

Identifying extreme COVID-19 mortality risks in English small areas: a disease cluster approach

A. Adin · P. Congdon · G. Santafé ·
M.D. Ugarte

Received: date / Accepted: date

Abstract The COVID-19 pandemic is having a huge impact worldwide and has highlighted the extent of health inequalities between countries but also in small areas within a country. Identifying areas with high mortality is important both of public health mitigation in COVID-19 outbreaks, and of longer term efforts to tackle social inequalities in health. In this paper we consider different statistical models and an extension of a recent method to analyze COVID-19 related mortality in English small areas during the first wave of the epidemic in the first half of 2020. We seek to identify hotspots, and where they are most geographically concentrated, taking account of observed area factors as well as spatial correlation and clustering in regression residuals, while also allowing for spatial discontinuities. Results show an excess of COVID-19 mortality cases in small areas surrounding London and in other small areas in North-East and North-West of England. Models alleviating spatial confounding show ethnic isolation, air quality and area morbidity covariates having a significant and broadly similar impact on COVID-19 mortality, whereas nursing home location seems to be slightly less important.

Keywords disease mapping · ecological regression · INLA · restricted regression · smoothing

A. Adin · G. Santafé · M.D. Ugarte
Department of Statistics, Computer Science and Mathematics, Public University of Navarre,
Spain.
Tel.: +34-348-169202, Fax: +34-948-169204
E-mail: lola@unavarra.es

A. Adin · G. Santafé · M.D. Ugarte
Institute for Advanced Materials and Mathematics (INAMAT2), Public University of
Navarre, Spain.

P. Congdon
School of Geography, Queen Mary University of London, UK.

1 Introduction

The COVID-19 epidemic has highlighted the extent of disease inequalities between different small areas within countries, and identifying higher risk areas is an important aspect both of public health mitigation in infectious disease outbreaks, and of longer term efforts to tackle social inequalities in health. Research into spatial inequalities in COVID-19 incidence and mortality draws on a longer tradition of ecological research into health inequities. Ecological research examines the impact of area social and physical environments on population health, and seeks to establish areas with high disease risk (Roux, 2016; Correa-Agudelo et al, 2021; Morenoff and Lynch, 2004; Berkowitz et al, 2020). Spatial clustering in area risk factors, whether observed or unobserved, is likely to produce geographic concentrations in excess risk. For example, in a study of COVID-19 mortality in Italian municipalities the authors Ciminelli and Garcia-Mandicó (2020) find that relatively few municipalities account for a disproportionate number of deaths. An official UK study into geographic concentrations of COVID-19 mortality (Office of National Statistics (ONS), 2020a) reported that “a few areas saw COVID-19 mortality more than seven times the expected level compared with the rest of the country”. Another UK study (Kontopantelis et al, 2021) reported disproportionate concentrations of excess mortality due to COVID-19 in some regions.

Diverse methodologies have contributed to recent developments in ecological research and to assessing area health risks, including Bayesian disease mapping or BDM (Kang et al, 2016). Disease mapping uses statistical models which recognize the spatial pattern present in disease rates (e.g. geographically close areas tend to have similar disease rates) through use of random effects, and offers methods to formally identify extreme risk (Stern and Cressie, 1999). One aim of such research is to smooth erratic fluctuations in risk arising from small populations and stochastic variation in disease counts. However, these procedures may sometimes produce over-smoothing, masking distinctive features in the disease risk surface, including sharp discontinuities (Duncan and Mengersen, 2020). Refinements of the basic BDM models to counter this include the use of different neighborhood matrix specifications for spatial and spatio-temporal model fitting (Briz-Redón et al, 2021) or attempts to identifying clusters or individual areas exhibiting discontinuity (Knorr-Held and Rasser, 2000; Anderson et al, 2014; Santafe et al, 2021). Ecological regression involving analysis of health outcomes should ideally use a relatively small area scale. Pinzari et al (2018) mention that, to avoid attenuating impacts of area characteristics, “units with greater social homogeneity would be appropriate for studying the associations between unit characteristics and a given health indicator”. In a study of geographic obesity variations the authors Procter et al (2008) argue that “operating at purely a global scale, say for a whole city, will ‘average out’ small areas of high prevalence such that the mean can be deemed acceptable and the pockets of problem areas are ignored, or rather, not noticed”. Geographically disaggregated models of COVID-19 outcomes have been quite widely applied (Karmakar et al, 2021; Gaudart et al, 2021; Ciminelli

and Garcia-Mandicó, 2020). For example, Karmakar et al (2021) consider variations in COVID-19 incidence and mortality between US counties; this scale of analysis has the caveat that US counties vary considerably in population size, meaning some counties may contain considerable outcome heterogeneity within their boundaries. Jalilian and Mateu (2021) study variations in the daily number of new COVID-19 confirmed cases in first-level administrative division units from Spain, Italy and Germany. Gaudart et al (2021) consider variation in COVID-19 across 96 administrative departments in France, while Ciminelli and Garcia-Mandicó (2020) consider a sample of 1161 Italian municipalities in the seven regions most severely hit by COVID-19. Several ecological regression models have been considered to estimate COVID-19 mortality risks at various geographic scales in the UK. The analysis of excess COVID-19 mortality by Kontopantelis et al (2021) uses ten regions in England and Wales, while Travaglio et al (2021) use data for English local authorities, averaging around 200 thousand population. The latter study found higher air pollution led to large increases in COVID-19 infectivity and mortality rates after controlling for demographic factors and health-related preconditions. Some UK analyses have been at small area level: for example, Harris (2020) considered COVID-19 mortality within the London region at the level of middle super output areas (MSOAs). MSOAs are census units averaging around 8300 population across England, with a 95th percentile population of 11900. The study by Daras et al (2021) was also at MSA level, but across all of England, and found COVID-19 area vulnerability to relate to ethnic composition, poverty, prevalence of long-term health conditions, living in care homes and living in overcrowded housing. As mentioned above, risk factors such as pollution, ethnic composition and poverty have been identified as area risk factors for COVID-19 in several studies. These are likely to be spatially concentrated (for example, pollution is higher in highly urbanized areas), and so one may anticipate spatial clustering in excess risk of COVID-19 outcomes. Discontinuities may also be present, due to factors such as location of food processing plants (Food and Environment Reporting Network (FERN), 2021; Davies, 2020-27-09); particular types of institution, such as prisons (Braithwaite et al, 2021); or segmented housing patterns, such as suburban social housing estates set in mainly owner occupied areas (White, 2000).

In this paper we consider several classical disease mapping models and an extension of the clustering method named DBSC (Santafe et al, 2021) to an analysis of COVID-19 related mortality in English small areas during the first wave of the epidemic in the first half of 2020. We seek to identify high risk areas, and where they are most geographically concentrated, taking account of observed area factors (e.g. pre-existing illness, ethnic composition) via regression, as well as spatial correlation and clustering in regression residuals, while also allowing for spatial discontinuities. Identifying associations between the risk of mortality and several covariates, alleviating spatial confounding, i.e., avoiding collinearity between fixed and random effects, is also of interest.

