

Bayesian modeling and clustering for spatio-temporal areal data: An application to Italian Fertility Rate

Federico Busetto, Daniele Fedrizzi, Niccolò Fontana,
Emre Kayalioglu, Riccardo Lazzarini, Anna Lunghi

Politecnico di Milano

February 14, 2023



Fertility in Italy

During the last years, Italy has been experiencing a significant decrease in the birth rate. For the first time ever, in 2021 the births were below 400.000.

The birth rate has declined across the country but with a different spread across geographic areas and provinces.

It's important to understand which socio-economic factors influence the most this trend and to explain in which way the geographic location has an impact on this phenomenon.

Fertility in Italy - ISTAT Data

ISTAT has collected 10 years worth of data (from 2011 to 2020) about the **fertility rate** (the average number of children per woman of childbearing age, 15-49 years) in each province of Italy.

Alongside the fertility rate, the **average age of the mother at delivery** and the **average age of the father at the birth of the child** were computed.

We enriched this dataset with some socio-economic indices, such as **employment** and **inactivity** rate, in order to better explain the birth rate decrease phenomenon.

Fertility in Italy - ISTAT Data

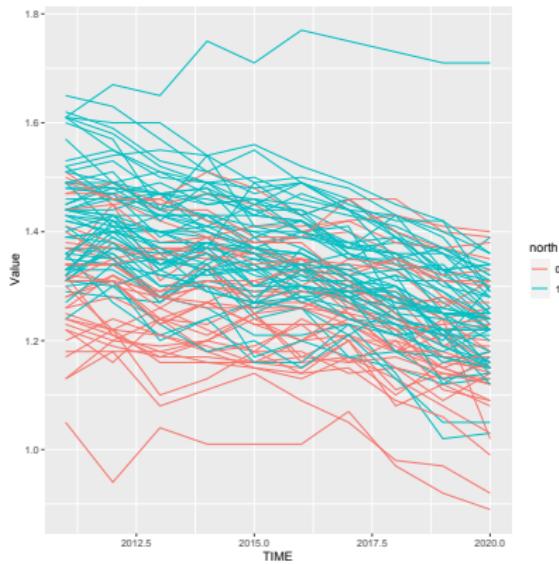
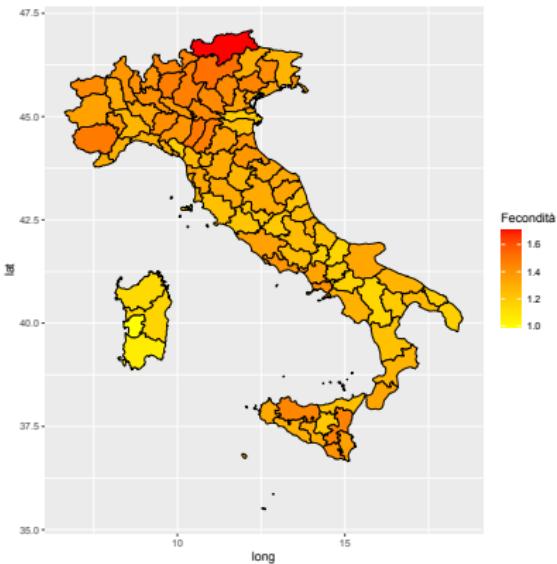


Figure: Average fertility rate in each province (left) and time-series of fertility rate divided into north and south (right)

