

Contents lists available at [ScienceDirect](#)

Spatial Statistics

journal homepage: www.elsevier.com/locate/spasta

Capturing spatial dependence of COVID-19 case counts with cellphone mobility data

Justin J. Slater^{a,b,*}, Patrick E. Brown^{a,b}, Jeffrey S. Rosenthal^a, Jorge Mateu^c^a Department of Statistical Sciences, University of Toronto, Canada^b Centre for Global Health Research, St. Michael's Hospital, Canada^c Department of Mathematics, University Jaume I of Castellon, Spain

ARTICLE INFO

Article history:

Received 12 April 2021

Received in revised form 26 August 2021

Accepted 6 September 2021

Available online xxxx

Keywords:

Bayesian hierarchical model

Besag York Mollié model

COVID-19

Gaussian Markov random field

Mobility data

ABSTRACT

Spatial dependence is usually introduced into spatial models using some measure of physical proximity. When analysing COVID-19 case counts, this makes sense as regions that are close together are more likely to have more people moving between them, spreading the disease. However, using the actual number of trips between each region may explain COVID-19 case counts better than physical proximity. In this paper, we investigate the efficacy of using telecommunications-derived mobility data to induce spatial dependence in spatial models applied to two Spanish communities' COVID-19 case counts. We do this by extending Besag York Mollié (BYM) models to include both a physical adjacency effect, alongside a mobility effect. The mobility effect is given a Gaussian Markov random field prior, with the number of trips between regions as edge weights. We leverage modern parametrizations of BYM models to conclude that the number of people moving between regions better explains variation in COVID-19 case counts than physical proximity data. We suggest that this data should be used in conjunction with physical proximity data when developing spatial models for COVID-19 case counts.

© 2021 Elsevier B.V. All rights reserved.

* Corresponding author at: Department of Statistical Sciences, University of Toronto, Canada.

E-mail address: justin.slater@mail.utoronto.ca (J.J. Slater).

<https://doi.org/10.1016/j.spasta.2021.100540>

2211-6753/© 2021 Elsevier B.V. All rights reserved.

1. Introduction

Spatial analyses of COVID-19 case data were first published as early as March of 2020 (Huang et al., 2020; Arab-Mazar et al., 2020; Giuliani et al., 2020), in an attempt to characterize, predict, and attenuate the severity of the pandemic. Subsequent studies have noted substantial spatial dependence in COVID-19 case counts (Kang et al., 2020; Bilal et al., 2020). This makes sense as regions that are close to each other likely have more people moving between them, spreading the disease to nearby regions.

Many groups have attempted to model COVID-19 case counts as a function of climate (Liu et al., 2020; Shi et al., 2020; Briz-Redón and Serrano-Aroca, 2020), healthcare quality (Sugg et al., 2021), socioeconomic factors (Baum and Henry, 2020) and more. More recently, mobility data has become more abundant and popular for modelling COVID-19 transmission. This makes sense because the disease spreads through human contact, meaning that case counts are likely to be a function of the number of people moving around. Such mobility data has been used to model the evolution of the epidemic in Spain (Arândiga et al., 2020; Iacus et al., 2020), assess the effectiveness of the Spanish lockdown (Orea and Álvarez, 2020), monitoring the epidemic in Switzerland (Persson et al., 2021), identify at-risk populations in France during a lockdown (Pullano et al., 2020), individual-level infection tracing in China (Kraemer et al., 2020), assess the timing of stay-home orders (Audirac et al., 2020), and evaluate the effectiveness of social distancing in the United States (Badr et al., 2020). This data can be found in many forms, but is commonly found in the form of aggregated areal *mobility matrices*. If we denote a mobility matrix \mathbf{M} , $[\mathbf{M}]_{ij}$ corresponds to the number of trips from region i to region j , and \mathbf{M}_{ii} represents the number of trips within region i .

These data have been applied in a variety of different models to answer numerous questions, but lack of available methods makes it difficult for researchers to use this data to its full potential. In this paper, we demonstrate a novel method for analysing this data, whereby the mobility data is used as edge weights in a Gaussian Markov random field (network) model. Previous work using network models have been applied to mobility data in the form of a network compartment model (Chang et al., 2021) which was used to conduct inference regarding societal inequities, and inform reopening. This work does not aim to make such claims, but rather demonstrate the efficacy of mobility data in modern parametrizations of Besag, York, and Mollié (BYM) models (Besag et al., 1991) and their extensions.

BYM models have been used frequently in the spatial analysis literature due to their effectiveness and computational efficiency. In these models, the spatial component is comprised of Conditional Autoregressive (CAR) (Besag, 1974) models and conventional random effects. This means that the spatial effect of region i depends only on its “neighbours”. Neighbours could be defined by any quantity the analyst has access to, but is most often defined by physical adjacency, i.e. if two regions share a common border, they are considered neighbours. Several ICAR/BYM models have been applied to COVID-19 data with neighbours defined in this way (DiMaggio et al., 2020; Huang and Brown, 2021; Brainard et al., 2020). Although these spatial model components based on physical adjacency are powerful and computationally efficient, it makes more sense to use mobility between regions to induce spatial dependence in COVID-19 models because the disease spreads via person-to-person contact.

In this paper, we build a BYM model where mobility data is used to induce spatial dependence between regions. Using mobility data within two Communities in Spain, Madrid and Castilla-Leon, we demonstrate the value of mobility data for COVID-19 spatial modelling applications. Furthermore, we extend modern parametrizations of BYM models to account for both physical adjacency and mobility simultaneously, and show that mobility data captures spatial variation in COVID-19 case counts much more accurately than physical adjacency alone.

This is a short focused paper with the following plan. Section 2 presents the data and the modelling strategy based on particular parametrizations of BYM models. The results come in Section 3, and the paper ends with a final discussion in Section 4.

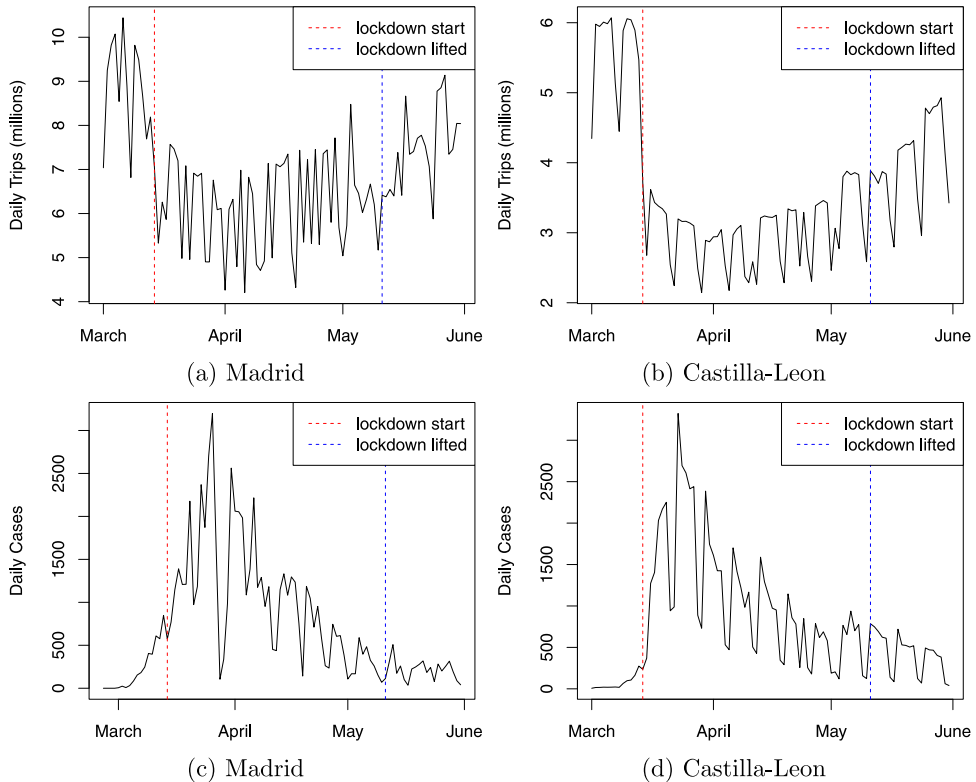


Fig. 1. Number of trips greater than 500 metres (a and b) and daily case counts (c and d) in the two Communities of Spain from March to June 2020.