The rest of the paper is structured as follows. Section 2 describes the methodology used to analyze the COVID-19 data. All the results are provided in Section 3. The paper ends with a discussion.

2 Methods

2.1 Predicting Relative Risk: Preliminary Regressions

We first consider conventional spatial regression of COVID-19 mortality as a first step in the analysis and to choose the best representation of the baseline spatial random effect structure. Choice of covariates and of the form of regression is important in defining regression residuals which are the input to the clustering stage used later. For area covariates, we consider results from a study of COVID-19 mortality (Congdon, 2021), across 6791 MSOAs (providing entire coverage of England), and covering March to June 2020 inclusive. This study found a measure of ethnic segregation to provide a better fitting model than one using simply the area percentage in ethnic groups. The following covariates have been included in the model as potential area-level risk factors: the Lieberman isolation index (ISOL) for measuring ethnic segregation; nursing home location (NH) to represent concentrations of frail elderly; a health deprivation and disability index (HDD) as a spatial measure of long term illness levels; and a measure of poor air quality (AIRQ). These variables are all coded in such a way as to be “positive” risk factors, with higher scores expected to be associated with higher mortality.

The COVID-19 deaths data is associated with the online article by the UK Office of National Statistics entitled “Deaths involving COVID-19 by local area and socioeconomic deprivation: deaths occurring between 1 March and 31 July 2020” (Office of National Statistics (ONS), 2020b). Data on ethnicity and nursing homes are from the UK Census, data on health deprivation are from a 2019 compendium of different types of small area deprivation (Ministry of Housing, Communities and Local Government (MHCLG), 2019), while data on air pollution are from the Access to Healthy Assets and Hazards small area indicators profile at <https://www.cdrc.ac.uk/new-update-access-to-healthy-assets-and-hazards-ahah-data-resource/> (Green et al, 2018).

Let O_i denote observed mortality data (counts of COVID-19 related deaths) in the i -th MSA during March to June 2020. The following spatial Poisson mixed model is considered

$$\begin{aligned} O_i | r_i &\sim \text{Poisson}(E_i r_i), \quad i = 1, \dots, n \\ \log(r_i) &= \alpha + \mathbf{x}_i' \boldsymbol{\beta} + \xi_i \end{aligned} \tag{1}$$

where E_i denotes the expected number of deaths. These are computed using age specific national COVID-19 mortality rates applied to MSA populations, so that impacts of population age structure on deaths are controlled for in

the analysis; α is an intercept term; r_i is the relative risk in area i (with the England wide relative risk being 1); $\mathbf{x}'_i = (x_{i1}, \dots, x_{i4})$ is the vector of standardized covariates in area i ; $\boldsymbol{\beta} = (\beta_1, \dots, \beta_4)'$ is the vector of fixed effects coefficients; and ξ_i is a spatial random effect.

The random effect is spatially structured to reflect possible geographic clustering in regression residuals, and a conditional autoregressive (CAR) prior is usually assumed (Besag et al, 1991; Lee, 2011). In CAR priors the spatial effect for area i given the spatial effects in neighbouring areas is based on the average in surrounding areas (the surrounding areas may be denoted as the neighborhood or locality of area i); estimated risks in area i are smoothed towards the locality average. Different prior distributions for the spatially structured random effect $\boldsymbol{\xi}$ have been proposed and there is still debate about which is the most effective at detecting risk variations (Lee, 2011; Riebler et al, 2016). Some variation in area disease risk may be spatially unstructured, and this motivates models allowing for unstructured heterogeneity as well as spatial clustering in risk. Thus, the long established intrinsic CAR (iCAR) prior (Besag et al, 1991) represents pure spatial dependence. However, greater flexibility may be provided by other area priors. For instance, the Leroux CAR (LCAR) prior (Leroux et al, 1999) includes a parameter λ_s , varying between 0 and 1, to represent the proportion of risk variation that is spatially structured. This LCAR prior only includes a single random effect. By contrast, the convolution CAR prior, often termed the BYM model (Besag et al, 1991), includes a random effect for unstructured heterogeneity as well as an spatially structured random effect. Additionally, a modification of the Dean et al (2001) model proposed by Riebler et al (2016), hereafter called the BYM2 model, addresses both identifiability and scaling issues of the BYM model. Scaling the model is crucial to ensure that the priors for the precision parameters have the same interpretation irrespective of the spatial graph. LCAR and BYM2 also deal with spatially structured and unstructured heterogeneity, although in a different way than the BYM model does: the LCAR model through the precision matrix, and the BYM2 model uses the covariance model instead (see for example Vicente et al, 2020 for details).

As well as the standard regression estimation used in disease mapping, we also consider restricted regression (RR) models (Reich et al, 2006). As mentioned above the form of regression is important to consider as it affects estimates of regression residuals. RR models ensure the random effects are orthogonal to the fixed effects, and so alleviate spatial confounding (see for example, Adin et al, 2021 and references therein). Spatial confounding (meaning that the spatial random effect is collinear with the observed covariates) may attenuate or bias regression coefficients on the observed covariates, and inflate the variance of the estimates of these coefficients (Prates et al, 2019). For completeness, results from a non-spatial generalized linear model (GLM) are shown also. As measures of fit, we adopt the Deviance Information Criterion (DIC) (Spiegelhalter et al, 2002), and the widely applicable information criterion (WAIC) (Watanabe, 2010). These are both lower for better fitting models. Predictive fit (cross-validation outside the sample) is measured using

two score statistics: Dawid-Sebastiani score, denoted DSS, and the logarithmic score, denoted LS (Czado et al, 2009). These are also lower for better fitting models. We initially consider a preliminary regression on the four area risk factors mentioned above, before investigating clustering in risks beyond that represented by the CAR random effect(s), for example due to risk discontinuities. All computations are made using the simplified Laplace approximation strategy of the R-INLA package (stable version 21.02.23). Regarding model hyperparameters, Normal prior distributions with mean 0 and variance equal to 1000 for fixed effects and uniform prior distributions on the positive real line for the standard deviations of the random effects have been adopted. A uniform(0,1) prior for the spatial autocorrelation parameter λ_s was also assumed.