Fertility in Italy - ISTAT Data

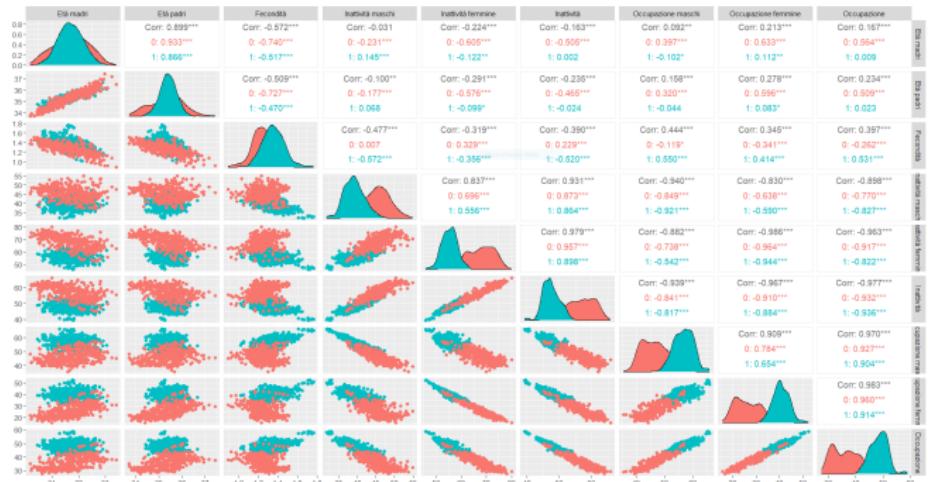


Figure: GGpairs of the dataset divided between north and south

Fertility rate is positively correlated with employment and negatively with inactivity rate. Fertility rate is negatively correlated with both the age of the mother and the age of the father. Employment and inactivity rate are strongly negatively correlated. north-south duality is reflected on the economic variables

Fertility in Italy - ISTAT Data

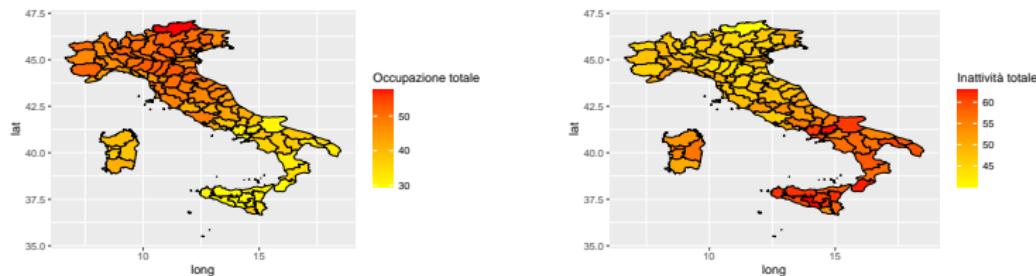


Figure: Average employment rate in each province (left) and average inactivity rate in each province (right)

	Moran's I	Geary's C
Fertility Rate	0.4378463	0.4286952
Age of the Mother	0.3901207	0.5001999
Female Inactivity	0.8496539	0.1227897
Male Inactivity	0.6619759	0.2628052
Female Employment	0.8588112	0.1089427
Male Employment	0.8010590	0.1487005

Full Model

Following paper (1) we consider the following model.

For all $i = 1, \dots, I$ and $t = 1, \dots, T$:

$$\begin{aligned} Y_{it} | \mathbf{x}_{it}, \boldsymbol{\beta}_{s_i}^*, w_{it}, \sigma^2, s_i &\stackrel{\text{ind}}{\sim} N(\mathbf{x}_{it}^\top \boldsymbol{\beta}_{s_i}^* + w_{it}, \sigma^2) \\ \mathbf{w}_t | \mathbf{w}_{t-1}, \xi, \tau^2, \rho, W &\sim N_I(\xi \mathbf{w}_{t-1}, \tau^2 Q(\rho, W)^{-1}) \\ \mathbf{w}_1 | \tau^2, \rho, W &\sim N_I(\mathbf{0}, \tau^2 Q(\rho, W)^{-1}) \\ \xi &\sim N(\mu_\xi, \sigma_\xi^2) \\ \sigma^2 &\sim \text{Inv-Gamma}(a_{\sigma^2}, b_{\sigma^2}) \\ \tau^2 &\sim \text{Inv-Gamma}(a_{\tau^2}, b_{\tau^2}) \\ \mathbf{s} | \alpha &\sim \text{PolyaUrn}(\mathbf{s} | \alpha) \\ \alpha &\sim \text{Gamma}(a_\alpha, b_\alpha) \\ \boldsymbol{\beta}_1^*, \dots, \boldsymbol{\beta}_{K_I}^* | \boldsymbol{\mu}_\beta, \Sigma_\beta &\stackrel{\text{iid}}{\sim} P_0, \\ P_0(d\boldsymbol{\beta}^*) &= N_{p+1}(d\boldsymbol{\beta}^* | \boldsymbol{\mu}_0, \Sigma_0) \end{aligned} \tag{1}$$

