



Small area estimation with spatio-temporal Fay–Herriot models

Yolanda Marhuenda^{a,*}, Isabel Molina^b, Domingo Morales^a

^a Centro de Investigación Operativa, Universidad Miguel Hernández de Elche, Spain

^b Departamento de Estadística, Universidad Carlos III de Madrid, Spain

ARTICLE INFO

Article history:

Received 17 November 2011

Received in revised form 13 July 2012

Accepted 4 September 2012

Available online 10 September 2012

Keywords:

Empirical best linear unbiased estimator

Poverty estimation

Small area estimation

Spatio-temporal model

ABSTRACT

Small area estimation is studied under a spatio-temporal Fay–Herriot model. Model fitting based on restricted maximum likelihood is described and empirical best linear unbiased predictors are derived under the model. A parametric bootstrap procedure is proposed for the estimation of the mean squared error of the small area estimators. The spatio-temporal model is compared with simpler models through simulation experiments, analyzing the gain in efficiency achieved by the use of the more complex model. The performance of the parametric bootstrap estimator of the mean squared error is also assessed. An application with Spanish EU-SILC data is carried out to obtain estimates of poverty indicators for Spanish provinces in 2008, making use of survey data from years 2004–2008.

© 2012 Elsevier B.V. All rights reserved.

1. Introduction

Many times estimators are demanded for domains or subgroups of the population for which the survey providing the data was not planned for. Then, for the domains with small sample size, direct estimators, obtained only using the sample data from the target domain, might not be reliable. Those domains are considered as “small domains” or “small areas”. This fact has led to the development of indirect estimators that make use of the information from related areas, by assuming either implicit or explicit models. These models link all the domains to enhance the estimation at a particular area, that is, they allow to borrow strength from other areas. Estimators based on explicit models, called typically model based estimators, are generally more precise than direct estimators, especially in the areas with smaller sample sizes.

When auxiliary variables related to our study variable are available at the area level (that is, population aggregates are available), a model that has been widely used in the literature of small area estimation and in numerous applications is the basic Fay–Herriot model. This model was first proposed by Fay and Herriot (1979) to obtain small area estimates of mean per capita income in US small areas using survey data. The model is composed of two-stages. In the first stage a model, called *sampling* model, is used to represent the sampling error of direct estimators. Let μ_d be the characteristic of interest in d -th area (typically the mean) and y_d be a direct estimator of μ_d . The sampling model indicates that direct estimators $\{y_d\}$ are unbiased and can be expressed as

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

where D is the total number of areas or domains target of inference. Here, $\{e_d\}$ are sampling errors, which, given μ_d , are independent and normally distributed with known variances, that is $e_d|\mu_d \sim N(0, \sigma_d^2)$, where σ_d^2 is the (assumed known)

* Correspondence to: Centro de Investigación Operativa, Universidad Miguel Hernández de Elche, Avda. de la Universidad s/n, 03202 Elche, Alicante, Spain. Tel.: +34 966658536; fax: +34 966658721.

E-mail addresses: y.marhuenda@umh.es (Y. Marhuenda), isabel.molina@uc3m.es (I. Molina), d.morales@umh.es (D. Morales).

design-based variance of direct estimator y_d , $d = 1, \dots, D$. In the second stage, the true area characteristics $\{\mu_d\}$ are assumed to vary linearly with a number p of area level auxiliary variables, that is,

$$\mu_d = \mathbf{x}'_d \boldsymbol{\beta} + u_d, \quad d = 1, \dots, D,$$

where \mathbf{x}_d is a (column) vector containing the aggregated (population) values of p auxiliary variables for area d , $\boldsymbol{\beta}$ is the vector of regression coefficients and $\{u_d\}$ are model errors, typically assumed to be i.i.d. from $N(0, \sigma_u^2)$ with variance σ_u^2 unknown and independent of $\{e_d\}$. Note that this model, called the *linking* model, links the target quantities μ_d of all the areas through the common regression parameter $\boldsymbol{\beta}$. The Fay–Herriot model can be expressed as a single model in the form

$$y_d = \mathbf{x}'_d \boldsymbol{\beta} + u_d + e_d, \quad d = 1, \dots, D.$$

Many different extensions of this model have been proposed in the literature. For example, a multivariate generalization was studied by González-Manteiga et al. (2008). When available, historical data offer precious information that can be used to improve the estimators at the current instant, that is, it is also possible to borrow strength from time. In this sense, Choudry and Rao (1989) extended the basic Fay–Herriot model including several time instants and considering an autocorrelated structure for sampling errors. More concretely, they considered the model