2.2 Excess Risk Detection with Models including Cluster-level Random Effects

Estimation of spatial regressions may show relatively high proportions of risk variation due to unstructured heterogeneity. This component of variation may contain important information about risk patterns that is not captured by the smoothly varying spatial random effect. There are also likely to be irregularities in disease patterns such that smoothing towards the locality average (a central feature in the spatial random effect modelling) is not appropriate. For example, a deprived area with relatively poor health may be surrounded by relatively affluent areas. Hence, a novel density-based spatial clustering (DBSC) algorithm was proposed in Santafe et al (2021) to deal with discontinuities and to smooth noisy risks in small areas. Here, we extend the proposed methodology to the context of ecological regression models as follows. In a first stage, the DBSC algorithm is used to obtain a single clustering partition $\mathbf{C} = \{C_1, \dots, C_k\}$ of the residuals of the non-spatial generalized linear model, that is,

$$\hat{\epsilon}_i = \log(O_i/E_i) - \mathbf{x}_i' \hat{\boldsymbol{\beta}}$$

by automatically detecting clustering centers and assigning each area to its nearest cluster centroid. Clusters consist of one or more areas, and are geographically connected if consisting of two or more areas. When the number of observed cases in a given small area is zero, a small constant is added to O_i .

The DBSC algorithm has a user-given parameter ℓ that defines the neighbourhood order used to detect cluster centroids. If $\ell = 1$ is considered, the neighborhood of an area contains only its adjacent neighbors (i.e. areas that share a common border). However, greater values of ℓ can be used to consider ℓ -order neighborhoods.

Then, in a second stage, the following model that includes both small area and cluster-level spatial random effects is fitted

$$\log(r_i) = \alpha + \mathbf{x}_i' \boldsymbol{\beta} + \xi_i + \psi_{j(i)} \quad (2)$$

where $\psi_{j(i)}$ is the cluster-level spatial random effect and $j(i)$ denotes that i th area is in cluster C_j . The same spatial CAR prior distributions are adopted for both ξ_i and $\psi_{j(i)}$, and the choice depends on which one provides the best fit in the preliminary regression analysis. To achieve identifiability, the following sum-to-zero constraints are placed over the area and cluster random effects respectively, namely

$$\sum_{i=1}^n \xi_i = \sum_{j=1}^k \psi_{j(i)} = 0.$$

Based on the evidence from the preliminary regressions, and considering that the clustering stage is computed using the residuals of a non-spatial model, restricted regression will be also applied in the second stage of the DBSC algorithm to alleviate confounding between fixed effects and the combined random effect $\xi_i + \psi_{j(i)}$, i.e, to avoid collinearity between fixed and random effects, making the fixed effect orthogonal to the random effects. The code to run the new version the DBSC algorithm together with the data set and complete documentation will be available at https://github.com/spatialstatisticsupna/RR_DBSC_article.

3 Results

3.1 Preliminary Regressions

Table 1 shows fit measures under different specifications of the spatial effects in the preliminary regressions. A simple Poisson model without random effects (named GLM in the table) together with mixed Poisson models incorporating different areal random priors (iCAR, LCAR, BYM, and BYM2) have been fitted. To deal with spatial confounding the corresponding restricted regression (RR) models have been also fitted. A summary of posterior distributions for the fixed effects and model hyperparameters is presented in Table 2 and Table 3 respectively. The posterior marginal distributions of the regression coefficients estimated for each model are plotted in the Appendix (Figure A.1 to Figure A.4). There, it can be clearly seen how restricted regression alleviates spatial confounding with regard to the spatial regressions.

Examination of the spatial autocorrelation parameter λ_s for the LCAR and the BYM2 models (see Table 3) suggests that the LCAR prior may overestimate the amount of spatial dependence in the random effect (the parameter takes values between 0 and 1, with 1 representing pure spatial dependence). In contrast, the BYM2 model gives a posterior mean estimate of 0.586 for its spatial autocorrelation parameter which also weights the spatially structured and spatially unstructured variability. In addition, the BYM2 model is better supported than the LCAR by the model fit criteria (see Table 1) with a significant improvement in mean deviance and better DIC, WAIC, LS, and DSS values. The BYM convolution model provides an estimate of 0.529 when computing the ratio between the marginal variance of the spatial error and the

Table 1 Model selection criteria for models fitted with INLA. GLM indicates the fit of a Generalized Linear Model (Poisson model) without random effects. iCAR, LCAR, BYM and BYM2 indicates the fit of a Poisson mixed model incorporating the intrinsic CAR, the Leroux CAR, the BYM and the BYM2 prior respectively, to the spatial random effect. The suffix RR is added to each name when restricted regression is applied.

	\bar{D}	p_D	DIC	WAIC	LS	DSS
GLM	44421.6	5.3	44426.9	44439.9	22219.9	34661.2
iCAR	31409.5	3296.4	34705.8	34779.1	18546.5	17765.2
iCAR-RR	31409.5	3296.6	34706.0	34779.1	18546.7	17765.1
LCAR	31355.4	3339.8	34695.1	34722.2	18548.5	17727.1
LCAR-RR	31355.2	3340.2	34695.4	34722.0	18548.6	17726.8
BYM	30910.7	3593.0	34503.7	34190.2	18268.3	17407.5
BYM-RR	30910.4	3593.6	34504.0	34189.8	18268.2	17407.1
BYM2	30911.0	3592.7	34503.7	34190.4	18268.0	17407.7
BYM2-RR	30910.8	3593.9	34504.6	34190.2	18268.1	17407.2

\bar{D} : mean deviance; p_D : effective number of parameters; **DIC**: deviance information criterion; **WAIC**: Watanabe-Akaike information criterion; **LS**: logarithmic score; **DSS**: David-Sebastiani score.

total of the marginal variances (one for spatial, one for heterogeneity), which could be interpreted as an approximation to the spatial autocorrelation parameter. However, we recommend to use the λ_s parameter in a scaled BYM2 model. The BYM2 and the convolution model have similar fit measures (and both provide an improved fit over the LCAR), suggesting that unstructured heterogeneity is important in the overall pattern of COVID-19 mortality risk variation.

The regression coefficients on all models show significant effects, namely 95% credible intervals confined to positive values, so all four postulated risk factors are relevant to explaining mortality variation. Covariate values are standardized, so comparing coefficients shows which are the most important risk factors: it seems that ethnic isolation, air quality and area morbidity have a broadly similar impact, whereas nursing home location is slightly less important. As can also be seen, similar posterior distributions to those obtained with the non-spatial generalized linear model (GLM) are obtained when fitting the restricted regression models. The Appendix plots show that the restricted regression coefficient estimates are more precise (narrower 95% credible intervals). This feature of the restricted regression option (in addition to controlling for spatial confounding) may be beneficial in establishing which observed risk factors are important for explaining mortality variation. It suggests that if restricted regression is not adopted, area risk factors that are relevant to risk variations are incorrectly assessed as not significant. Table 2 also suggests that impacts of ethnic isolation and poor air quality may be attenuated when there is no control for spatial confounding. Maps with the posterior median estimates of relative risks and posterior exceedence

Table 2 Posterior mean, posterior standard deviation, and 95% credible intervals of the regression coefficients for GLM, iCAR, LCAR, BYM and BYM2 models (and the corresponding RR versions) fitted with INLA.