where $Q(\rho, W) = \rho(\text{diag}(W\mathbf{1}) - W) + (1 - \rho)\mathbb{I}_I$

Equivalent Formulation

For all $i = 1, \dots, I$ and $t = 1, \dots, T$:

$$\begin{aligned} \mathbf{Y}_i | X_i, \mathbf{w}_i, \sigma^2, P &\stackrel{\text{ind}}{\sim} \int_{\mathbb{R}^{p+1}} N_T(X_i \boldsymbol{\beta} + \mathbf{w}_i, \sigma^2 I) P(d\boldsymbol{\beta}) \\ P &\sim DP(\alpha, P_0) \\ \mathbf{w}_t | \mathbf{w}_{t-1}, \xi, \tau^2, \rho, W &\sim N_I(\xi \mathbf{w}_{t-1}, \tau^2 Q(\rho, W)^{-1}) \\ \mathbf{w}_1 | \tau^2, \rho, W &\sim N_I(\mathbf{0}, \tau^2 Q(\rho, W)^{-1}) \\ \xi &\sim N(\mu_\xi, \sigma_\xi^2) \\ \sigma^2 &\sim \text{Inv-Gamma}(a_{\sigma^2}, b_{\sigma^2}) \\ \tau^2 &\sim \text{Inv-Gamma}(a_{\tau^2}, b_{\tau^2}) \\ \alpha &\sim \text{Gamma}(a_\alpha, b_\alpha) \\ P_0(d\boldsymbol{\beta}) &= N_{p+1}(d\boldsymbol{\beta} | \boldsymbol{\mu}_0, \Sigma_0) \end{aligned} \tag{2}$$

where $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{iT})'$, $\mathbf{w}_i = (w_{i1}, \dots, w_{iT})'$ and X_i is a matrix of T rows, which are the vector \mathbf{x}'_{it} for $t = 1, \dots, T$.

Code Implementation

- Since the implementation of the MCMC algorithm for sampling from posterior distributions was performed in STAN, we used formulation 2 that did not include a Pólya Urn scheme as a prior for s .
- The Pólya Urn scheme is in fact a discrete distribution and therefore incompatible with the Hamiltonian Monte Carlo, that is the MCMC method used by STAN to sample from the posteriors.
- Using Stick-Breaking construction for the first line of 2 and truncating it to H equal to 10, a compromise value between the performance of the code and the goodness of the approximation, we obtained a finite mixture model

Code Optimization

- Fixing $\rho = 0.95$ implies fixing the covariance matrix of the random effects, that can be factorized with the Cholesky decomposition, obtaining L such that $L^\top L = Q(\rho, W)^{-1}$
- We can reparametrize the random effects in the following more efficient way

$$\boldsymbol{w}_1 = L\boldsymbol{w}_1^{raw}$$

$$\boldsymbol{w}_t = \xi \boldsymbol{w}_{t-1} + L\boldsymbol{w}_t^{raw}$$

$$\boldsymbol{w}_t^{raw} \stackrel{\text{iid}}{\sim} N_I(\mathbf{0}, \tau^2 \mathbb{I}_{p+1})$$

- We used data structures that made it easy to use STAN *vectorization*, both in terms of matrix operations and vectorized probability functions
- We avoided using transpositions within the definition loops and adopted the most efficient solutions for accessing rows and columns of random effects

Simulation Study

To verify the correctness of the implementation and show the ability of the model to correctly estimate parameters and cluster areal locations we performed a study on simulated data.

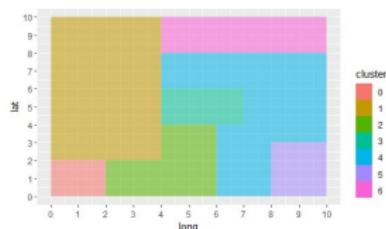


Figure: Underlying cluster structure for simulated data.