2. Methods

2.1. Data

This paper is focused on two regions in Spain. Castilla-Leon is the largest Community in Spain by area and is located in the northwest part of Spain, with a population of 2.5 million. The Community of Madrid is located in the central part of Spain and has a population of around 6.8 million, and it is home of the capital of the country, Madrid City, with 3.3 million inhabitants.

The human mobility data was obtained from Barcelona Supercomputing Center Flow-map dashboard (Valencia, 2021). Trips within Madrid and Castilla-Leon were extracted from over 13 million phone records provided by a Spanish cellphone company. Both passive (GPS) and active (text messages, calls etc.) data were aggregated to construct daily movement matrices in each of the Communities, prior to the authors acquisition of the data. Given that trips were only recorded from one cellphone company, adjustment was made to estimate the number of total trips between each region. As a result, the entries of the mobility matrices are non-integer values.

Figs. 1(a) and 1(b) show the total daily movement between regions in Madrid, and Castilla-Leon, respectively. There is a sharp drop in the number of trips around March 14th 2020, which corresponds to a nation-wide lockdown. Lockdown restrictions began to ease around May 11th, where the number of trips slowly began to rise. Figs. 1(c) and 1(d) show the number of cases of COVID-19 cases in both Communities. COVID-19 daily cases data were retrieved from the open data portal of Castilla-Leon (General Directorate of Information Systems, Quality and Pharmaceutical

Provision, 2021) and from the Epidemiological Surveillance Network of Madrid (Anon, 2021). Notice that the movement drops as cases rise, because a lockdown was implemented in response to the increasing severity of the epidemic. In order to avoid this potential “reverse causality” problem, we will only use movement data in the first week of March. Our justification for this is that there is a time lag between when the virus spreads and the resulting COVID cases are confirmed. That is, the “first wave” of the epidemic was likely influenced mostly by the movement that occurred prior to the peak in cases, and less by the movement that occurred during it.

Fig. 2 shows the spatial distribution of the COVID-19 case rates up until May 31, 2020. The cases per thousand people range from (approximately) 0–30 in Madrid, and 0–100 for Castilla-Leon. We can see that there is substantial variation in the case rates within each of these Communities. Note that the extreme values in these plots are mostly small regions, which makes sense since the variance of case rates is higher when population is small. In the north of Madrid, there is a cluster of municipalities that have very low case rates. In Castilla-Leon, case rates are highest near the southeast border, which is the border to Madrid.

Fig. 3 shows the number of trips to, from, and within each Municipality of Madrid (there are 179 of these small regions), and Castilla-Leon (there are 245 health zones). Madrid and Castilla-Leon are considered separately throughout this paper. Although they are adjacent, data on movements between the two communities are not available. In Madrid, there is a lot of movement in and around Madrid City, and less movement in the more rural areas. Castilla-Leon shows a less predictable movement pattern, as there is not a single capital city that accounts for most of the movement. This movement data will be used to induce spatial correlation between regions, as described in Section 2.3.

2.2. Spatial autoregressive models

Besag, York, and Mollié (BYM) models (Besag et al., 1991) are widely used in spatial epidemiology and disease mapping due to their simplicity and computational efficiency. They assume the incidence of disease in region i follows a Poisson distribution

$$Y_i \sim \text{Pois}(E_i \lambda_i)$$

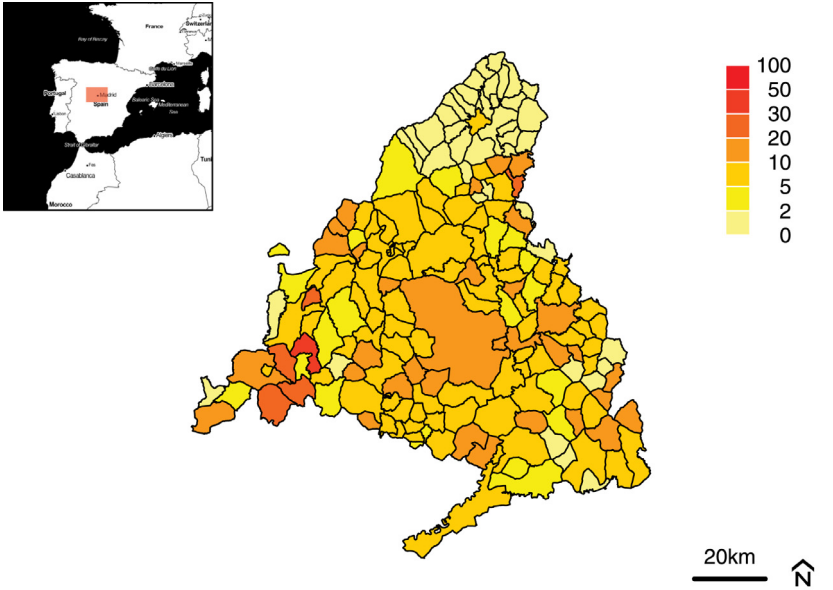
where Y_i is the number of infected cases in region i , and E_i is some form of expected count or offset, which could be the at-risk population, exposure time, etc. The log-relative risk, λ_i , is often modelled as

$$\begin{aligned} \log(\lambda_i) &= \mu + \beta X + \phi_i + \theta_i \\ \phi_i | \phi_{-i} &\sim N\left(\frac{1}{\sum_j w_{ij}} \sum_j w_{ij} \phi_j, \frac{\sigma_\phi^2}{\sum_j w_{ij}}\right) \\ \theta_i &\stackrel{i.i.d.}{\sim} N(0, \sigma_\theta^2) \end{aligned} \tag{1}$$

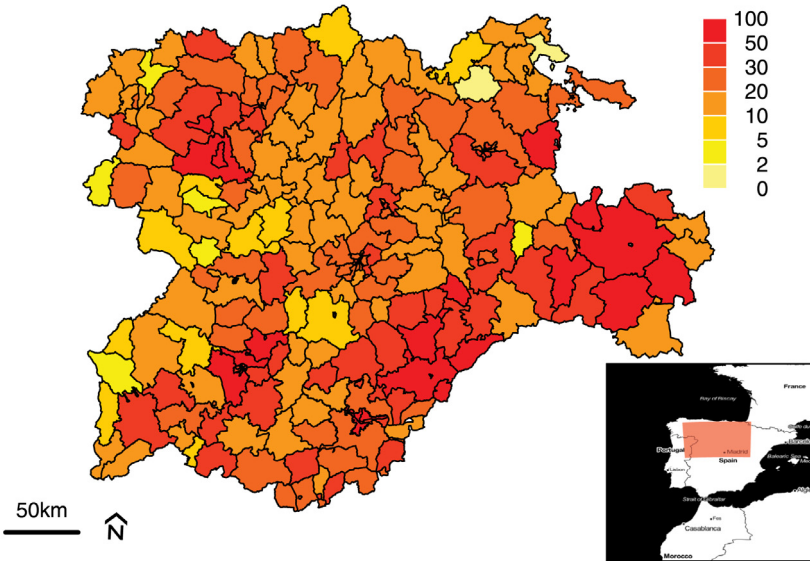
where μ is the overall intercept, β is the effect of spatial covariates, ϕ_i is the structured spatial random effect, and θ_i is the unstructured spatial random effect which allows for overdispersion in the response. In the spatial formulation of the BYM model, $w_{ij} = 1$ when regions i and j share a common border, and 0 otherwise. That is, region i 's structured spatial effect is only conditionally dependent on its neighbours, given all other regions. The distributions $\{\phi_i | \phi_{-i}\}_{i=1}^n$ are known as the *full conditionals*, where ϕ_{-i} is short hand for the set $\{\phi_1, \phi_2, \dots, \phi_{i-1}, \phi_{i+1}, \dots, \phi_n\}$. We can see from (1) that $E(\phi_i | \phi_{-i})$ is a weighted average of its neighbours, resulting in spatial smoothing. These full conditionals correspond to the joint distribution of the ϕ 's being a Gaussian Markov random field (GMRF) (Rue and Held, 2005), with