INLA models		mean	sd	$q_{0.025}$	median	$q_{0.975}$
Lieberson index (ISOL)	GLM	0.162	0.005	0.151	0.162	0.173
	iCAR	0.102	0.015	0.072	0.102	0.132
	LCAR	0.105	0.015	0.075	0.105	0.135
	BYM	0.113	0.013	0.087	0.113	0.139
	BYM2	0.113	0.013	0.087	0.113	0.139
	iCAR-RR	0.153	0.006	0.142	0.153	0.164
	LCAR-RR	0.153	0.006	0.142	0.153	0.164
	BYM-RR	0.153	0.006	0.142	0.153	0.164
	BYM2-RR	0.153	0.006	0.142	0.153	0.164
Nursing homes (NH)	GLM	0.087	0.004	0.079	0.087	0.095
	iCAR	0.115	0.007	0.101	0.115	0.129
	LCAR	0.114	0.007	0.100	0.114	0.128
	BYM	0.107	0.007	0.092	0.107	0.121
	BYM2	0.107	0.007	0.092	0.107	0.121
	iCAR-RR	0.100	0.004	0.092	0.100	0.108
	LCAR-RR	0.100	0.004	0.092	0.100	0.108
	BYM-RR	0.101	0.004	0.092	0.101	0.109
	BYM2-RR	0.101	0.004	0.092	0.101	0.109
Health Deprivation and Disability (HDD)	GLM	0.163	0.005	0.154	0.163	0.172
	iCAR	0.195	0.013	0.170	0.195	0.220
	LCAR	0.188	0.012	0.164	0.188	0.213
	BYM	0.192	0.012	0.169	0.192	0.216
	BYM2	0.192	0.012	0.169	0.192	0.216
	iCAR-RR	0.167	0.005	0.158	0.167	0.176
	LCAR-RR	0.167	0.005	0.158	0.167	0.176
	BYM-RR	0.166	0.005	0.157	0.166	0.175
	BYM2-RR	0.166	0.005	0.157	0.166	0.175
Air quality (AIRQ)	GLM	0.170	0.006	0.158	0.170	0.181
	iCAR	0.071	0.038	-0.003	0.071	0.145
	LCAR	0.121	0.036	0.049	0.121	0.189
	BYM	0.082	0.027	0.028	0.082	0.135
	BYM2	0.082	0.027	0.028	0.082	0.135
	iCAR-RR	0.145	0.006	0.133	0.145	0.157
	LCAR-RR	0.145	0.006	0.133	0.145	0.157
	BYM-RR	0.147	0.006	0.135	0.147	0.159
	BYM2-RR	0.147	0.006	0.135	0.147	0.159

probabilities $P(r_i > 1|\mathbf{O})$ obtained with the BYM2 model are available at https://emi-sstcdapp.unavarra.es/England_MSOA/.

3.2 Risk clustering

The model selection criteria of the preliminary regressions (see Table 1) support the BYM and BYM2 options (involving an unstructured heterogeneity

Table 3 Posterior mean, posterior standard deviation, and 95% credible intervals of the model hyperparameters for GLM, iCAR, LCAR, BYM and BYM2 models (and the corresponding restricted regression (RR) versions) fitted with INLA.

INLA models		mean	sd	$q_{0.025}$	median	$q_{0.975}$
iCAR and iCAR-RR models	τ_s	1.338	0.046	1.250	1.336	1.433
LCAR and LCAR-RR models	τ_s	1.331	0.046	1.242	1.330	1.425
	λ_s	0.936	0.039	0.839	0.948	0.981
BYM and BYM-RR models	τ_u	7.706	0.435	6.925	7.677	8.633
	τ_v	4.417	0.398	3.688	4.398	5.252
BYM2 and BYM2-RR models	τ_s	3.658	0.139	3.393	3.655	3.938
	λ_s	0.586	0.032	0.464	0.528	0.591

Table 4 Model selection criteria for DBSC models fitted with INLA considering different neighborhood orders (parameter ℓ)

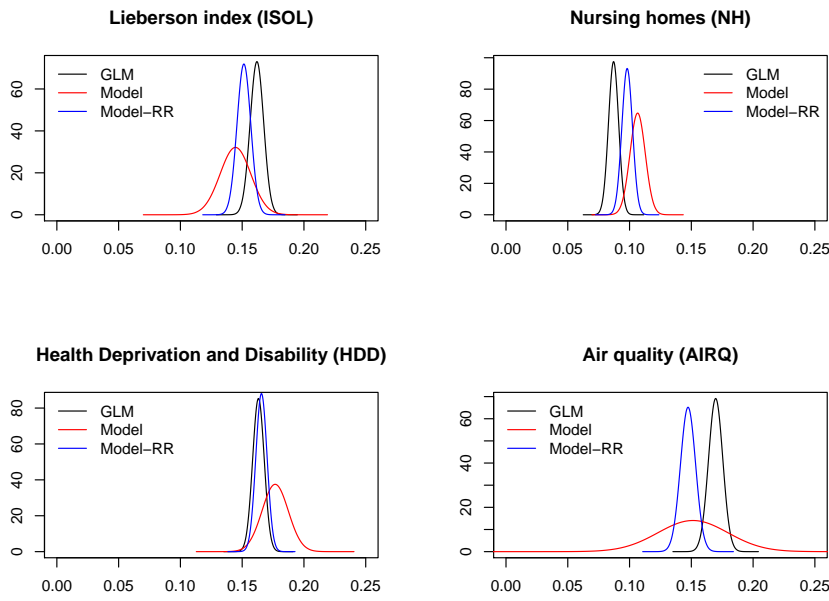
	\bar{D}	p_D	DIC	WAIC	LS	DSS
$\ell = 1$	30548.9	2417.4	32966.3	32879.5	16970.6	17222.3
$\ell = 2$	30665.7	2655.1	33320.8	33204.1	17263.6	17282.1
$\ell = 3$	30699.9	2724.8	33424.7	33297.4	17377.9	17304.5
$\ell = 4$	30733.9	2713.5	33447.4	33328.4	17374.8	17329.1
$\ell = 5$	30710.8	2829.9	33540.7	33363.9	17505.2	17296.1
$\ell = 6$	30750.0	2790.7	33540.7	33393.1	17491.0	17338.1
$\ell = 7$	30716.8	2810.1	33526.9	33379.3	17457.3	17302.4
$\ell = 8$	30769.0	2912.4	33681.3	33516.7	17636.6	17332.4

random effect, as well as a spatial effect that pools towards the locality average). The DBSC clustering aims to better elucidate the sources of this variability and also assesses when the principle of locality smoothing may need to be modified. Having identified clusters using the regression residuals of the simple Poisson model, we then apply an extended spatial regression model including cluster random effects $\psi_{j(i)}$. Finally, we also carry out restricted regression to alleviate spatial confounding. The results obtained for different values of the ℓ parameter are shown in Table 4. In general, considerably improved model fit criteria are obtained as compared to the preliminary BYM2 and BYM models which do not take account of risk clustering (compare Table 1 and Table 4). The option $\ell = 1$ is clearly preferred.