- Each areal location had three covariates and four regression coefficients
- ρ is set to 0.95
- $\tau^2 = \sigma^2 = 1$
- we used no prior for μ_0 and σ_0^2 , set equal to zero and one respectively, as there was no reason to further facilitate the detection of clusters.

Simulation Study

- The density plots show that the true β coefficients are accurately recovered, while the trace plots imply good mixing and convergence.
- Estimating the clustering with the Binder loss and the VI we obtained the same results and we were able to find all the 7 clusters from which the data came, which still underlines the goodness of the used model.

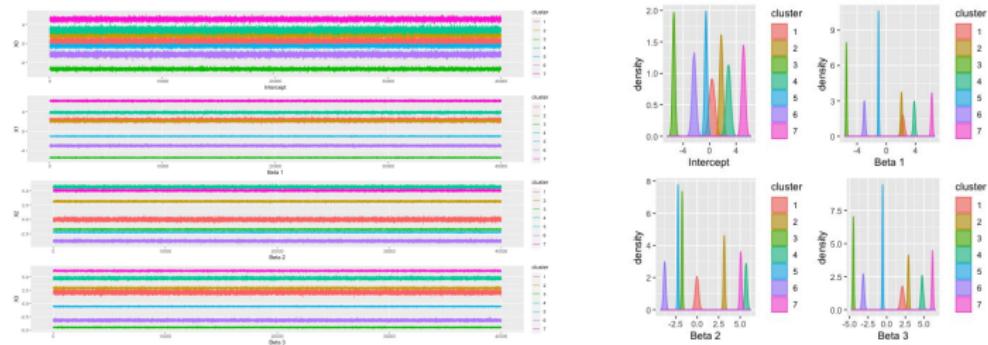


Figure: Posterior inference of the cluster-specific regression coefficients β .

Prior Elicitation

To help the model in recognizing clusters, we put priors on the hyperparameters of the base measure of the Dirichlet Process:

$$P_0(d\beta^*) = N_{p+1}(d\beta^* | \mu_0, \Sigma_0)$$

$$\mu_0 \sim N_{p+1}(\eta_0, \tau_0^2 \mathbb{I}_{p+1})$$

$$\Sigma_0 = \sigma_0^2 \mathbb{I}_{p+1}$$

$$\sigma_0^2 \sim \text{Inv-Gamma}(a_{\sigma_0^2}, b_{\sigma_0^2})$$

Following (1), we fixed the hyperparameters of α to $a_\alpha = 3$ and $b_\alpha = 2$, in order to obtain a prior expected number of clusters of 6.75 with prior variance of number of clusters equal to 7.32.

Posterior Inference

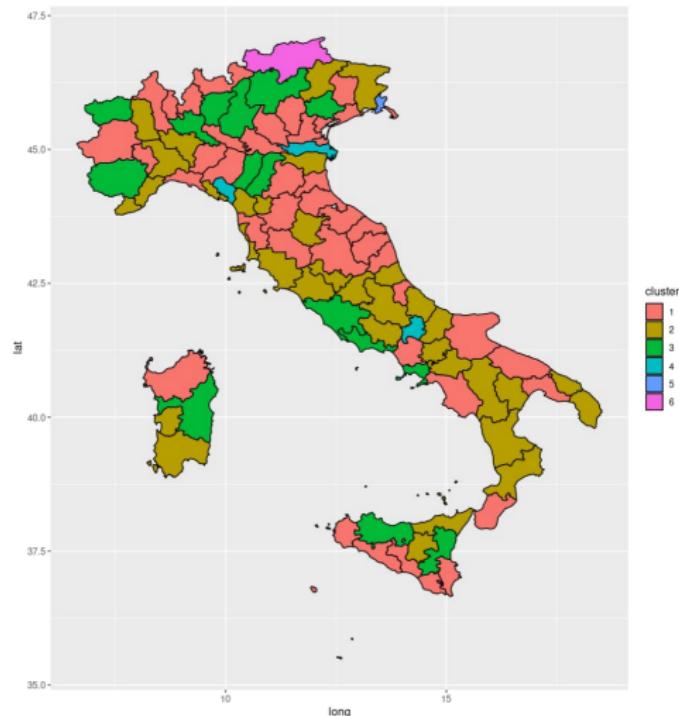


Figure: Partition of the Italian provinces obtained by minimizing the posterior expectation of the VI loss function

