

Spatial Disaggregation in Epidemiology

Diego Ayma ¹ María Durbán ¹ Dae-Jin Lee ²

¹Department of Statistics, Universidad Carlos III de Madrid, Spain

²BCAM - Basque Center for Applied Mathematics, Spain

January 20, 2015

Motivation

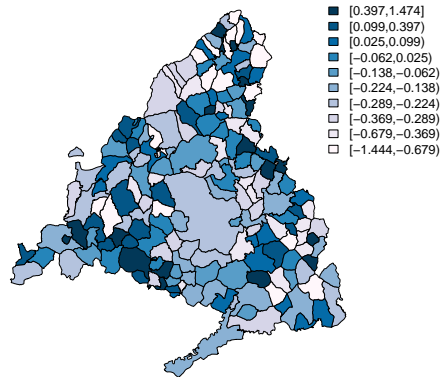
- Disease mapping studies deal with data that are, in general, **aggregated** over geographical units.
- A **choropleth map** provides an essential tool for the description of these data:
 - ✓ Quick visualization of spatial distribution of rates or risks.
 - ⚠ Choice of the colour scale in the legend map is arbitrary!
Epidemiologists prefer to use quantiles as class boundaries.

Motivation

Example

- Standardized mortality ratio (SMR) for female deaths by cardiovascular diseases in the Community of Madrid (2001).
- $\log(SMR_k) = \log(o_k / e_k)$, where o_k and e_k are the observed and expected numbers of female deaths in municipality k , for $k = 1, \dots, 179$.
- The class boundaries correspond to the deciles of raw $\log(SMR)$ s.

Raw $\log(SMR)$ by municipalities



Motivation

- Raw $\log(SMR)$ s vary abruptly between geographical units, making hard to detect meaningful underlying patterns.
- ✓ **Solution:** use **smoothing methods** to enhance the visibility of underlying trends in noisy data.
- One approach: use a spatial two dimensional P-spline (represented as a mixed model), at the centroids of the geographical units, to model that trend (Lee and Durbán, 2009).
 - ★ We will apply the Poisson version of this approach (**Poisson P-GLMM**) to the previous raw data.

Motivation

Example (continued)

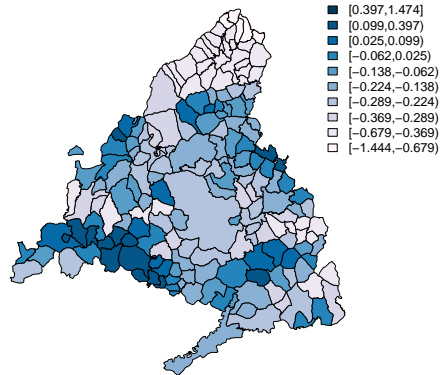
- **Poisson P-GLMM**. Given $\mathbf{o} \sim \text{Poiss}(\boldsymbol{\mu})$,

$$\boldsymbol{\mu} = \mathbf{e} \exp(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha}),$$

with $\boldsymbol{\alpha} \sim \mathcal{N}(\mathbf{0}, \mathbf{G}(\lambda_1, \lambda_2))$, where \mathbf{X} and \mathbf{Z} are constructed from the centroid coordinates of municipalities.

- Parameters $\boldsymbol{\beta}$, $\boldsymbol{\alpha}$, λ_1 , and λ_2 are estimated using the Penalized quasi-likelihood (PQL) approach of Breslow and Clayton (1993).

Smooth log(SMR) by municipalities



Motivation

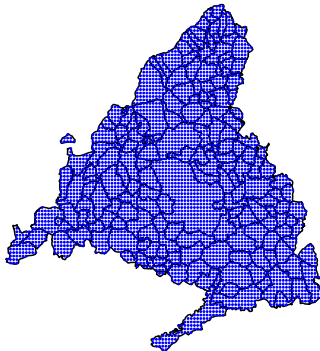
Some limitations:

- Choropleth maps present the common biases visual perception that larger rural and sparsely populated areas are of greater importance.
- Usually, there exists a mismatch of spatial supports between health data and relevant explanatory variables, making impossible the direct analysis of their possible relationship.

Motivation

Goal: Estimate a *latent* spatial trend for $\log(SMR)$ s at a desirable *fine* scale, using available data. For example:

From municipalities (179) to a fine grid (4359)



Motivation

- **Proposal:** Extend the penalized composite link model (P-CLM, Eilers, 2007) to the spatial case, into a mixed model framework.
- The resulting model will be called (spatial) penalized composite link mixed model, or **Poisson P-CLMM**.