Posterior marginal distributions of the regression coefficients under the confounded and restricted regression options are shown in Table 5 and Figure 1. As before, it can be seen that restricted regression produces more precise estimates of the regression coefficients, with ethnic isolation, area morbidity, and poor air quality again figuring as the most important observed predictors of risk variability.

Table 5 Posterior mean, posterior standard deviation, and 95% credible intervals of the regression coefficients for GLM and DBSC models ($\ell = 1$) fitted with INLA.

INLA models		mean	sd	$q_{0.025}$	median	$q_{0.975}$
Lieberson index (ISOL)	GLM	0.162	0.005	0.151	0.162	0.173
	DBSC	0.145	0.012	0.120	0.154	0.169
	DBSC-RR	0.151	0.006	0.140	0.151	0.162
Nursing homes (NH)	GLM	0.087	0.004	0.079	0.087	0.095
	DBSC	0.107	0.006	0.094	0.107	0.119
	DBSC-RR	0.098	0.004	0.090	0.098	0.106
Health Deprivation and Disability (HDD)	GLM	0.163	0.005	0.154	0.163	0.172
	DBSC	0.177	0.011	0.156	0.177	0.198
	DBSC-RR	0.166	0.005	0.157	0.166	0.174
Air quality (AIRQ)	GLM	0.170	0.006	0.158	0.170	0.181
	DBSC	0.151	0.028	0.096	0.151	0.207
	DBSC-RR	0.147	0.006	0.135	0.147	0.159

**Fig. 1** Posterior marginal distributions of the regression coefficients (DBSC model).

3.3 Risk classifications

The geographic distribution of extreme relative risk is of importance for prioritizing areas for intervention and countering excess morbidity in future epidemics. To this end we use an eight-fold categorization of MSOAs according to both their urban-rural location (Office of National Statistics (ONS), 2013) and region of location (nine regions). We define extreme relative risk in the i -th MSOA using high posterior probabilities that r_i exceeds 1.5, namely the

exceedence probabilities $P(r_i > 1.5|\mathbf{O}) > 0.9$. In words, there is a high probability that the excess of risk in the i -th MSOA will be at least 50% more when compared with the global risk in England. We also consider overlapping relative risk: where risk is elevated both in an MSOA itself, and also in surrounding (adjacent) MSOAs. We define this as occurring when both relevant probabilities are high, namely when $Pr(r_i > 1|\mathbf{O}) > 0.9$, and when $Pr(R_i > 1|\mathbf{O}) > 0.9$, where R_i is the average relative risk in surrounding MSOAs (the locality of area i).

Table 6 shows the number of MSOAs with $Pr(r_i > 1.5|\mathbf{O}) > 0.9$ according to urban-rural category, under the conventional model (without a cluster random effect) as in Equation (1), and under the DBSC based regression model (which includes a cluster random effect), as in Equation (2). It can be seen that the latter produces a higher number of MSOAs with extremely high risk, especially in highly urbanized settings, namely 518 as against 396. The cluster regression also provides a classification of extreme risk that includes more deaths (11086 out of a total of 49232, or 22.5%) as against the conventional regression.

Table 7 classifies MSOAs by region. Here extremely high risk is concentrated in London, and to a lesser extent the two most Northerly regions. This feature is more clearly apparent under the clustering regression, especially in the North West region. The clustering regression also identifies more overlapping risk, again in Northern regions. High (but not necessarily extremely high) risk can be assessed on the basis of probabilities $Pr(r_i > 1|\mathbf{O}) > 0.9$. We may also consider extremely low risk areas, with relative risk below $1/1.5 = 0.667$, on the basis of the probabilities $Pr(r_i < 1/1.5|\mathbf{O}) > 0.9$, and low relative risk areas with $Pr(r_i < 1|\mathbf{O}) > 0.9$.

Table 8 shows that the conventional model classifies a much higher proportion of areas as having intermediate risk, with lower numbers of extreme high or low risk. Discrepant classifications of risk can be defined based on comparing $Pr(r_i > 1.5|\mathbf{O})$ between the cluster-adjusted and conventional models. A discrepant high risk classification is defined when $Pr(r_i > 1.5|\mathbf{O}) > 0.9$ under the DBSC model, but $Pr(r_i > 1.5|\mathbf{O}) < 0.8$ under the BYM2 model. There are 93 such areas, in which total deaths are 952, and total expected are 501.8, giving a point estimate of 1.90 for the standard mortality ratio, so meriting the classification as high risk. A discrepant low risk classification is defined when $Pr(r_i < 0.667|\mathbf{O}) > 0.9$ under the DBSC model, but $Pr(r_i < 0.667|\mathbf{O}) < 0.8$ under the BYM2 model. There are 176 such areas, with a total of 340 deaths against 683.5 expected, giving a point estimate of 0.5 for the standardized mortality ratio in these areas, and so clearly low risk areas.

Maps with the posterior median estimates of relative risks and posterior exceedence probabilities $Pr(r_i > 1|\mathbf{O})$ obtained with the $\ell = 1$ DBSC model are available at https://emi-sstcdapp.unavarra.es/England_MSOA/.

Table 6 Total MSOAs Classified as Extreme Relative Risk by Urban-Rural Category: BYM2 vs. DBSC Models.

Urban-Rural Category	Observed SMR	Number of areas with high probability $r_i > 1.5$ (BYM2)	Number of areas with high probability $r_i > 1.5$ (DBSC)	Deaths in areas with high probability $r_i > 1.5$ (BYM2)	Deaths in areas with high probability $r_i > 1.5$ (DBSC)	Total areas with high probability of overlapping risk (BYM2)	Total areas with high probability of overlapping risk (DBSC)
Urban: Major Conurbation	1.43	396	518	6471	7651	437	478
Urban: Minor Conurbation	1.16	14	22	296	399	1	2
Urban: City & Town	0.89	115	155	2395	2838	31	51
Urban: City/Town in Sparse Setting	0.47	0	0	0	0	0	0
Rural: Town & Fringe	0.70	5	7	118	130	1	2
Rural: Town & Fringe in Sparse Setting	0.51	0	0	0	0	0	0
Rural: Village & Dispersed	0.59	2	3	49	68	0	1
Rural: Village & Dispersed in Sparse Setting	0.37	0	0	0	0	0	0
All MSOAs	1.00	532	705	9329	11086	470	534

BYM2: conventional regression; DBSC: regression adjusted for discontinuity clustering.

Table 7 Total MSOAs Classified as Extreme Relative Risk by English Region: BYM2 vs DBSC Models.

