

Restricted spatial regression in practice: geostatistical models, confounding, and robustness under model misspecification

Ephraim M. Hanks^{a*}, Erin M. Schliep^b, Mevin B. Hooten^{c,d} and Jennifer A. Hoeting^d

In spatial generalized linear mixed models (SGLMMs), covariates that are spatially smooth are often collinear with spatially smooth random effects. This phenomenon is known as spatial confounding and has been studied primarily in the case where the spatial support of the process being studied is discrete (e.g., areal spatial data). In this case, the most common approach suggested is restricted spatial regression (RSR) in which the spatial random effects are constrained to be orthogonal to the fixed effects. We consider spatial confounding and RSR in the geostatistical (continuous spatial support) setting. We show that RSR provides computational benefits relative to the confounded SGLMM, but that Bayesian credible intervals under RSR can be inappropriately narrow under model misspecification. We propose a posterior predictive approach to alleviating this potential problem and discuss the appropriateness of RSR in a variety of situations. We illustrate RSR and SGLMM approaches through simulation studies and an analysis of malaria frequencies in The Gambia, Africa. Copyright © 2015 John Wiley & Sons, Ltd.

Keywords: generalized linear mixed model; spatial confounding; random effects; restricted regression

1. INTRODUCTION

In the traditional spatial generalized linear mixed model (SGLMM), $\mathbf{y} = (y_1, \dots, y_n)'$ is a set of observations of the random field of interest, where y_i is the observation at spatial location \mathbf{s}_i , $i = 1, 2, \dots, n$. For a given link function $g(\cdot)$, the transformed conditional mean of \mathbf{y} is $\mathbf{z} = (z_1, \dots, z_n)' = (g(E(y_1)), \dots, g(E(y_n)))'$, and the linear predictor of the SGLMM can be written as follows:

$$\mathbf{z} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta}, \quad \boldsymbol{\eta} \sim N(\mathbf{0}, \boldsymbol{\Sigma}) \quad (1)$$

where $\boldsymbol{\beta}$ are regression parameters on covariates in \mathbf{X} , and $\boldsymbol{\eta} = (\eta_1, \dots, \eta_n)'$ is a zero-mean random effect with spatial covariance matrix $\boldsymbol{\Sigma}$.

The term 'spatial confounding' has been used to describe multicollinearity among spatial covariates \mathbf{X} and the spatial random effect $\boldsymbol{\eta}$ in (1). Hodges and Reich (2010) note that confounding can be strong enough that fixed effects $\boldsymbol{\beta}$ that are important under a non-spatial linear model may be non-significant when a spatial random effect is included. Paciorek (2010) showed that this confounding can lead to bias in estimation, especially when the spatial random effect $\boldsymbol{\eta}$ is spatially smooth and has a large effective range of spatial autocorrelation.

The most common proposed approach in recent literature to alleviate potential spatial confounding is to constrain the spatial random effect $\boldsymbol{\eta}$ to be orthogonal to the fixed effects in \mathbf{X} (Hodges and Reich, 2010; Hughes and Haran, 2013). In the restricted spatial regression (RSR) approach, the linear predictor of the SGLMM (1) is replaced by

$$\mathbf{z} = \mathbf{X}\boldsymbol{\delta} + (\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}, \quad \boldsymbol{\eta} \sim N(\mathbf{0}, \boldsymbol{\Sigma}) \quad (2)$$

where $\mathbf{P}_X \equiv \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ projects onto the space spanned by the columns of \mathbf{X} . Then $\boldsymbol{\eta}^* \equiv (\mathbf{I} - \mathbf{P}_X)\boldsymbol{\eta}$ is a spatial random effect that is orthogonal to \mathbf{X} .

* Correspondence to: Ephraim M. Hanks, Department of Statistics, The Pennsylvania State University, 326 Thomas Building, University Park, PA 16801, U.S.A. E-mail: hanks@psu.edu

a Department of Statistics, The Pennsylvania State University, PA 16801, U.S.A.

b Duke University, NC 27708, U.S.A.

c US Geological Survey, Colorado Cooperative Fish and Wildlife Research Unit, CO 80523, U.S.A.

d Colorado State University, CO 80523, U.S.A.

In the existing literature (e.g., Hodges and Reich, 2010; Hughes and Haran, 2013), RSR has been proposed as a way to recover the unbiased estimates of the regression parameters under a non-spatial (NS) model with linear predictor

$$\mathbf{z} = \mathbf{X}\boldsymbol{\delta}. \quad (3)$$

Hughes and Haran (2013) suggest a reduced rank approach to spatial analysis under the RSR framework and suggest that the Moran eigenvectors form a natural set of basis functions for spatial random effects orthogonal to the fixed effects in \mathbf{X} .