Posterior Inference

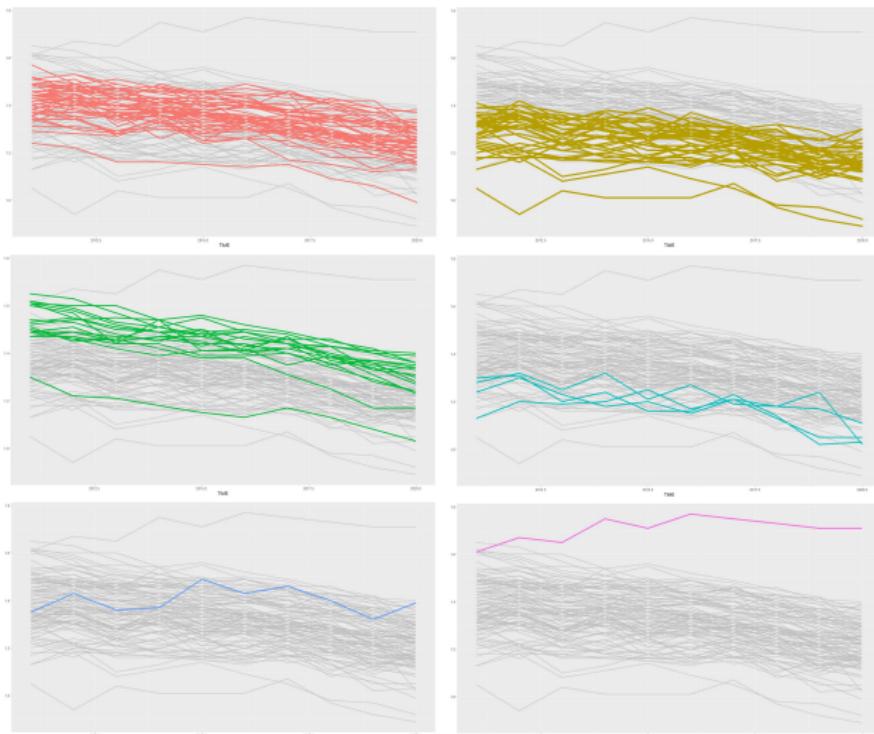


Figure: Fertility rate time-series divided for each cluster

Posterior Inference

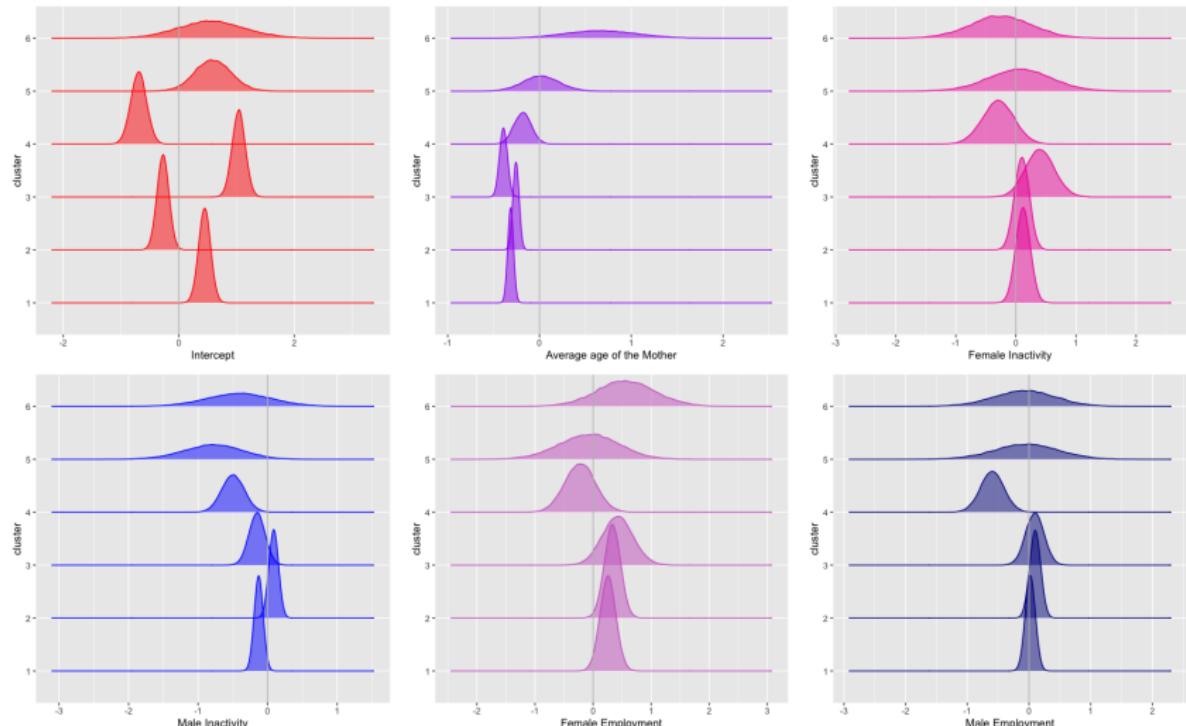