Region	Observed SMR	Number of areas with high probability $r_i > 1.5$ (BYM2)	Number of areas with high probability $r_i > 1.5$ (DBSC)	Deaths in areas with high probability $r_i > 1.5$ (BYM2)	Deaths in areas with high probability $r_i > 1.5$ (DBSC)	Total areas with high probability of overlapping risk (BYM2)	Total areas with high probability of overlapping risk (DBSC)
North East	1.16	42	48	909	954	11	18
North West	1.24	89	132	1616	2118	50	76
Yorkshire-Humberside	1.01	50	67	904	1103	16	23
West Midlands	1.13	69	86	1138	1303	68	80
East Midlands	0.91	21	29	407	472	10	18
East	0.86	23	33	456	597	13	13
South East	0.84	28	40	590	717	1	1
London	1.58	204	262	3154	3637	301	304
South West	0.49	6	8	155	185	0	1
All Areas	1.00	532	705	9329	11086	470	534

BYM2: conventional regression; DBSC: regression adjusted for discontinuity clustering.

Table 8 Relative Risk Categories by Model: BYM2 vs DBSC.

	BYM2 Model	% All MSOAs	DBSC Model	% All MSOAs
Extreme High Relative Risk, $\Pr(r_i > 1.5 O) > 0.9$	532	7.8	705	10.4
Elevated (excl Extremely High) Relative Risk, $\Pr(r_i > 1 O) > 0.9$	1055	15.5	1186	17.5
Intermediate Relative Risk, $0.9 > \Pr(r_i > 1 O) > 0.1$	3166	46.6	2408	35.5
Low Relative Risk (excl Extremely Low), $\Pr(r_i < 1 O) > 0.9$	1242	18.3	1386	20.4
Extreme Low Relative Risk, $\Pr(r_i < 0.67 O) > 0.9$	796	11.7	1106	16.3
All Categories	6791	100.0	6791	100.0

4 Discussion

Delineation of high risk areas is a primary aim in disease mapping. A disease mapping model that underpredicts cases or deaths in high risk areas may lead to resourcing decisions that do not match health need. It may also be relevant to consider distinctively low risk areas, not so much on resourcing grounds, but because the location of low risk is important in assessing which environments are favorable from the viewpoint of reducing health risk. The preceding section has compared risk classifications under a conventional disease mapping model (without a clustering term) and a model including a cluster random effect to account for discrepancies in the conventional model, which can be called a clustering adjusted model. It can be seen from Table 6 that the latter produces a higher number of MSOAs with extreme risk, especially in highly urbanized settings. The clustering adjusted regression model (and hence its classifications of risk) is considerably better supported by fit measures than the conventional model. This implies, *inter alia*, that the conventional model is understating extreme relative risk in urban areas. We may also consider areas with low risk, as defined by high probabilities that the relative risk is under 1, or decisively under 1, namely below 1/1.5.

Table 8 shows that the conventional disease mapping model tends to classify a noticeably higher proportion of areas (47%) as having intermediate risk, as compared to the clustering adjusted model (36%), and to understate the numbers of definitively high risk and definitively low risk areas. This suggests over-smoothing under the conventional model. If the classification produced by conventional disease mapping was used as a basis for resource allocation, then it would tend to disadvantage areas with the highest need. An examination of areas with discrepant risk classifications shows that the conventional disease mapping model provides estimated probabilities $Pr(r_i > 1.5|\mathbf{O}) < 0.8$ for 93 areas, despite such areas having an SMR of 1.9. Similarly the conventional disease mapping model provides estimated probabilities $Pr(r_i < 0.667|\mathbf{O}) < 0.8$ for 176 areas, despite such areas having an SMR of 0.5. Examination of the discrepant risk areas suggest that the clustering adjusted model corrects for misclassification, which may occur when an area has relatively high (low) mortality as compared to its locality of surrounding areas. An unadjusted spatial smoothing mechanism (which pools to the locality average) may mean that an above average mortality area may have a relative risk estimate of below 1 if its locality has comparatively lower mortality. The DBSC clustering approach will tend to allocate such an area to a high mortality cluster to compensate for the spatial smoothing effect. This is not to discount the utility of spatial smoothing in disease mapping, but to suggest that this smoothing principle may need to be modified when there are risk discontinuities. The clustering model will also tend to adjust when an area is classified as relatively low risk on the basis of observed area risk factors (e.g. when the four predictors used to predict COVID-19 mortality are all below average), whereas other indications are of high mortality. This could be a high ratio of O_i to E_i in both the

area itself (especially when E_i is relatively high, say above 5), and also in its locality of surrounding areas.

The other methodological feature of the analysis of this paper is the benefit of comparing restricted spatial regression, which controls for spatial confounding (i.e., avoids the collinearity between fixed and random effects), with conventional disease mapping. The four area risk factors used to predict COVID-19 mortality have more precisely identified effects under restricted regression attributing all the competing explanatory effect to the covariates and considering random effects as smoothing devices. In some situations, where an area risk factor has a less clear cut effect, the gain in precision may mean the difference between deciding whether a regression effect is significant or not.

The main substantive conclusions to emerge from the analysis of this paper are a pronounced metropolitan vs rural delineation of risk, which much outweighs any North-South divide. In fact, the main difference is between London and other regions. This stands opposite to longer term health contrasts between Northern and Southern England (Buchan et al, 2017). As to establishing such contrasts, a spatial regression model incorporating a clustering stage to identify risk continuities provides a risk classification that provides a clear advantage on the basis of a range of fit measures. The conventional spatial regression, here provided by the BYM2 model of Riebler et al (2016), tends to classify a much higher proportion of areas as having intermediate risk, and is subject to apparent misclassification of some areas. The latter is exemplified by subsets of areas with clearly elevated (or depressed) risk based on standard mortality ratios, but not classified as extreme risk. Some limitations of the analysis here may be mentioned. No risk classifications are perfect, and are subject to stochastic uncertainty. Furthermore risk classifications of areas are of population aggregates, and more localized analysis may be needed to establish which sub-areas of high risk MSOAs show most adverse risk.

Acknowledgements This work has been supported by Projects MTM2017-82553-R (AEI/FEDER, UE) and PID2020-113125RB-I00 (AEI)

Ethics declarations

Conflict of interest

The authors declare that they have no conflict of interest.

Appendix

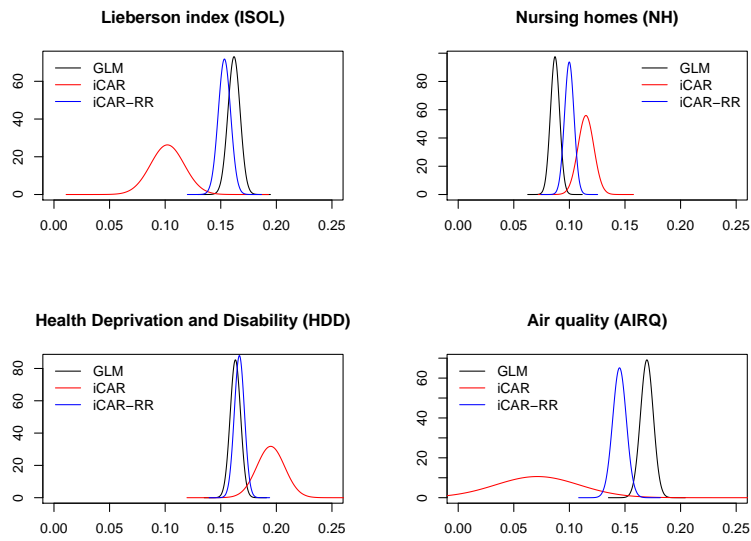


Fig. A.1 Posterior marginal distributions of the regression coefficients (iCAR model).

