data that is not explained by the auxiliary variables. Often auxiliary individual information is not available, but data aggregated to the small areas can be found in administrative registers. Then the model can be stated at the small area level. An area-level linear mixed model with random area effects was first proposed by Fay and Herriot (1979) to estimate average per-capita income in small places of the US. Since then, the Fay–Herriot model has been one of the most widely used models in small area estimation.

In recent years, many researchers have investigated applications of the Fay–Herriot model to small area estimation problems. Without being exhaustive, we cite some papers dealing with the Fay–Herriot model. Prasad and Rao (1990); Datta and Lahiri (2000); Das et al. (2004); González-Manteiga et al. (2010); Jiang et al. (2011); Datta et al. (2011a) and Kubokawa (2011) gave tools for measuring the uncertainty of model-based small area estimators. Datta et al. (2011b), Bell et al. (2013) and Pfeiffermann et al. (2014) studied the problem of benchmarking. Ybarra and Lohr (2008) proposed a new small area estimator that accounts for sampling variability in the auxiliary information. Herrador et al. (2011) treated situations where small areas are divided into two groups and domain random effects have different variances across the groups. Slud and Maiti (2011) were interested on small area estimation based on left censored survey data.

Statisticians are often required to estimate correlated descriptive measures, like poverty or unemployment indicators. Multivariate models take into account for the correlation of several variables and typically fit to this kind of situations.

\* Corresponding author.

E-mail address: [d.morales@umh.es](mailto:d.morales@umh.es) (D. Morales).

Some papers can be found in the literature of small area estimation where multivariate linear mixed models are employed. Fay (1987) and Datta et al. (1991) compared the precision of small area estimators obtained from univariate models for each response variable with the ones obtained by a multivariate model. Datta et al. (1996) used also a multivariate Fay–Herriot model for obtaining hierarchical Bayes estimates of median income of four-person families for the US states. González-Manteiga et al. (2008b) studied a class of multivariate Fay–Herriot model with a common random effect for all the components of the target vector. They further introduced bootstrap approximations to prediction errors. This paper introduces a class of multivariate Fay–Herriot models with one random effect per component of the target vector and allowing for different covariance patterns between the components of the vector of random effects. This is a new and flexible class of multivariate models that does not contain the models of González-Manteiga et al. (2008b) as particular cases.

Historical data give relevant information that can be used to obtain better small area estimators. Several authors have proposed extensions of the Fay–Herriot model that borrow strength from time. Choudry and Rao (1989) introduced a model including several time instants and considering an autocorrelated structure for sampling errors. Rao and Yu (1994) proposed a model that borrows information across areas and over time. Ghosh et al. (1996) proposed a time correlated area level model to estimate the median income of four-person families for American states. Datta et al. (1999), You and Rao (2000), Datta et al. (2002), Esteban et al. (2011, 2012), Marhuenda et al. (2013) and Morales et al. (2015) gave some extensions of the Rao–Yu model with applications to the estimation of labor or poverty indicators. Singh et al. (2005) and Pfeiffermann and Burck (1990) considered models with time-varying random slopes obeying an autoregressive process. This paper applies multivariate Fay–Herriot models to time correlated data. In this setup, the introduced multivariate models contain the models proposed by Esteban et al. (2011) as particular cases.

The paper is organized as follows. Section 2 introduces the multivariate Fay–Herriot model and gives a residual maximum likelihood (REML) fitting algorithm. Unlike the model introduced by González-Manteiga et al. (2008b) with a common random effect for all the components of the target variable, the new models have multivariate vectors of random effects with the same dimension as the target variable and allowing for different correlation structures. Section 3 approximates the matrix of mean squared crossed prediction errors (MSE) of the multivariate empirical best predictor (EBLUP) and gives some estimators. The first MSE estimator is a plug-in derivation of the MSE approximation. The remaining MSE estimators combine parametric bootstrap with analytic terms appearing in the MSE approximation. Section 4 presents three simulation experiments. The first simulation studies the behavior of the multivariate EBLUPs under different correlation structures. The second simulation compares the performances of the MSE estimators proposed in Section 3. The third simulation studies the robustness of the EBLUPs against departures from normality. Section 5 applies the developed methodology to data from the Spanish Living Conditions surveys of 2005 and 2006. Two applications are presented. The target of the first application is the estimation of 2006 poverty proportions and gaps. The second application jointly estimates 2005 and 2006 poverty proportions. Section 6 gives some concluding remarks. The Appendix contains detailed proofs of main results.