The model

Poisson P-CLMM. Given $\mathbf{o} \sim \text{Poiss}(\boldsymbol{\mu})$,

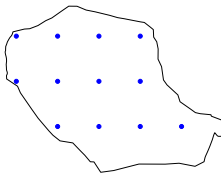
$$\boldsymbol{\mu} = \mathbf{C}\boldsymbol{\gamma} \leftarrow \boldsymbol{\gamma} = \mathbf{e} \exp(\boldsymbol{\eta}) \leftarrow \boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha},$$

with $\boldsymbol{\alpha} \sim \mathcal{N}(\mathbf{0}, \mathbf{G}(\lambda_1, \lambda_2))$. Here:

- $\boldsymbol{\gamma}$: Latent expectation at the fine scale.
- \mathbf{C} : Composition matrix that describes how latent expectations was combined before generating the data (must be known).
- $\boldsymbol{\eta}$: Latent spatial trend at the fine scale, where \mathbf{X} and \mathbf{Z} are the fixed and random effects matrices, constructed from fine grid coordinates.

The model

- e : expected values at the fine scale. If it is not known, we may split e_k into equal parts, over the corresponding fine scale (\hat{e}_{naive}).
- For example, for municipality 4, we have $e_4 = 49.9158$. Then, we split it in 12 equal parts.



Parameters estimation and approximate errors

Parameters β , α , λ_1 , and λ_2 are estimated using a modified PQL approach. The estimated fixed and random parameters are:

$$\begin{aligned}\hat{\beta} &= \left(\check{\mathbf{X}}^T \mathbf{V}^{-1} \check{\mathbf{X}} \right)^{-1} \check{\mathbf{X}}^T \mathbf{V}^{-1} \mathbf{z}, \\ \hat{\alpha} &= \mathbf{G} \check{\mathbf{Z}}^T \mathbf{V}^{-1} \left(\mathbf{z} - \check{\mathbf{X}} \hat{\beta} \right),\end{aligned}$$

where $\mathbf{V} = \mathbf{W}^{-1} + \check{\mathbf{Z}} \mathbf{G} \check{\mathbf{Z}}^T$, $\check{\mathbf{X}} = \mathbf{W}^{-1} \mathbf{C} \mathbf{\Gamma} \mathbf{X}$, $\check{\mathbf{Z}} = \mathbf{W}^{-1} \mathbf{C} \mathbf{\Gamma} \mathbf{Z}$, with matrices $\mathbf{W} = \text{diag}(\mu)$ and $\mathbf{\Gamma} = \text{diag}(\gamma)$ and working vector $\mathbf{z} = \check{\mathbf{X}} \beta + \check{\mathbf{Z}} \alpha + \mathbf{W}^{-1}(\mathbf{o} - \mu)$.

Parameters estimation and approximate errors

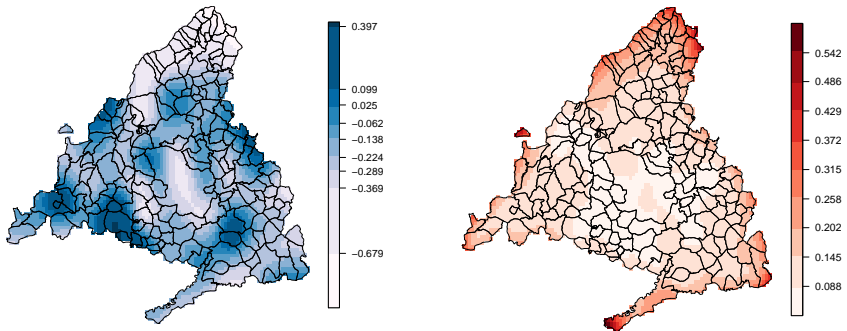
The approximate variance-covariance matrix for $\mathbf{x}\hat{\beta} + \mathbf{z}\hat{\alpha}$ is:

$$\text{Var}(\mathbf{x}\hat{\beta} + \mathbf{z}\hat{\alpha}) \approx [\mathbf{x} \mid \mathbf{z}] \begin{bmatrix} \check{\mathbf{x}}^T \mathbf{W} \check{\mathbf{x}} & \check{\mathbf{x}}^T \mathbf{W} \check{\mathbf{z}} \\ \check{\mathbf{z}}^T \mathbf{W} \check{\mathbf{x}} & \check{\mathbf{z}}^T \mathbf{W} \check{\mathbf{z}} + \mathbf{G}^{-1} \end{bmatrix}^{-1} [\mathbf{x} \mid \mathbf{z}]^T.$$

The approximate errors for $\mathbf{x}\hat{\beta} + \mathbf{z}\hat{\alpha}$ correspond to the squared root of the diagonal values of this matrix.