$$\begin{aligned} \phi &\sim \text{MVN}(\mathbf{0}, \mathbf{Q}^{-1}) \\ \mathbf{Q} &= \sigma_\phi^{-2} \mathbf{D}(\mathbf{I} - \mathbf{W}) \end{aligned}$$

where \mathbf{W} is a matrix of weights such that $w_{ij} > 0$ for $i \neq j$ and $w_{ii} = 0$, and σ_ϕ^2 is a variance parameter to be estimated. \mathbf{D} is a diagonal matrix such that $\mathbf{D}_{ii} = \sum_j w_{ij}$. This definition ensures



(a) Madrid



(b) Castilla-Leon

Fig. 2. COVID-19 cases per thousand, up to May 31 2020 for two communities in Spain. Background map ©Stamen Design.

that the precision matrix, \mathbf{Q} , is both symmetric and positive definite. In addition to the 0–1 weights based on regions being adjacent, other weighting schemes, such as inverse of Euclidean distance

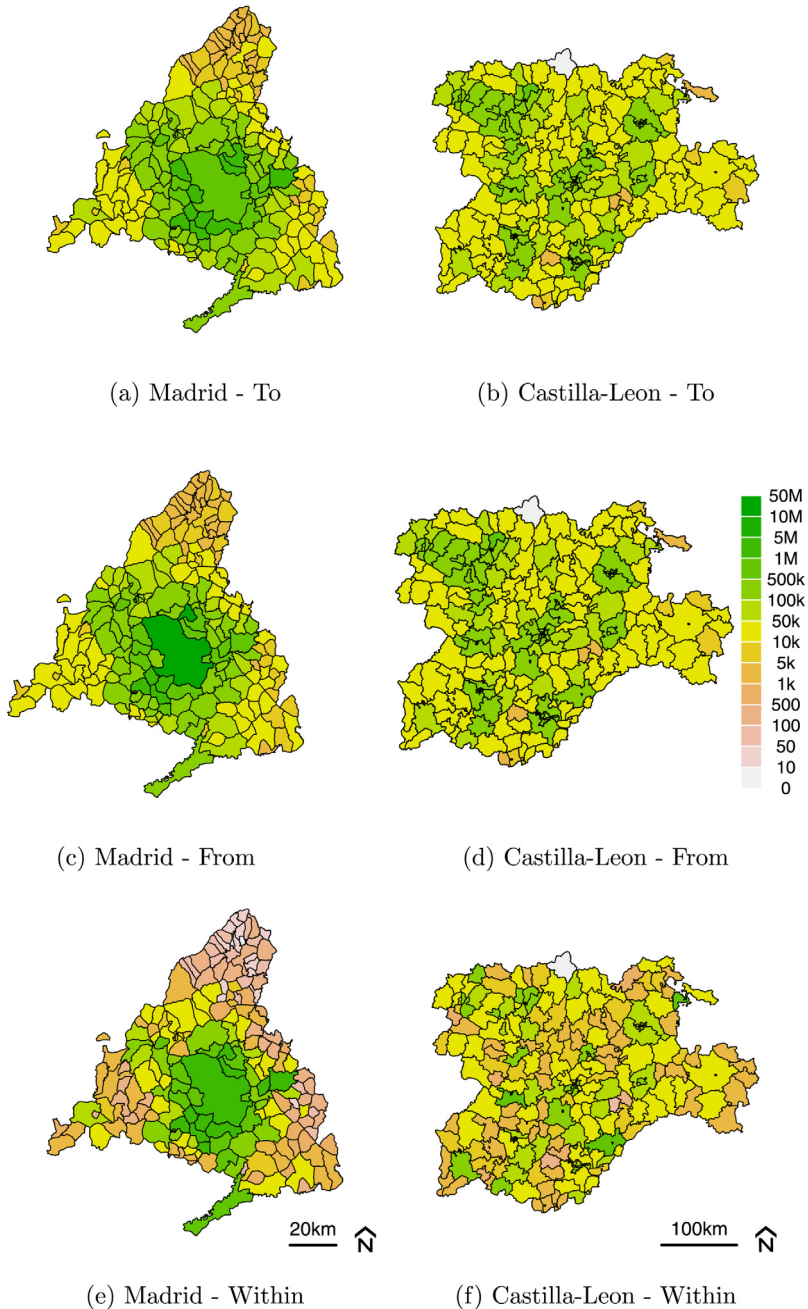


Fig. 3. Number of trips (incoming, outgoing, and within) the 179 regions of Madrid, and 245 health zones of Castilla-Leon, for the period March 1 to March 7 2020.

between regions, have been used. For a comparison of common weighting schemes, see [Duncan et al. \(2017\)](#). When we specify Q in this way, we refer to this as an Intrinsic Autoregressive (ICAR)

model for ϕ . The joint density function has a computationally convenient form with

$$p(\phi) \propto \exp\left[-\frac{1}{2\sigma_\phi^2} \sum_{i<j} w_{ij}(\phi_i - \phi_j)^2\right]$$

which is sometimes referred to as *the pairwise difference formula*. Notice that this density is invariant to the addition of a constant to each ϕ_i , leaving the spatial random effects unidentifiable up to a constant. This is typically remedied by imposing the constraint $\sum_i \phi_i = 0$ (Duncan et al., 2017). We will now modify this BYM model to account for movement between regions, in addition to physical adjacency.

2.3. Movement augmented BYM model

In order to extend the BYM model to allow for spatial correlation based on movement data, a second ICAR term, γ_i , with dependence structure governed by the movement data is added to the model. We also retain an adjacency-determined spatial effect ϕ_i in order to infer the relative importance of mobility-based and adjacency-based spatial dependence in determining COVID-19 case counts. The resulting model is

$$\begin{aligned} \log(\lambda_i) &= \mu + \beta X_i + \phi_i + \gamma_i + \theta_i \\ \phi_i | \phi_{-i} &\sim N\left(\frac{1}{\sum_j w_{ij}} \sum_j w_{ij} \phi_j, \frac{\sigma_\phi^2}{\sum_j w_{ij}}\right) \\ \gamma_i | \gamma_{-i} &\sim N\left(\frac{1}{\sum_j v_{ij}} \sum_j v_{ij} \gamma_j, \frac{\sigma_\gamma^2}{\sum_j v_{ij}}\right) \\ \theta_i &\sim N(0, \sigma_\theta^2) \end{aligned}$$