However, RSR (2) carries strong assumptions for the fixed effects in spatial models. Under the NS and RSR models, the regression parameters $\boldsymbol{\delta}$ in (2) and (3) model the *unconditional* relationship between the transformed mean \mathbf{z} of the response \mathbf{y} and the predictors in \mathbf{X} , where all variability in the direction of \mathbf{X} is explained by the linear function $\mathbf{X}\boldsymbol{\delta}$. In contrast, under the SGLMM (1), the regression parameters $\boldsymbol{\beta}$ model the *conditional* relationship between \mathbf{z} and the predictors \mathbf{X} , conditional on the spatial random effect $\boldsymbol{\eta}$. The implications of choosing among inference on the unconditional regression parameters and inference on the conditional regression parameters have not been fully explored.

We compare RSR and SGLMM modeling approaches in the presence of potential confounding between first order (mean-structure) and second order (covariance) modeling components. While spatial confounding has been the subject of considerable recent work (Reich *et al.*, 2006; Hodges and Reich, 2010; Paciorek, 2010; Hughes and Haran, 2013), we specifically consider spatial confounding in spatial models with continuous support (i.e., geostatistical models). Spatial confounding has perhaps seen less attention in the continuous spatial support literature as the original purpose of Kriging (e.g., Cressie, 1993) was for prediction. Increasingly, however, geostatistical spatial models are used to account for spatial autocorrelation in studies where the main focus is on the interpretation of the relationship between spatial covariates and a response (e.g., Diggle and Ribeiro, 2007). In this setting, spatial confounding can have a significant effect on the estimation and interpretation of regression parameters in an SGLMM.

In Section 2, we briefly review RSR for models with discrete (areal) spatial support. We then propose an approach for fitting geostatistical RSR models using Markov chain Monte Carlo and note the importance of considering the support of the process being modeled when constraining random effects in RSR. In Section 3 we consider the effects of model misspecification and show that posterior credible intervals for regression parameters under RSR can be inappropriately narrow if the true model is the SGLMM. In Section 4, we show that, conditional on the spatial random effect, RSR is a reparameterization of the SGLMM. We use this relationship between RSR and the SGLMM to propose a posterior predictive approach to appropriately expand RSR credible intervals under model misspecification. In Section 5, we conduct a detailed simulation study of RSR and SGLMM approaches for a wide variety of spatial models with Matern covariance. In Section 6, we apply RSR and SGLMM approaches to modeling malaria incidence in The Gambia, Africa. In Section 7, we conclude with a discussion of the relative advantages and disadvantages of RSR and SGLMM under different study aims.

2. SPATIAL SUPPORT AND RESTRICTED SPATIAL REGRESSION

2.1. Restricted spatial regression for models with discrete spatial support

In the existing literature, RSR has been employed almost exclusively in spatial models for areal data where observation locations occur on a lattice. In this case, a common model for spatial autocorrelation in $\boldsymbol{\eta}$ is the intrinsic conditional autoregressive (ICAR) model. Let the support of $\boldsymbol{\eta}$ be the nodes (or vertices) $\mathbf{V} = 1, 2, \dots, n$ of a graph $\mathbf{G} = (\mathbf{V}, \mathbf{A})$ with adjacency matrix \mathbf{A} , where $A_{ij} = 1$ indicates an edge exists between nodes i and j , and $A_{ij} = 0$ indicates no edge. In the standard ICAR model, the spatial precision (inverse covariance) matrix $\mathbf{Q} = \boldsymbol{\Sigma}^{-1}$ is assumed known, up to a multiplicative parameter σ^2

$$\mathbf{Q} = \frac{1}{\sigma^2} (\text{diag}(\mathbf{A}\mathbf{1}) - \mathbf{A}) \quad (4)$$

One considerable computational advantage of the ICAR model for $\boldsymbol{\eta}$ is that the precision matrix \mathbf{Q} is typically sparse and fixed up to a multiplicative constant. Restricting the spatial random effect $\boldsymbol{\eta} \sim N(\mathbf{0}, \mathbf{Q}^{-1})$ to be orthogonal to the fixed effects \mathbf{X} can be accomplished by representing the random effect $\boldsymbol{\eta}$ as a linear combination of basis vectors orthogonal to \mathbf{X} . Hughes and Haran (2013) use the Moran eigenvectors, which are the eigenvectors of $(\mathbf{I} - \mathbf{P}_\mathbf{X})\mathbf{A}(\mathbf{I} - \mathbf{P}_\mathbf{X})$. The resulting eigenvectors are orthogonal to \mathbf{X} (Boots and Tiefelsdorf, 2000) and Hughes and Haran (2013) propose a reduced rank approach to RSR by keeping only a subset of the Moran eigenvectors. Of particular note is that the Moran eigenvectors only depend on \mathbf{A} and \mathbf{X} and thus only need to be computed once if an iterative algorithm is used for model fitting or inference. Similar approaches using Moran or related spatial eigenvectors to account for spatial autocorrelation or missing spatial covariates have been used (e.g., Griffith, 2000; Griffith and Peres-Neto, 2006; Tiefelsdorf and Griffith, 2007).