## 2. Multivariate Fay–Herriot models

Let  $U$  be a finite population partitioned into  $D$  domains  $U_1, \dots, U_D$ . Let  $\mu_d = (\mu_{d1}, \dots, \mu_{dR})'$  be a vector of characteristics of interest in the domain  $d$  and let  $y_d = (y_{d1}, \dots, y_{dR})'$  be a vector of direct estimators of  $\mu_d$ . The multivariate Fay–Herriot model is defined in two stages. The sampling model is

$$y_d = \mu_d + e_d, \quad d = 1, \dots, D, \quad (1)$$

where the vectors  $e_d \sim N(0, V_{ed})$  are independent and the  $R \times R$  covariance matrices  $V_{ed}$  are known. Moreover, it is assumed that the  $\mu_{dr}$ 's are linearly related to  $p_r$  explanatory variables associated to the  $r$ th characteristic in the domain  $d$ . Let  $x_{dr} = (x_{dr1}, \dots, x_{drp_r})'$  be a row vector containing the  $p_r$  explanatory variables for  $\mu_{dr}$  and let  $x_d = \text{diag}(x_{d1}, \dots, x_{dR})_{R \times p}$  with  $p = \sum_{r=1}^R p_r$ . Let  $\beta_r$  be a column vector of size  $p_r$  containing the regression parameters for  $\mu_{dr}$  and let  $\beta = (\beta'_1, \dots, \beta'_R)'_{p \times 1}$ . González-Manteiga et al. (2008b) considered the linking model

$$\mu_d = x_d \beta + 1_R v_d, \quad v_d \stackrel{\text{ind}}{\sim} N(0, \sigma_v^2), \quad d = 1, \dots, D, \quad (2)$$

where  $1_n$  is the  $n \times 1$  vector with all elements equal to 1. This paper introduces multivariate Fay–Herriot models by assuming (1) and substituting the condition (2) by the more realistic linking model

$$\mu_d = x_d \beta + u_d, \quad u_d \stackrel{\text{ind}}{\sim} N(0, V_{ud}), \quad d = 1, \dots, D, \quad (3)$$

where the vectors  $u_d$ 's are independent of the vectors  $e_d$ 's. The  $R \times R$  covariance matrices  $V_{ud}$  depend on  $m$  unknown parameters,  $\theta_1, \dots, \theta_m$ , with  $1 \leq m \leq \frac{R(R-1)}{2} + R$ . Let  $I_n$  be the  $n \times n$  identity matrix,  $\delta_{\ell d}$  be the Kronecker delta,  $y = (y_1, \dots, y_D)'$  be the vector of response variables and define

$$\begin{aligned} u &= \text{col}_{1 \leq d \leq D}(u_d), & e &= \text{col}_{1 \leq d \leq D}(e_d), & u_d &= \text{col}_{1 \leq r \leq R}(u_{dr}), & e_d &= \text{col}_{1 \leq r \leq R}(e_{dr}), \\ X &= \text{col}_{1 \leq d \leq D}(x_d), & Z_d &= \text{col}_{1 \leq \ell \leq D}(\delta_{\ell d} I_R), & Z &= \text{col}'_{1 \leq d \leq D}(Z_d) = I_{DR}, & V_u &= \text{diag}_{1 \leq d \leq D}(V_{ud}), \end{aligned}$$

where  $\text{col}$  and  $\text{col}'$  are matrix operators stacking by columns and rows respectively.

Download English Version:

<https://daneshyari.com/en/article/6869467>

Download Persian Version:

<https://daneshyari.com/article/6869467>

[Daneshyari.com](https://daneshyari.com)