where ϕ_i and γ_i are the spatial random effects with priors based on the physical data and movement data respectively. The geographically-defined process ϕ_i has weights $w_{ij} = 1$ if regions i and j share a common border and are 0 otherwise, while the movement-defined process γ_i has weights v_{ij} representing the number of trips between regions i and j . Using mobility as edge weights in network models has shown to be effective in the context of infectious diseases (Schrödle et al., 2012; Volkova et al., 2010; Geilhufe et al., 2014). Schrödle et al. (2012) used mobility weights in an autoregressive term, which allowed the weights matrices to be asymmetric. However, given that our mobility data is being used in a Gaussian prior for a random effect, the precision matrices of ϕ and γ , Q_ϕ and Q_γ , must be symmetric. Therefore we require $w_{ij} = w_{ji}$ and $v_{ij} = v_{ji}$. While the first equality will always be true, the mobility matrices are not perfectly symmetric, thus symmetry was induced by defining v_{ij} as the sum of the numbers of trips from i to j and from j to i . The GRMF does not account for the movement within a region, so the movement within a region was included in the model as a spatial covariate X_i (fixed effect). That is, X_i was computed as

$$X_i = \frac{\frac{v_{ii}}{E_i} - \text{mean}\left(\frac{v_{jj}}{E_j}\right)}{\text{sd}\left(\frac{v_{jj}}{E_j}\right)}$$

where v_{ii}/E_i is the number of trips per person within a region, and $\text{mean}(v_{jj}/E_j)$ and $\text{sd}(v_{jj}/E_j)$ are the mean and standard deviations of the trips per person in all other regions. This model was run on both the Madrid and Castilla-Leon data.

There are two main drawbacks with the formulations of BYM models presented thus far. Firstly, the interpretation of the parameters σ_γ and σ_ϕ depend on the average number of neighbours and the total number of trips for each region, and hence their magnitudes are not comparable (Sørbye and Rue, 2014). Secondly, σ_ϕ , σ_γ , and σ_θ are hard to estimate without very careful choices of hyperpriors (Leroux et al., 2000). We will now address these shortcomings via reparametrizations.

2.4. Reparametrizations and priors

In order to solve issues with comparability, interpretability, and estimation, we apply a reparameterization of our model that is inspired by [Riebler et al. \(2016\)](#) with

$$\begin{aligned} \sigma^2 &\approx \text{Var}(\phi_i + \gamma_i + \theta_i) \\ \phi_i^* | \phi_{-i}^* &\sim N\left(\frac{1}{\sum_j w_{ij}} \sum_j w_{ij} \phi_j^*, \frac{\rho_\phi \sigma^2}{s_\phi \sum_j w_{ij}}\right) \\ \gamma_i^* | \gamma_{-i}^* &\sim N\left(\frac{1}{\sum_j v_{ij}} \sum_j v_{ij} \gamma_j^*, \frac{\rho_\gamma \sigma^2}{s_\gamma \sum_j v_{ij}}\right) \\ \theta_i &\sim N(0, \rho_\theta \sigma^2) \end{aligned}$$

where $\rho_\phi + \rho_\gamma + \rho_\theta = 1$ and $0 < \rho_\gamma, \rho_\phi, \rho_\theta < 1$. The priors for σ and ρ are

$$\begin{aligned} \sigma &\sim N_+(0, 1) \\ \rho &\sim \text{Dirichlet}(1, 1, 1) \end{aligned}$$

Note that

$$\begin{aligned} \phi_i^* &= \sigma \left(\sqrt{\rho_\phi / s_\phi}\right) \phi_i \\ \gamma_i^* &= \sigma \left(\sqrt{\rho_\gamma / s_\gamma}\right) \gamma_i. \end{aligned}$$

Here, σ^2 is the combined variance of the spatial effects, and the ρ 's are mixing parameters, interpreted as the proportion of the combined spatial variance explained by each model component. Note that $\rho_\theta = 1$ reduces the spatial component to purely overdispersion, $\rho_\phi = 1$ reduces the spatial component of the model to an adjacency ICAR model for the spatial effects, and $\rho_\gamma = 1$ reduces the spatial component to a mobility ICAR model. Most importantly, if $\rho_\gamma > \rho_\phi$ then this means that the mobility data better explains variation in COVID-19 case counts than the adjacency data. As long as the spatial weights matrix and the mobility weights matrix are linearly independent, then having both spatial and mobility terms in our model present no issues with identifiability ([Rodrigues and Assunção, 2012](#)). Finally, s_γ and s_ϕ are scaling factors, such that the geometric means of $s_\gamma^{-1} \text{Var}(\gamma_i)$ and $s_\phi^{-1} \text{Var}(\phi_i)$ are both ≈ 1 for each i , meaning that γ_i^* and ϕ_i^* are the log relative risk contributions from the movement data and physical data respectively ([Sørbye and Rue, 2014](#)). Scaling is absolutely necessary in order to conduct inference on the ρ 's. We compute the scaling factors as follows

$$s = \exp\left(\frac{1}{n} \sum_{i=1}^n \log[\mathbf{Q}^-]_{ii}\right)$$

where \mathbf{Q}^- is the generalized inverse of the $n \times n$ precision matrix ([Freni-Sterrantino et al., 2018](#)). In order to scale the precision matrices of the spatial effects, the generalized inverse for sparse matrices from ([Rue et al., 2017](#)) was used. The diagonal elements, $[\mathbf{Q}^-]_{ii}$, of \mathbf{Q}^- are referred to as the *marginal variances* of the structured spatial effects, i.e $\text{var}(\phi_i) = [\mathbf{Q}_\phi^-]_{ii}$ and $\text{var}(\gamma_i) = [\mathbf{Q}_\gamma^-]_{ii}$.

As was the case with the ICAR model in (1), we can derive the full conditionals of the combined spatial effect, $\tau_i = \phi_i^* + \gamma_i^* + \theta_i^*$, for the model described in Section 2.3

$$\tau_i | \tau_{-i} \sim N\left[\frac{\sum_j \left(\frac{\rho_\phi}{s_\phi} w_{ij} + \frac{\rho_\gamma}{s_\gamma} v_{ij}\right) \tau_j}{\frac{\rho_\phi}{s_\phi} \sum_j w_{ij} + \frac{\rho_\gamma}{s_\gamma} \sum_j v_{ij} + \rho_\theta}, \frac{\sigma^2}{\frac{\rho_\phi}{s_\phi} \sum_j w_{ij} + \frac{\rho_\gamma}{s_\gamma} \sum_j v_{ij} + \rho_\theta}\right] \quad (2)$$

These full conditionals can help provide some intuition as to the mechanism by which this model provides spatial smoothing. As $\rho_\gamma \rightarrow 1$, τ_i is simply the weighted sum of the other regions, where the weights are the proportion of region i 's total movement between each other region. If $\rho_\phi \rightarrow 1$, the conditional mean of τ_i reduces to the arithmetic average of the spatial effects of its neighbours.

If $\rho_\theta \rightarrow 1$, then the conditional mean shrinks to 0 (remember that $\rho_\phi + \rho_\gamma + \rho_\theta = 1$). Given that ρ_θ is positive, the conditional mean is always shrunk towards 0, resulting in spatial smoothing. In practice, the conditional mean will be a weighted average of the estimates smoothed by the movement GMRF, the physical GMRF, and 0. It is important to note here that the w_{ij}/s_ϕ and v_{ij}/s_γ are relative measures due to the scaling factors. That is, doubling the total amount of movement has no effect on the conditional mean or variance of τ_i . This is in contrast to the combined spatial effects in the commonly used Leroux model (Leroux et al., 2000). Additionally, the variance of $\tau_i|\tau_{-i}$ is lower when region i has a lot of movement or many neighbours, relative to the other regions.