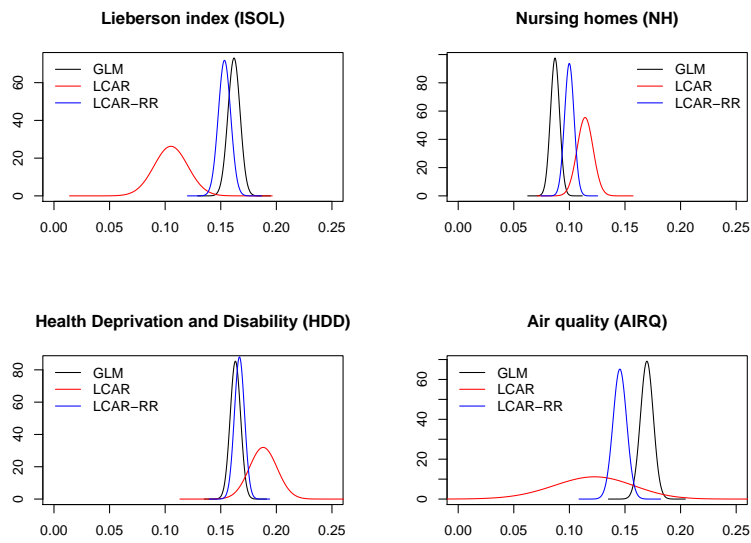


Fig. A.2 Posterior marginal distributions of the regression coefficients (LCAR model).

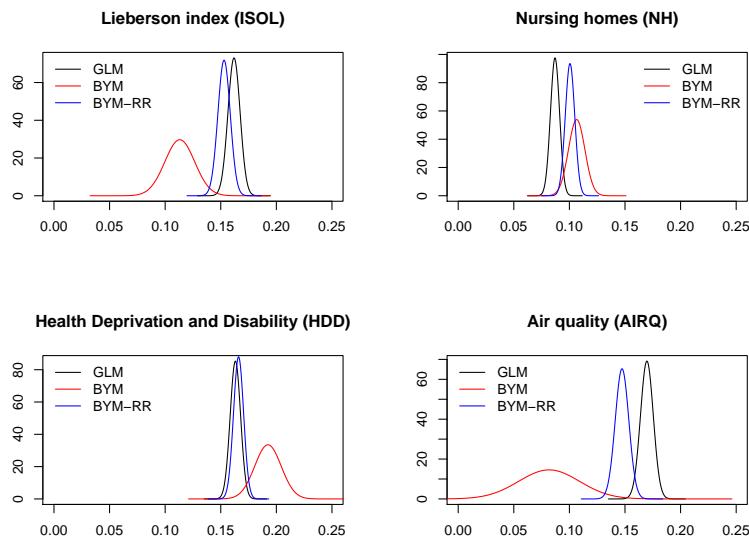


Fig. A.3 Posterior marginal distributions of the regression coefficients (BYM model).

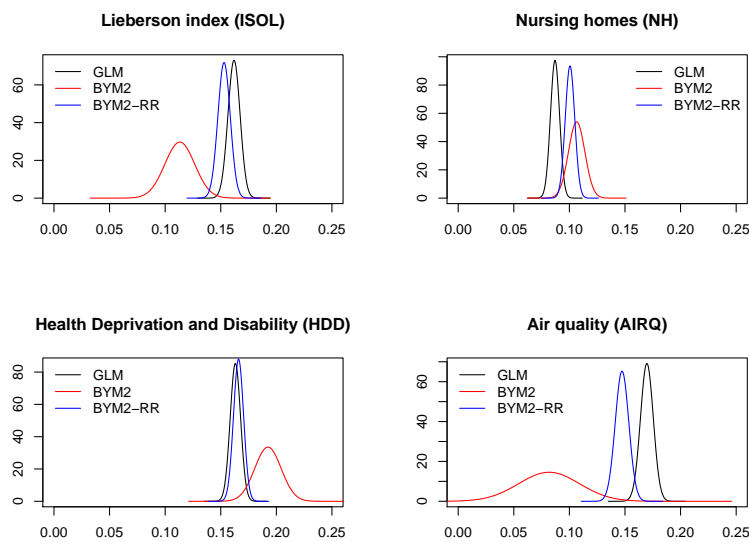


Fig. A.4 Posterior marginal distributions of the regression coefficients (BYM2 model).

References

- Adin A, Goicoa T, Hodges JS, Schnell P, Ugarte MD (2021) Alleviating confounding in spatio-temporal areal models with an application on crimes against women in India. *Statistical Modelling* p accepted, DOI <https://doi.org/10.1177/1471082X211015452>
- Anderson C, Lee D, Dean N (2014) Identifying clusters in Bayesian disease mapping. *Biostatistics* 15(3):457–469
- Berkowitz RL, Gao X, Michaels EL, Mujahid MS (2020) Structurally vulnerable neighborhood environments and racial/ethnic COVID-19 inequities. *Cities and Health* pp 1–4
- Besag J, York J, Mollié A (1991) Bayesian image restoration, with two applications in spatial statistics. *Annals of the Institute of Statistical Mathematics* 43(1):1–20
- Braithwaite I, Edge C, Lewer D, Hard J (2021) High COVID-19 death rates in prisons in England and Wales, and the need for early vaccination. *The Lancet Respiratory Medicine* 9(6):569–570
- Briz-Redón A, Iftimi A, Correcher JF, De-Andrés J, Lozano M, Romero-García C (2021) A comparison of multiple neighborhood matrix specifications for spatio-temporal model fitting: a case study on covid-19 data. *Stochastic Environmental Research and Risk Assessment* p Published online
- Buchan IE, Kontopantelis E, Sperrin M, Chandola T, Doran T (2017) North-South disparities in English mortality 1965-2015: longitudinal population study. *Journal Of Epidemiology And Community Health* 71(9):928–936
- Ciminelli G, Garcia-Mandicó S (2020) COVID-19 in Italy: an analysis of death registry data. *Journal of Public Health* 42(4):723–730
- Congdon P (2021) COVID-19 mortality in English neighborhoods: the relative role of socioeconomic and environmental factors. *J* 4(2):131–146
- Correa-Agudelo E, Mersha TB, Branscum AJ, MacKinnon NJ, Cuadros DF (2021) Identification of vulnerable populations and areas at higher risk of COVID-19-related mortality during the early stage of the epidemic in the United States. *International Journal of Environmental Research and Public Health* 18(8):4021
- Czado C, Gneiting T, Held L (2009) Predictive model assessment for count data. *Biometrics* 65(4):1254–1261
- Daras K, Alexiou A, Rose TC, Buchan I, Taylor-Robinson D, Barr B (2021) How does vulnerability to COVID-19 vary between communities in England? Developing a Small Area Vulnerability Index (SAVI). *J Epidemiol Community Health*
- Davies R (2020-27-09) Covid cases at UK food factories could be over 30 times higher than reported. *The Guardian* URL <https://www.theguardian.com/environment/2020/sep/27/covid-cases-at-food-factories-in-uk-could-be-over-30-times-higher-than-reported>
- Dean CB, Ugarte MD, Militino AF (2001) Detecting interaction between random regions and fixed age effects in disease mapping. *Biometrics* 57(1):197–202

