

## LR5\_3 Area-level small area estimation methods for domain statistics

Area-level small area estimation methods for domain statistics	
General setting	
<b>Configuration</b>	Configuration 5
<b>Type of sources</b>	Combination of aggregate-level survey and register data
<b>Statistic of interest</b>	Domain estimates in social or business statistics
<b>Type of errors</b>	Sampling and measurement errors
<b>Quality measure</b>	Estimation of mean squared error for the statistics of interest
Method (use one block per method)	
<b>Description</b>	<p>Design-based or direct survey estimates for small subpopulations often have large sampling variances due to small sample sizes. By linking the subpopulation quantities of interest through a model it is usually possible to obtain estimates with smaller mean squared errors. This is the subject of small area estimation (Rao and Molina, 2015). The subpopulations are not restricted to geographical areas, but can be anything as long as they are mutually exclusive, there are not too few of them, and the quantities of interest show some degree of similarity between the areas.</p> <p>A popular model in small area estimation is the basic area-level or Fay-Herriot model (Fay and Herriot, 1979). The input data for this model consist of a set of survey estimates <math>y_i</math> and corresponding variance estimates <math>\psi_i</math>, as well as auxiliary data <math>X_i</math>, where each <math>X_i</math> can be a vector. The index <math>i</math> refers to the subpopulations (areas) of interest in a particular study, such as municipalities, and <math>y_i</math> are direct or initial estimates for the unknown quantities of interest <math>\theta_i</math>. The aim is to find estimates for <math>\theta_i</math>, which improve on the initial estimates <math>y_i</math>. The model is</p> $y_i = \theta_i + e_i \tag{1}$ $\theta_i = \beta'X_i + v_i$ <p>with independent errors</p> $e_i \sim N(0, \psi_i) \quad \text{and} \quad v_i \sim N(0, \sigma_v^2)$ <p>The first line of (1) is like a measurement equation relating <math>y_i</math> to the latent quantities <math>\theta_i</math>, while the second line is the structural part of the model relating the quantities of interest to auxiliary information and to each other via common regression coefficients <math>\beta</math>.</p> <p>Equations (1) can be combined into</p> $y_i = \beta'X_i + v_i + e_i$

	<p>which is a special case of a linear mixed model. The Empirical Best Linear Unbiased Predictor (EBLUP) based on this model is</p> $\hat{\theta}_i = \hat{\beta}'X_i + \hat{v}_i = \hat{\gamma}_i y_i + (1 - \hat{\gamma}_i) \hat{\beta}'X_i$ <p>where</p> $\hat{\beta} = \left( \sum_i \hat{\gamma}_i X_i X_i' \right)^{-1} \sum_i \hat{\gamma}_i X_i y_i$ $\hat{\gamma}_i = \frac{\hat{\sigma}_v^2}{\hat{\sigma}_v^2 + \psi_i}$ <p>Various methods are available to find <math>\hat{\sigma}_v^2</math>, including (restricted) maximum likelihood, see Rao and Molina (2015).</p> <p>Given <math>\sigma_v^2 = \hat{\sigma}_v^2</math>, the mean squared error (mse) for the EBLUP is</p> $mse(\hat{\theta}_i) = \hat{\gamma}_i \psi_i + \hat{\sigma}_v^2 (1 - \hat{\gamma}_i)^2 X_i' \left( \sum_j \hat{\gamma}_j X_j X_j' \right)^{-1} X_i$ <p>An additional contribution to the mse comes from uncertainty in the estimate <math>\hat{\sigma}_v^2</math>. This contribution is usually relatively small, but analytical expressions can be found in Rao and Molina (2015).</p>
<b>Assumptions</b>	<p>The method is based on the assumption that the model holds. In particular,</p> <ul style="list-style-type: none"> <li>• the input estimates <math>y_i</math> are unbiased for <math>\theta_i</math></li> <li>• the errors <math>v_i</math> and <math>e_i</math> are normally distributed</li> <li>• the sampling variances <math>\psi_i</math> are assumed known, whereas in practice they must be estimated</li> </ul>
<b>Advantages (benefits)</b>	<ul style="list-style-type: none"> <li>• the method borrows strength over areas, thereby achieving much smaller variances for areas with small sample sizes, at the cost of some bias</li> <li>• estimates can be obtained even for areas without observations</li> <li>• the basic area level model is appealing to many survey statisticians since the initial estimates can account for complex sampling design features and differential response rates.</li> <li>• as the model is specified at an aggregate level, it is typically not very computationally demanding to fit</li> </ul>
<b>Disadvantages (costs)</b>	<ul style="list-style-type: none"> <li>• only auxiliary information at the area level can be used in the model. The model cannot correct for unit-level selection bias. Such correction should already be incorporated in the input estimates.</li> <li>• mean squared errors may be underestimated because the sampling errors are assumed to be known</li> </ul>
<b>Case study (per method if needed)</b>	
<b>Agency – country</b>	Statistics Netherlands
<b>Topic</b>	Estimation of municipal unemployment based on Labour Force Survey data. Several estimation methods compared in a simulation study.
<b>Data sets used</b>	LFS data from several years was used to create an artificial population. Auxiliary

	variables from several registrations have been used, including registered unemployment.
<b>Results (e.g. different methods)</b>	Estimates based on the area-level model were found to be more accurate than survey regression estimates, which are similar to direct GREG estimates. Small area methods specified at the unit-level performed even better, for one thing because they better exploit auxiliary information which in this study was available at the unit-level.
<b>Report</b>	See Boonstra et al. (2009).
<b>Final remarks</b>	
<b>Gap analysis</b>	The method has already been extended in most ways one can think of. These extensions include time-series extensions, spatial correlation, modelling of sampling error, non-linear specifications, etc. It would be interesting to combine this line of parametric aggregate-level modelling with macro-integration methods for incorporating equality and inequality restrictions.
<b>Other remarks</b>	
<b>References</b>	<p>Boonstra, H.J., van den Brakel, J., Buelens, B., Krieg, S. and M. Smeets (2008). Towards small area estimation at Statistics Netherlands. <i>Metron</i> 66 (1), pp. 21-49.</p> <p>Fay, R. E., and R.A. Herriot (1979). Estimates of income for small places: an application of James-Stein procedures to census data. <i>Journal of the American Statistical Association</i> 74 (366a), pp. 269-277.</p> <p>Rao, J.N.K. and I. Molina (2015), <i>Small Area Estimation</i>. John Wiley.</p>