2.5. Inference, computation, and validation

Four chains each with 3000 iterations of No U-Turn Sampling were used for parameter estimation within Stan (Stan Development Team, 2021). The first 1500 iterations were used as a warm-up, the 1500 remaining iterations from each chain were thinned by a factor of 10, leaving 600 total posterior samples to perform inference. As mentioned in Section 2.2, we require $\sum_i \phi_i = 0$. In practice, we use the soft constraint

$$\sum_i \phi_i \sim N(0, 0.001)$$

for computation purposes (as recommended by the Stan team Morris et al., 2019). To complete the model, priors for β and μ were $N(0, 1)$. To ensure the robustness of our results, we also ran BYM models using the adjacency data and the movement data separately. That is, for both Madrid and Castilla-Leon, we ran a model where we assumed $\rho_\gamma = 0$, and a separate model where $\rho_\phi = 0$. The results of these four models are presented in Section 3.2.

Our code and posterior samples are posted at https://github.com/cghr-toronto/public/tree/mast/covid/spain_public_code.

3. Results

3.1. Joint model

Table 1 shows posterior medians and credible intervals for the mixing parameters for the model with both movement and adjacency spatial effects. For both Madrid and Castilla-Leon, the proportion of spatial variation explained by γ is much higher than that of ϕ and θ . The posterior probability that $\rho_\gamma > \rho_\phi$ was 0.997 for Madrid, and 0.998 for Castilla-Leon. However, ϕ does seem to account for a non-trivial amount of spatial variation in both Madrid and Castilla-Leon. This means that although movement data is likely more explanatory, adjacency data can help with explaining variation in COVID-19 cases. Additionally, there is a substantial amount of spatial variation explained by the unstructured spatial effect for Madrid. This is not the case for Castilla-Leon, as most of the mass of the posterior of ρ_θ is near 0. This makes sense given that Madrid has a large metropolitan centre surrounded by a mix of suburbs and rural areas, so there are probably spatial confounders that our model is missing. For a plot of the posterior densities of ρ , see Appendix A.

Figs. 4(a) through 4(d) show the spatial distribution γ^* and ϕ^* , plotted using the same colour scale for comparability. We can see that γ 's log-relative risks have a lot more spatial variation in both Communities. The log-relative risks for ϕ tend to have smooth spatial gradients, while γ tends to identify clusters of regions as high-risk areas. As seen in Eq. (2), the expectation of the combined spatial effects are a weighted average of these spatial effects, and 0 (notice that the numerator can be rewritten as $\sum_j (\frac{\rho_\phi}{s_\phi} w_{ij} + \frac{\rho_\gamma}{s_\gamma} v_{ij} + \rho_\theta \cdot 0) \tau_j$ where $\rho_\theta > 0$). Figs. 4(e) and 4(f) show the predicted cases per 1000 people per region, showing highly similar patterns to the observed values in Fig. 2.

The standard deviation was slightly larger for Castilla-Leon than it was for Madrid. Fig. B.2 shows the spatial distribution of the standard deviation of the cases per thousand people in both communities. Here, we can see that the standard deviation is pretty small in and around Madrid-city, because the movement to and from Madrid-city is causing a high-degree of spatial smoothing

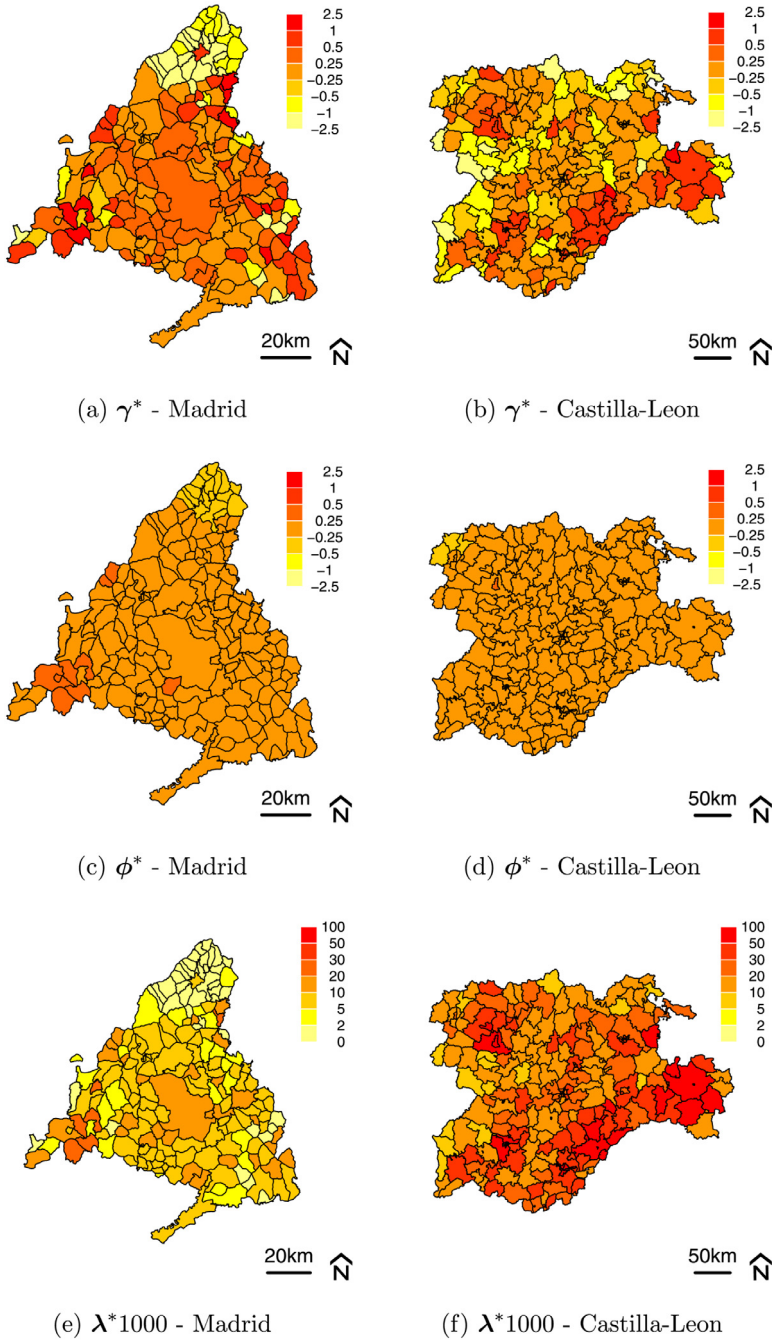


Fig. 4. Log-relative risk contributions (a–d) from the movement effects (γ^*) and spatial effects (ϕ^*). The predicted cases per thousand people are also presented (e–f).

Table 1

Posterior medians, and 95% credible intervals for ρ in BYM models using movement and physical (adjacency) data in the same model.

	Parameter	Madrid Est (95% CrI)	Castilla-Leon Est (95% CrI)
ρ	Movement	0.76 (0.54, 0.89)	0.88 (0.66, 0.98)
	Neighbour	0.13 (0.01, 0.39)	0.09 (0.01, 0.30)
	Independent	0.10 (0.02, 0.25)	0.02 (0.00, 0.09)
	μ	-5.36 (-5.51, -5.24)	-3.75 (-3.78, -3.73)
	β	0.12 (0.05, 0.20)	-0.01 (-0.04, 0.02)
	σ	0.65 (0.55, 0.78)	0.72 (0.63, 0.83)

Table 2

Posterior medians, and 95% credible intervals for ρ in BYM models using movement and physical (adjacency) data in separate models.