Figure: Posterior kernel density estimates of the regression coefficients

Posterior Inference

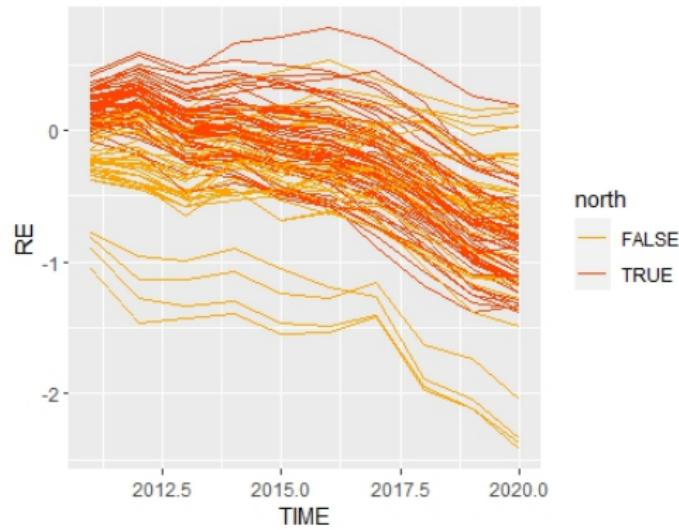
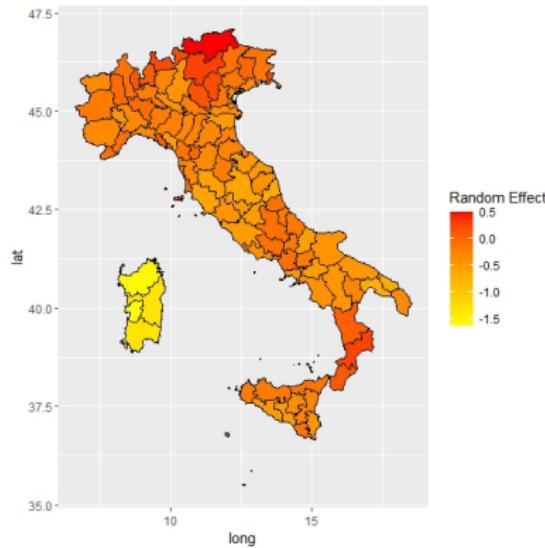