Application

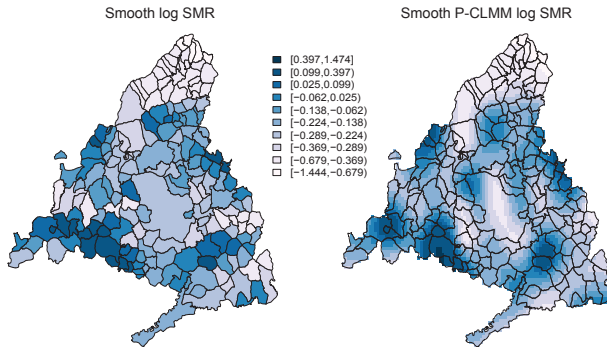
From municipalities to a fine grid:



Here, we use equispaced cut-points for the error map legend.

Application

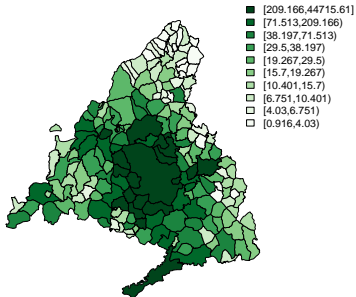
Compare (coarse) smooth P-GLMM $\log(SMR)$ with (continuous) smooth P-CLMM $\log(SMR)$...



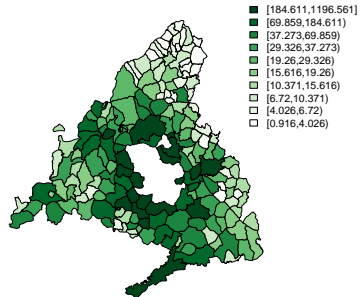
Issues

- Other alternatives to estimate e ?
- We tried to estimate it using a Poisson and Gaussian spatial P-CLMM approach, but we had numerical problems. There are extreme raw expected values (for example, Madrid) ...

Raw expected values

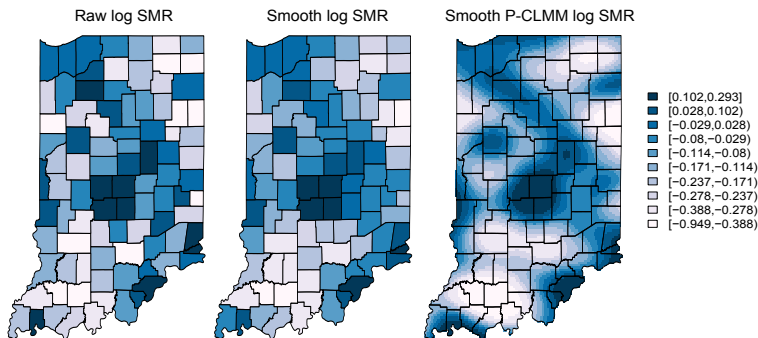


Raw expected values (without Madrid)



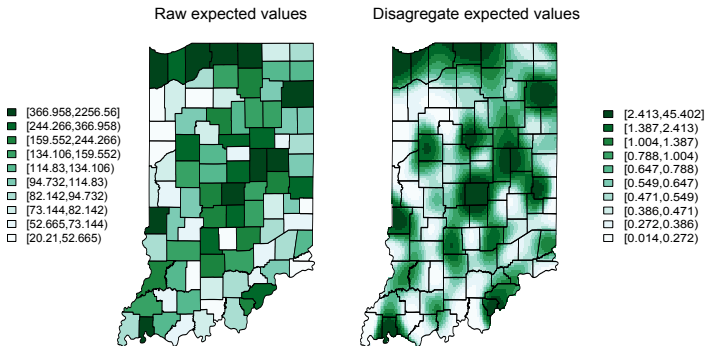
Indiana data set

But, look at this new data set. The Raw $\log(SMRs)$ correspond to white female deaths by lung cancer in Indiana, USA, by (92) states. The smooth P-CLMM $\log(SMRs)$ were obtained using \hat{e}_{naive}