	Parameter	Madrid Est (95% CrI)	Castilla-Leon Est (95% CrI)
ρ	Movement	0.82 (0.66, 0.91)	0.95 (0.89, 0.98)
	Neighbour	0.56 (0.22, 0.83)	0.77 (0.58, 0.91)
μ	Movement	-5.34 (-5.48, -5.23)	-3.75 (-3.78, -3.73)
	Neighbour	-5.18 (-5.30, -5.09)	-3.74 (-3.78, -3.70)
β	Movement	0.12 (0.05, 0.18)	-0.02 (-0.05, 0.02)
	Neighbour	0.13 (0.01, 0.24)	-0.01 (-0.05, 0.04)
σ	Movement	0.63 (0.55, 0.76)	0.74 (0.65, 0.83)
	Neighbour	0.66 (0.56, 0.83)	0.58 (0.51, 0.66)

in the surrounding area. The effect of movement within regions, β , is associated with larger case counts in Madrid, but this is not the case for Castilla-Leon. This small covariate effect could result in more variance being attributable to the random effects, potentially contributing to the larger σ in Castilla-Leon.

3.2. Model validation - individual models

Table 2 shows posterior medians and credible intervals for the ρ parameter from the movement and physical BYM models described in Section 2.5, fit separately to Madrid and Castilla-Leon (four models total). In both regions, the model where spatial smoothing is induced by population movement explains a higher proportion of the variation in the outcome, indicated by the posterior density of ρ having more mass near 1. Additionally, the BYM model that used physical adjacency as a spatial smoother had a much wider credible interval for ρ , indicating more model uncertainty. Both models show more uncertainty in the region of Madrid than for Castilla-Leon, likely due to the fact that Madrid is more heterogeneous in terms of population density and other factors. For full posterior densities of the ρ parameter, see Fig. A.2.

4. Discussion

In this paper, we have demonstrated that there is much value in using mobility data in combination with geographical proximity for defining correlation structures COVID-19 incidence data. We showed that even while using only one week of movement data, we were able to explain the spatial variation in COVID-19 counts better than using the classic BYM model. Additionally, we showed that the model can be re-parametrized so that the means by which smoothing occurs in these mobility models is intuitive.

A key limitation of this work is that the models presented in this paper do not serve as individual-level infectious disease models, as correlation is induced by a latent effect rather than direct dependence between the counts. However, this will be a natural extension of this work and would

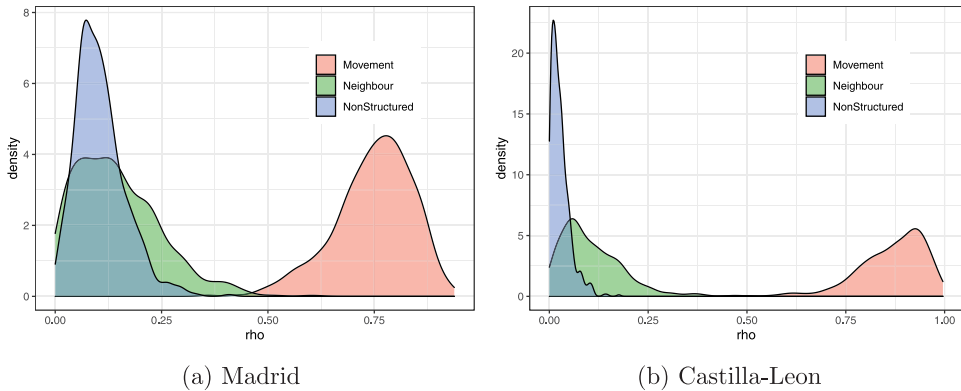


Fig. A.1. Posterior Density of the proportion of variance explained by each of the 3 spatial parameters when adjacency and movement data are included in the same model.

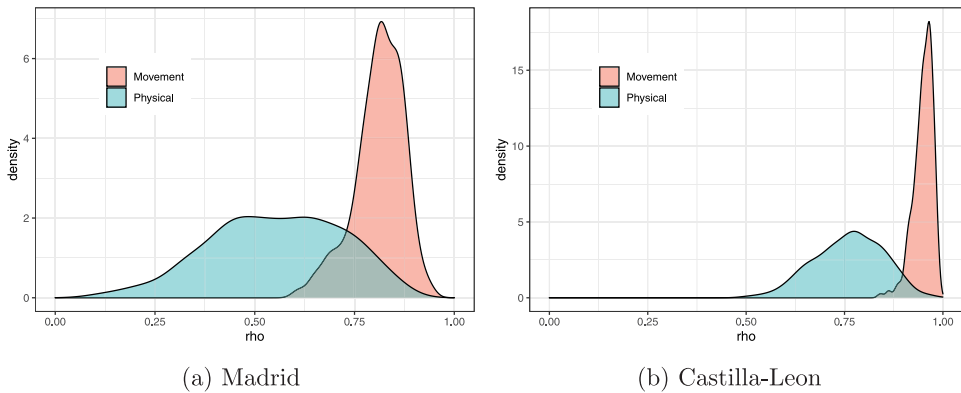
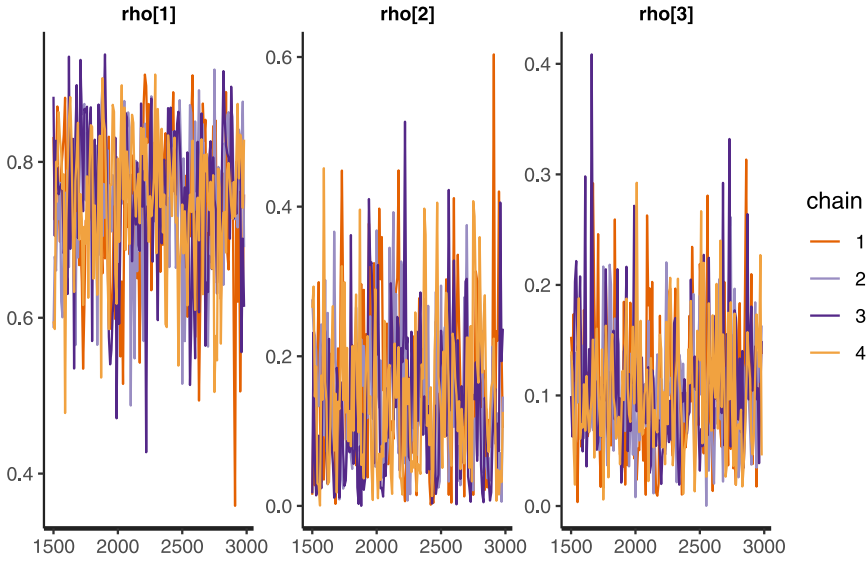


Fig. A.2. Posterior Density of the proportion of variance explained by spatial components when adjacency and movement data are used in separate models (model validation).

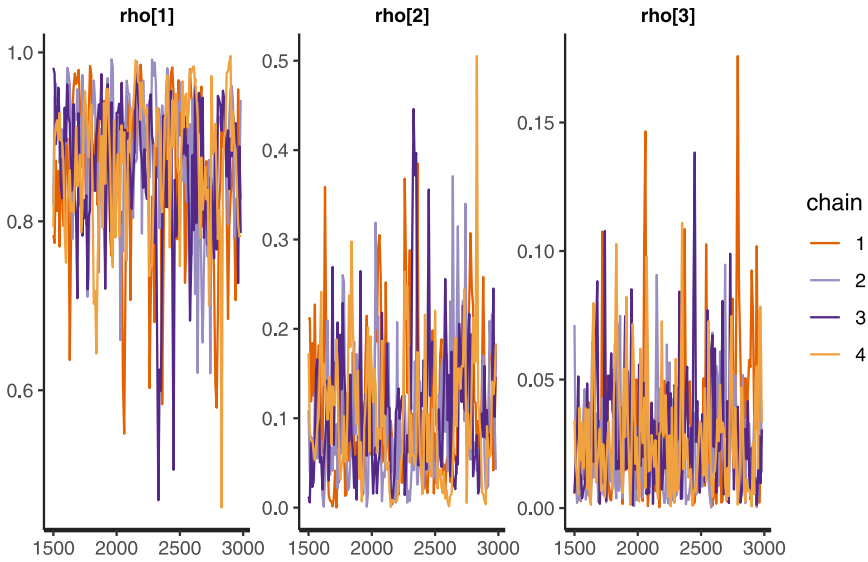
require the addition of many more parameters, including multiple mobility network components at various time points. This will ultimately pose a computational challenge as well.