Figure: *Estimated spatio-temporal random effects*

MCMC Diagnostics

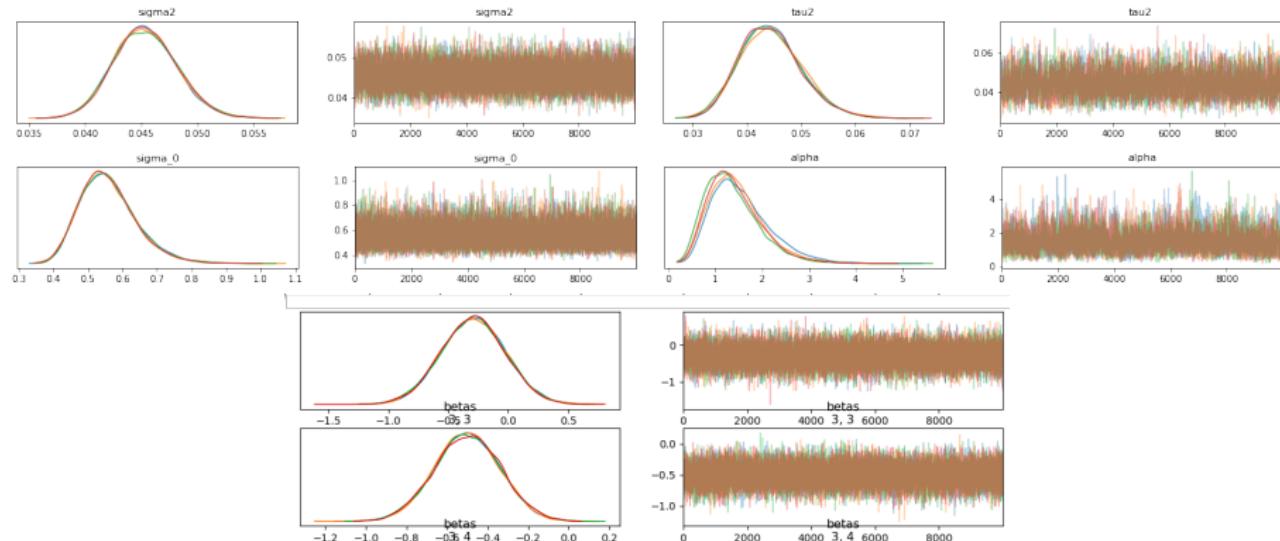


Figure: Traceplots for some variables in the state of the MCMC

Competitor Model

ST.CARar() model in R-package CARBayesST:

$$\begin{aligned} Y_{it} | \mathbf{x}_{it}, \boldsymbol{\beta}, w_{it}, \sigma^2 &\stackrel{\text{ind}}{\sim} N(\mathbf{x}_{it}^\top \boldsymbol{\beta} + w_{it}, \sigma^2) \\ \mathbf{w}_t | \mathbf{w}_{t-1}, \xi, \tau^2, \rho, W &\sim N_I(\xi \mathbf{w}_{t-1}, \tau^2 Q(\rho, W)^{-1}) \\ \mathbf{w}_1 | \tau^2, \rho, W &\sim N_I(\mathbf{0}, \tau^2 Q(\rho, W)^{-1}) \\ \boldsymbol{\beta} &\sim N_{p+1}(\mu_\beta, \Sigma_\beta) \\ \xi &\sim U([0, 1]) \\ \rho &\sim U([0, 1]) \\ \sigma^2 &\sim \text{Inv-Gamma}(a_{\sigma^2}, b_{\sigma^2}) \\ \tau^2 &\sim \text{Inv-Gamma}(a_{\tau^2}, b_{\tau^2}) \end{aligned} \tag{3}$$

	WAIC	LPML
ST.CARar()	-29	-191
Model (1)	4	-62

References

- [1] Alexander Mozdzen, Andrea Cremaschi, Annalisa Cadonna, Alessandra Guglielmi, and Gregor Kastner. Bayesian modeling and clustering for spatio-temporal areal data: An application to italian unemployment. *Spatial Statistics*, 2022
- [2] Stan Development Team. Stan user's guide, 2011–2022
- [3] Duncan Lee, Alastair Rushworth, Gary Napier, and William Pettersson. Carbayest version 3.3: Spatio- temporal areal unit modelling in r with conditional autoregressive priors
- [4] Rushworth A, Lee D, and Mitchell R. A spatio-temporal model for estimating the long- term effects of air pollution on respiratory hospital admissions in greater london. *Spatial and Spatio-temporal Epidemiology*, 2014
- [5] Redazione OpenPolis. Il calo delle nascite dopo l'emergenza covid, May 2022