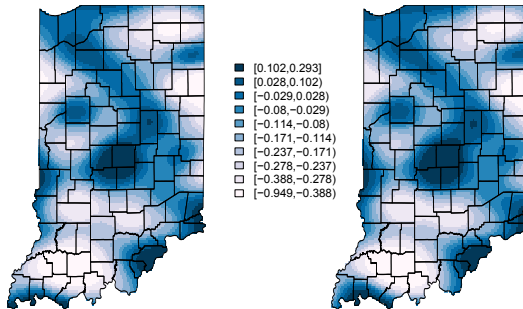
Indiana data set

In this case, we were able to estimate an underlying surface for expected values using a P-CLMM approach. Look at the raw expected values and its disaggregation:



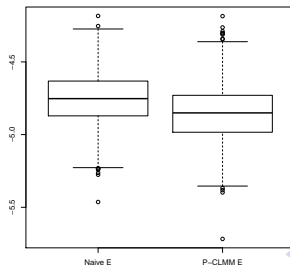
Indiana data set

If we use this estimate (\hat{e}_{P-CLMM}), we obtain the smooth surface in left panel. Compare with the smooth surface obtained using \hat{e}_{naive} (right panel):



Simulation

If we consider the smooth surface obtained using \hat{e}_{P-CLMM} , we can simulate $n = 1000$ observed values (at the fine grid) as Poisson realizations, with intensity $\hat{e}_{P-CLMM} * \hat{\eta}_{\hat{e}_{P-CLMM}}$, aggregated them, and then apply our method using \hat{e}_{naive} and \hat{e}_{P-CLMM} as expected values at the fine scale. The boxplot of the MSEs (in log scale) for each case is the following:



Thanks!

References

- **Lee, D.-J. and Durbán, M. (2009)** Smooth-CAR mixed models for count data. Computational Statistics & Data Analysis, 53:2958-2979.
- **Eilers, P. H. C. (2007).** Ill-posed problems with counts, the composite link model and penalized likelihood. Statistical Modelling, 7:239-254.
- **Breslow, N. E. and Clayton, D. G. (1993).** Approximate inference in generalized linear mixed models. Journal of the American Statistical Association, 88:9-25.

Appendix

Mixed model formulation:

Let \mathbf{x}_1 and \mathbf{x}_2 be the lon-lat coordinates, at the fine scale. Consider B-spline bases $\mathbf{B}_i = \mathbf{B}(\mathbf{x}_i)$, and difference matrices $\mathbf{D}_i = \mathbf{D}_i(q_i)$, $i = 1, 2$. Using the SVD of $\mathbf{D}_i^T \mathbf{D}_i$,

$$\mathbf{U}_i \mathbf{\Sigma}_i \mathbf{U}_i^T,$$

with partitions $\mathbf{U}_i = [\mathbf{U}_{in} : \mathbf{U}_{is}]$ and $\mathbf{\Sigma}_i = \text{BlockDiag}(\mathbf{0}_{q_i}, \tilde{\mathbf{\Sigma}}_i)$, we can define:

$$\mathbf{X} = \mathbf{X}_2 \square \mathbf{X}_1 \text{ and } \mathbf{Z} = [\mathbf{Z}_2 \square \mathbf{X}_1 : \mathbf{X}_2 \square \mathbf{Z}_1 : \mathbf{Z}_2 \square \mathbf{Z}_1],$$

where $\mathbf{X}_i = \mathbf{B}_i \mathbf{U}_{in}$ and $\mathbf{Z}_i = \mathbf{B}_i \mathbf{U}_{is}$.

Appendix

It can be shown that $\mathbf{G}(\lambda_1, \lambda_2)$ is reformulated as $\mathbf{G}(\lambda_1, \lambda_2) = \mathbf{F}^{-1}$, where:

$$\mathbf{F} = \begin{bmatrix} \lambda_2 \tilde{\Sigma}_2 \otimes \mathbf{I}_{q_1} & & \\ & \lambda_1 \mathbf{I}_{q_2} \otimes \tilde{\Sigma}_1 & \\ & & \lambda_1 \mathbf{I}_{c_2 - q_2} \otimes \tilde{\Sigma}_1 + \lambda_2 \tilde{\Sigma}_2 \otimes \mathbf{I}_{c_1 - q_1} \end{bmatrix}$$

with $c_1 = \text{ncol}(\mathbf{B}_1)$ and $c_2 = \text{ncol}(\mathbf{B}_2)$.