An additional limitation of this work is that the availability and structure of mobility data will vary across data sources, and may only be available in higher income countries. Furthermore, there is selection bias in the movement data, as it only tracks those who actually have a cellphone, which may tend to be younger and more economically advantaged individuals. Given potential differences in quality of these data, its efficacy in spatial models may need to be assessed on a case by case basis.

Furthermore, the models presented in this paper may suffer from overfitting. A potential remedy for this would be to put a penalized complexity prior (Simpson et al., 2017) on the mixing parameters, which may improve inference by shrinking ρ_γ (and perhaps ρ_ϕ) towards 0. An interesting area for future work would be to combine Dirichlet and penalized complexity priors to specify a joint prior for the mixing parameters as described in Fuglstad et al. (2020), which can be implemented using the *makemyprior* R package (Hem et al., 2021). This was deemed unnecessary for this work, as we were mainly interested in comparing ρ_γ to ρ_ϕ , and felt that our prior should not favour either one of these terms.



(a) Madrid



(b) Castilla-Leon

Fig. A.3. Traceplots of ρ .

Despite these limitations, this work demonstrates the value of mobility data and provides the foundation for various extensions and future work. This data is only becoming more abundant as time passes, and methods that allow for efficient use of this data are essential to model the current epidemic, and any spatial epidemiological application where population movement is likely a predictor of disease.

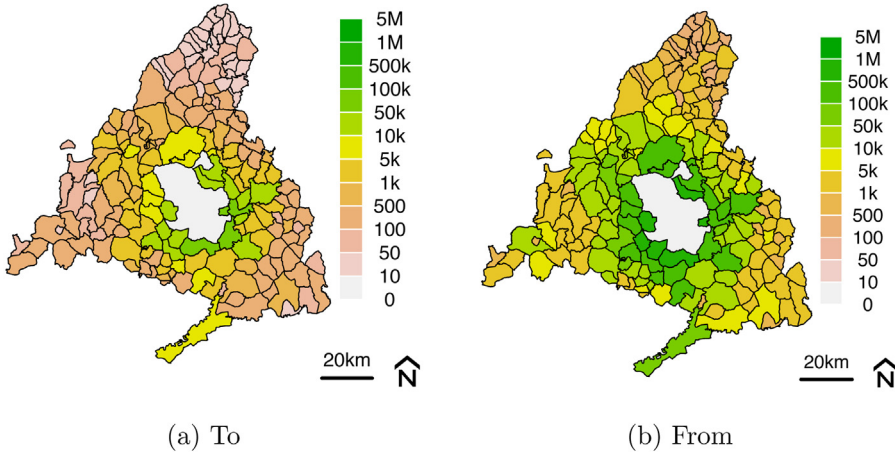


Fig. B.1. Number of trips to and from Madrid City (white).

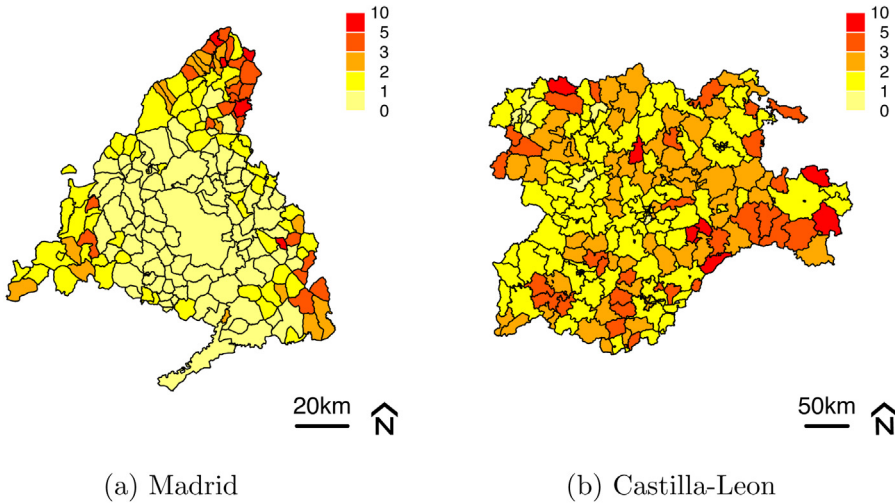


Fig. B.2. Standard deviations of predicted cases per thousand people.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Funding information

This research was funded by the Natural Sciences and Engineering Research Council (RGPIN-2017-06856) and the Ministry of Science and Innovation (PID2019-107392RB-I00).

Appendix A. Posterior densities of ρ for various models

See [Figs. A.1–A.3](#).

Appendix B. Additional spatial plots

See Figs. B.1 and B.2.