- Duncan EW, Mengersen KL (2020) Comparing Bayesian spatial models: Goodness-of-smoothing criteria for assessing under-and over-smoothing. *PLOS ONE* 15(5):e0233,019
- Food and Environment Reporting Network (FERN) (2021) Mapping covid-19 outbreaks in the food system. <https://thefern.org/2020/04/mapping-covid-19-in-meat-and-food-processing-plants/>
- Gaudart J, Landier J, Huiart L, Legendre E, Lehot L, Bendianeand MK, Chiche L (2021) Factors associated with the spatial heterogeneity of the first wave of COVID-19 in France: a nationwide geoepidemiological study. *The Lancet Public Health* 6(4):e222–e231
- Green MA, Daras K, Davies A, Barr B, Singleton A (2018) Developing an openly accessible multi-dimensional small area index of ‘Access to Healthy Assets and Hazards’ for Great Britain, 2016. *Health and Place* 54:11–19
- Harris R (2020) Exploring the neighborhood-level correlates of covid-19 deaths in London using a difference across spatial boundaries method. *Health and Place* 66:102,446
- Jalilian A, Mateu J (2021) A hierarchical spatio-temporal model to analyze relative risk variations of covid-19: a focus on spain, italy and germany. *Stochastic Environmental Research and Risk Assessment* 35:797–812
- Kang S, Cramb S, White N, Ball S, Mengersen D (2016) Making the most of spatial information in health: A tutorial in Bayesian disease mapping for areal data. *Geospatial Health* 11(2):190–198
- Karmakar M, Lantz PM, Tipirneni R (2021) Association of social and demographic factors with COVID-19 incidence and death rates in the US. *JAMA Network Open* 4(1):e2036,462–e2036,462
- Knorr-Held L, Rasser G (2000) Bayesian detection of clusters and discontinuities in disease maps. *Biometrics* 56(1):13–21
- Kontopantelis E, Mamas MA, Deanfield J, Asaria M, Doran T (2021) Excess mortality in England and Wales during the first wave of the COVID-19 pandemic. *Journal of Epidemiology & Community Health* 75(3):213–223
- Lee D (2011) A comparison of conditional autoregressive models used in Bayesian disease mapping. *Spatial and Spatio-temporal Epidemiology* 2(2):79–89
- Leroux BG, Lei X, Breslow N (1999) Estimation of disease rates in small areas: A new mixed model for spatial dependence. In: Halloran M, Berry D (eds) *Statistical Models in Epidemiology, the Environment and Clinical Trials*, New York: Springer-Verlag, pp 179–191
- Ministry of Housing, Communities and Local Government (MHCLG) (2019) *English Indices of Deprivation 2019*. MHCLG, London, UK
- Morenoff D, Lynch J (2004) What Makes a Place Healthy? Neighborhood Influences on Racial/Ethnic Disparities in Health over the Life Course. In: Bulatao R, Cohen B (eds) *Critical Perspectives on Racial and Ethnic Differences in Health in Late Life*, National Academies Press, Washington, p Chapter 11
- Office of National Statistics (ONS) (2013) *The 2011 Rural-Urban Classification for Small Area Geographies: A User Guide and Frequently Asked Questions*

- v1.0. ONS, London, UK
- Office of National Statistics (ONS) (2020a) Analysis of geographic concentrations of COVID-19 mortality over time, England and Wales: deaths occurring between 22 February and 28 August 2020. ONS, London, UK
- Office of National Statistics (ONS) (2020b) Deaths involving COVID-19 by local area and socioeconomic deprivation: deaths occurring between 1 March and 31 July 2020. *Statistical Bulletin*. ONS, London, UK
- Pinzari L, Mazumdar S, Girosi F (2018) A framework for the identification and classification of homogeneous socioeconomic areas in the analysis of health care variation. *International Journal of Health Geographics* 17(1)
- Prates MO, Assunção RM, Rodrigues EC (2019) Alleviating spatial confounding for areal data problems by displacing the geographical centroids. *Bayesian Analysis* 14(2):623–647
- Procter KL, Clarke GP, Ransley JK, Cade J (2008) Micro-level analysis of childhood obesity, diet, physical activity, residential socioeconomic and social capital variables: where are the obesogenic environments in Leeds? *Area* 40(3):323–340
- Reich BJ, Hodges JS, Zadnik V (2006) Effects of residual smoothing on the posterior of the fixed effects in disease-mapping models. *Biometrics* 62(4):1197–1206
- Riebler A, Sørbye SH, Simpson D, Rue H (2016) An intuitive Bayesian spatial model for disease mapping that accounts for scaling. *Statistical Methods in Medical Research* 25(4):1145–1165
- Roux A (2016) Neighborhoods and health: what do we know? What should we do? *American Journal of Public Health* 106(3):430–431
- Santafe G, Adin A, Lee D, Ugarte MD (2021) Dealing with risk discontinuities to estimate cancer mortality risks when the number of small areas is large. *Statistical Methods in Medical Research* 30(1):6–21
- Spiegelhalter DJ, Best NG, Carlin BP, van der Linde A (2002) Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 64(4):583–639
- Stern H, Cressie N (1999) Inference for extremes in disease mapping. In: Lawson A, Biggeri A, Böhning D, Lesaffre E, Viel J, Bertollini R (eds) *Disease mapping and risk assessment for public health*, Chichester, UK: John Wiley, pp 63–84
- Travaglio M, Yu Y, Popovic R, Selley L, Leal NS, Martins LM (2021) Links between air pollution and COVID-19 in England. *Environmental Pollution* 268:115,859
- Vicente G, Goicoa T, Fernandez-Rasines P, Ugarte MD (2020) Crime against women in India: unveiling spatial patterns and temporal trends of dowry deaths in the districts of Uttar Pradesh. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 183(2):655–679
- Watanabe S (2010) Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research* 11(12):3571–3594

White PE (2000) Who lives in deprived areas in British cities? / Qui habite les quartiers de grande pauvreté des villes britanniques? *Géocarrefour* 75(2):107–116