$$y_{dt} = \mathbf{x}'_{dt} \boldsymbol{\beta} + u_d + e_{dt}, \quad d = 1, \dots, D, \quad t = 1, \dots, T,$$

where here, y_{dt} and \mathbf{x}_{dt} are respectively the response and the vector of auxiliary variables for area d at time instant t , with $\mu_{dt} = \mathbf{x}'_{dt} \boldsymbol{\beta} + u_d$ being the target characteristic for the same area and time instant. For each domain d , the errors $\{e_{dt}\}_{t=1}^T$ were assumed to follow an autoregressive process of order 1, AR(1), that is,

$$e_{dt} = \rho e_{d,t-1} + \epsilon_{dt}, \quad |\rho| < 1, \quad \epsilon_{dt} \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma_\epsilon^2).$$

This model is not allowing for time variation in the area characteristics $\{\mu_{dt}\}$ that is not explained by auxiliary variables.

Another simple model that borrows information across areas and over time, and which includes unexplained area-time variation, was proposed by Rao and Yu (1994), and is given by

$$y_{dt} = \mathbf{x}'_{dt} \boldsymbol{\beta} + u_{1d} + u_{2dt} + e_{dt}, \quad d = 1, \dots, D, \quad t = 1, \dots, T, \quad (1)$$

with $\mu_{dt} = \mathbf{x}'_{dt} \boldsymbol{\beta} + u_{1d} + u_{2dt}$ being now the target characteristic for area d and time instant t . In this model, the area effects $\{u_{1d}\}$ are constant over time following the usual assumptions in the basic Fay–Herriot model and, for each area d , $\{u_{2dt}\}_{t=1}^T$ are time-varying effects that follow an AR(1), but they are independent across areas. Sampling errors $\{e_{dt}\}_{t=1}^T$ are also independent across areas and normally distributed with zero mean vector and general covariance matrix Σ_d , assumed to be known. Sampling errors $\{e_{dt}\}$ are also independent of area and area-time random effects, $\{u_d\}$ and $\{u_{dt}\}$. Rao and Yu (1994) estimated the variance parameters by a method of moments and, for known autocorrelation parameter ρ , they gave a second order approximation to the mean squared error (MSE) of the empirical best liner unbiased predictor (EBLUP) obtained from that model. They also proposed several alternative estimators of ρ together with corresponding MSE estimators. In simulation studies, they report significant gains in efficiency of the EBLUP based on the temporal model when the between-time variation relative to sampling variation was small and area variation was large. Their results indicated that the efficiency was growing with the number of available time instants T .

Other models with temporal correlation have been proposed. Ghosh et al. (1996) proposed a slightly more complicated time correlated area level model to estimate the median income of four-person families for the fifty American states and the district of Columbia. You and Rao (2000) and Datta et al. (2002) used the Rao–Yu model (1), but replacing the AR(1) process by a random walk. Datta et al. (1999) considered a similar model but added extra terms to the linking models to reflect seasonal variation in their application. They applied their model to estimate monthly unemployment rates for nine American states and the district of Columbia. You et al. (2001) applied the Rao–Yu model to estimate monthly unemployment rates for census metropolitan areas in Canada. Finally, Pfeiffermann and Burck (1990) and Singh et al. (1991) considered a model with time-varying random slopes obeying an autoregressive process.

Apart from making use of historical data and including the temporal correlation in the model, it is well known that when there is unexplained spatial correlation in the data, not considering it in our model will lead to erroneous inferences (Cressie, 1993). In small area estimation, when areas are properly delimited regions of the population, closer areas tend to have more similar socio-economic characteristics. This was taken into account in the basic Fay–Herriot model by Singh et al. (2005), Petrucci and Salvati (2006) and Pratesi and Salvati (2008), who considered an extension of the Fay–Herriot model by assuming that area effects $\{u_d\}$ follow a simultaneously autoregressive process of order 1 or SAR(1). When data from neighboring areas are correlated, considering this kind of spatial correlation in the model leads to more efficient small area estimators, see Molina et al. (2009). Bayesian spatial models have been considered by Moura and Migon (2002) and You and Zhou (2011). Thus, taking into account the spatial correlation among data from different areas allows to borrow even more strength from the areas.

Here we consider a model that incorporates historical data similarly as in the Rao–Yu model, and at the same time includes spatial correlation among data from neighboring areas. We study small area estimation under this spatio-temporal model, deriving first model parameter estimators, then obtaining EBLUPs of the area means and finally describing a parametric bootstrap procedure for estimation of the MSE of the EBLUPs. Simulation studies analyze the performance of

Download English Version:

<https://daneshyari.com/en/article/417506>

Download Persian Version:

<https://daneshyari.com/article/417506>

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