References

- Anon, 2021. Epidemiological surveillance network of madrid. <https://datos.gob.es>. (Accessed 10 Jan 2021).
- Arab-Mazar, Z., Sah, R., Rabaan, A.A., Dhama, K., Rodriguez-Morales, A.J., 2020. Mapping the incidence of the COVID-19 hotspot in Iran – implications for travellers. *Travel Med. Infect. Dis.* 34, 101630.
- Arándiga, F., Baeza, A., Cordero-Carrión, I., Donat, R., Martí, M.C., Mulet, P., Yáñez, D.F., 2020. A spatial-temporal model for the evolution of the COVID-19 pandemic in Spain including mobility. *Mathematics* 8 (10), 1677.
- Audirac, M., Tec, M., Meyers, L.A., Fox, S., Zigler, C., 2020. How timing of stay-home orders and mobility reductions impacted first-wave COVID-19 deaths in US counties. medRxiv. <https://doi.org/10.1101/2020.11.24.20238055>.
- Badr, H.S., Du, H., Marshall, M., Dong, E., Squire, M.M., Gardner, L.M., 2020. Association between mobility patterns and COVID-19 transmission in the USA: a mathematical modelling study. *Lancet Infect. Dis.* 20 (11), 1247–1254.
- Baum, C.F., Henry, M., 2020. Socioeconomic factors influencing the spatial spread of COVID-19 in the United States. Preprints with the Lancet. <http://dx.doi.org/10.2139/ssrn.3559569>.
- Besag, J., 1974. Spatial interaction and the statistical analysis of lattice systems. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 36 (2), 192–225.
- Besag, J., York, J., Mollié, A., 1991. Bayesian image restoration, with two applications in spatial statistics. *Ann. Inst. Statist. Math.* 43 (1), 1–20.
- Bilal, U., Barber, S., Tabb, L., Diez-Roux, A.V., 2020. Spatial inequities in COVID-19 testing, positivity, incidence and mortality in 3 US cities: a longitudinal ecological study. medRxiv. <https://doi.org/10.1101/2020.05.01.20087833>.
- Brainard, J.S., Rushton, S., Winters, T., Hunter, P.R., 2020. Spatial risk factors for Pillar 1 COVID-19 case counts and mortality in rural eastern England, U.K. medRxiv. <https://doi.org/10.1101/2020.12.03.20239681>.
- Briz-Redón, A., Serrano-Aroca, A., 2020. A spatio-temporal analysis for exploring the effect of temperature on COVID-19 early evolution in Spain. *Sci. Total Environ.* 728, 138811.
- Chang, S., Pierson, E., Koh, P.W., Gerardin, J., Redbird, B., Grusky, D., Leskovec, J., 2021. Mobility network models of COVID-19 explain inequities and inform reopening. *Nature* 589 (7840), 82–87.
- DiMaggio, C., Klein, M., Berry, C., Frangos, S., 2020. Blacks/African Americans are 5 times more likely to develop COVID-19: spatial modeling of New York city ZIP code-level testing results. p. 2020, medRxiv, vol. 14.
- Duncan, E.W., White, N.M., Mengersen, K., 2017. Spatial smoothing in Bayesian models: a comparison of weights matrix specifications and their impact on inference. *Int. J. Health Geogr.* 16 (1), 1–16.
- Freni-Sterrantino, A., Ventrucci, M., Rue, H., 2018. A note on intrinsic conditional autoregressive models for disconnected graphs. *Spatial Spatio-Temporal Epidemiol.* 26, 25–34.
- Fuglstad, G.-A., Hem, I.G., Knight, A., Rue, H., Riebler, A., 2020. Intuitive joint priors for variance parameters. *Bayesian Anal.* 15 (4), 1109–1137.
- Geilhufe, M., Held, L., Skråvseth, S.O., Simonsen, G.S., Godtliebsen, F., 2014. Power law approximations of movement network data for modeling infectious disease spread. *Biom. J.* 56 (3), 363–382.
- General Directorate of Information Systems, Quality and Pharmaceutical Provision, 2021. Open data of castile and leon. <https://datosabierto.jcyl.es/web/es/datos-abiertos-castilla-leon.html>. (Accessed 10 Jan 2021).
- Giuliani, D., Dickson, M.M., Espa, G., Santi, F., 2020. Modelling and predicting the spatio-temporal spread of Coronavirus Disease 2019 (COVID-19) in Italy. Preprints with the Lancet. <http://dx.doi.org/10.2139/ssrn.3559569>.
- Hem, I.G., Fuglstad, G.-A., Riebler, A., 2021. makemyprior: Intuitive construction of joint priors for variance parameters in R. arXiv preprint arXiv:2105.09712.
- Huang, G., Brown, P.E., 2021. Population-weighted exposure to air pollution and COVID-19 incidence in Germany. *Spatial Stat.* 41, 100480.
- Huang, H., Wang, Y., Wang, Z., Liang, Z., Qu, S., Ma, S., Mao, G., Liu, X., 2020. Epidemic features and control of 2019 novel coronavirus pneumonia in Wenzhou, China. Preprints with the Lancet. <http://dx.doi.org/10.2139/ssrn.3550007>.
- Iacus, S.M., Santamaría, C., Sermi, F., Spyros, S., Tarchi, D., Vespe, M., 2020. Human mobility and COVID-19 initial dynamics. *Nonlinear Dynam.* 101 (3), 1901–1919.
- Kang, D., Choi, H., Kim, J.-H., Choi, J., 2020. Spatial epidemic dynamics of the COVID-19 outbreak in China. *Int. J. Infect. Dis.* 94, 96–102.
- Kraemer, M.U., Yang, C.-H., Gutierrez, B., Wu, C.-H., Klein, B., Pigott, D.M., Du Plessis, L., Faria, N.R., Li, R., Hanage, W.P., et al., 2020. The effect of human mobility and control measures on the COVID-19 epidemic in China. *Science* 368 (6490), 493–497.
- Leroux, B.G., Lei, X., Breslow, N., 2000. Estimation of disease rates in small areas: a new mixed model for spatial dependence. In: *Statistical Models in Epidemiology, the Environment, and Clinical Trials*. Springer, pp. 179–191.
- Liu, J., Zhou, J., Yao, J., Zhang, X., Li, L., Xu, X., He, X., Wang, B., Fu, S., Niu, T., et al., 2020. Impact of meteorological factors on the COVID-19 transmission: A multi-city study in China. *Sci. Total Environ.* 726, 138513.
- Morris, M., Wheeler-Martin, K., Simpson, D., Mooney, S.J., Gelman, A., DiMaggio, C., 2019. Bayesian hierarchical spatial models: Implementing the Besag York Mollié model in Stan. *Spatial Spatio-Temporal Epidemiol.* 31, 100301.
- Orea, L., Álvarez, I.C., 2020. How effective has the spanish lockdown been to battle COVID-19? A spatial analysis of the coronavirus propagation across provinces. *Documento Trabajo FEDEA* 3, 1–33.

- Persson, J., Parie, J.F., Feuerriegel, S., 2021. Monitoring the COVID-19 epidemic with nationwide telecommunication data. arXiv preprint arXiv:2101.02521.
- Pullano, G., Valdano, E., Scarpa, N., Rubrichi, S., Colizza, V., 2020. Evaluating the effect of demographic factors, socioeconomic factors, and risk aversion on mobility during the COVID-19 epidemic in France under lockdown: a population-based study. *Lancet Digit. Health* 2 (12), e638–e649.
- Riebler, A., Sørbye, S.H., Simpson, D., Rue, H., 2016. An intuitive Bayesian spatial model for disease mapping that accounts for scaling. *Stat. Methods Med. Res.* 25 (4), 1145–1165.
- Rodrigues, E.C., Assunção, R., 2012. Bayesian spatial models with a mixture neighborhood structure. *J. Multivariate Anal.* 109, 88–102.
- Rue, H., Held, L., 2005. *Gaussian Markov Random Fields: Theory and Applications*. CRC Press.
- Rue, H., Riebler, A., Sørbye, S.H., Illian, J.B., Simpson, D.P., Lindgren, F.K., 2017. Bayesian computing with INLA: a review. *Annu. Rev. Stat. Appl.* 4, 395–421.
- Schrödle, B., Held, L., Rue, H., 2012. Assessing the impact of a movement network on the spatiotemporal spread of infectious diseases. *Biometrics* 68 (3), 736–744.
- Shi, P., Dong, Y., Yan, H., Zhao, C., Li, X., Liu, W., He, M., Tang, S., Xi, S., 2020. Impact of temperature on the dynamics of the COVID-19 outbreak in China. *Sci. Total Environ.* 728, 138890.
- Simpson, D., Rue, H., Riebler, A., Martins, T.G., Sørbye, S.H., 2017. Penalising model component complexity: A principled, practical approach to constructing priors. *Statist. Sci.* 32 (1), 1–28.
- Sørbye, S.H., Rue, H., 2014. Scaling intrinsic Gaussian Markov random field priors in spatial modelling. *Spatial Stat.* 8, 39–51.
- Stan Development Team, 2021. *Stan modeling language users guide and reference manual, version 2.26*. <https://mc-stan.org>.
- Sugg, M.M., Spaulding, T.J., Lane, S.J., Runkle, J.D., Harden, S.R., Hege, A., Iyer, L.S., 2021. Mapping community-level determinants of COVID-19 transmission in nursing homes: A multi-scale approach. *Sci. Total Environ.* 752, 141946.
- Valencia, A., 2021. COVID-19 flow maps. <https://flowmaps.life.bsc.es/flowboard/>. (Accessed 10 Jan 2021).
- Volkova, V.V., Howey, R., Savill, N.J., Woolhouse, M.E., 2010. Sheep movement networks and the transmission of infectious diseases. *PLoS One* 5 (6), e11185.