2.2. Restricted spatial regression for models with continuous spatial support

When spatial locations occur in continuous space, it is common to assume that the spatial covariance matrix $\boldsymbol{\Sigma}$ is a function of the spatial locations $(\mathbf{s}_1, \dots, \mathbf{s}_n)$ and parameters $\boldsymbol{\theta}$ that govern spatial autocorrelation. A common model for continuous spatial autocorrelation is the Matern class of covariance functions in which the (i, j) -th element of $\boldsymbol{\Sigma}$ is given by

$$\Sigma_{ij} = \sigma^2 C_\nu(d_{ij}; \phi) = \sigma^2 \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\sqrt{2\nu} \frac{d_{ij}}{\phi} \right)^\nu K_\nu \left(\sqrt{2\nu} \frac{d_{ij}}{\phi} \right), \quad (5)$$

where d_{ij} is the Euclidean distance between the spatial locations of the i -th and j -th observations, σ^2 is the partial sill parameter, ν is the Matern smoothness parameter, ϕ is a range parameter, and $K_\nu(\cdot)$ is the modified Bessel function of the second kind (e.g., Cressie, 1993).

In contrast to the ICAR model (4), in which the precision matrix is fixed up to a multiplicative constant, updating the covariance parameters $\theta \equiv (\nu, \sigma^2, \phi)$ of a random effect with Matern covariance changes the covariance matrix nonlinearly. Updating θ in an iterative procedure (e.g., an MCMC sampler or numerical optimization routine) requires inverting the covariance matrix $\Sigma(\theta)$ to evaluate the likelihood. If a Moran eigenvector approach to RSR is used, a new eigenvector decomposition will need to be computed each time ϕ is updated. One possible approach to making this computationally feasible is to consider a discrete uniform prior on the range parameter ϕ (e.g., Diggle and Ribeiro Jr, 2002) and pre-compute the Moran eigenvector decomposition for each value of ϕ in the support of the prior. A reduced rank approach following that of Hughes and Haran (2013) could potentially be used for continuous spatial data by keeping only Moran eigenvectors for each value of ϕ that correspond to positive eigenvalues. However, the number of positive eigenvalues may change with ϕ , leading to variation in the dimensionality of the Moran basis.

If ϕ is allowed to have continuous support, a computationally efficient approach to implementing RSR without dimension reduction within an MCMC sampler is to constrain η to be orthogonal to the fixed effects by ‘conditioning by Kriging’ (Rue and Held, 2005). In this approach, a sample from $\eta \sim N(\mu, \Sigma)$ with the RSR constraint (4) that $\mathbf{X}'\eta = \mathbf{0}$ is obtained by first sampling from the unconstrained distribution

$$\eta^* \sim N(\mu, \Sigma) \quad (6)$$

and then applying the transformation

$$\eta = \eta^* - \Sigma \mathbf{X} (\mathbf{X}' \Sigma \mathbf{X})^{-1} \mathbf{X}' \eta^* \quad (7)$$

If \mathbf{X} has p columns representing p covariates, then the transformation (7) requires only the inversion of the $p \times p$ matrix $\mathbf{X}' \Sigma \mathbf{X}$, leading to a computationally efficient approach to implementing RSR for models with continuous spatial support.

2.3. Restricted spatial regression for partially-observed random fields

The idea of restricting the spatial random effect η to be orthogonal to fixed effects (4) is intuitive when the support of the observations y_1, \dots, y_n is identical to the spatial support of η , but care is required in the geostatistical case where the spatial field is not fully observed. As an example, consider the case where the spatial support is discrete, and the response is only observed at a set of n spatial locations $\{s_1, s_2, \dots, s_n\}$, while the full support of the random field is the set of n observed locations together with a set of m unobserved locations $(s_{n+1}, s_{n+2}, \dots, s_{n+m})$. We assume here that the covariates \mathbf{X} are observed on the full support of the random field. Under the SGLMM (1), the joint transformed mean of the response at the observed (\mathbf{z}_o) and unobserved (\mathbf{z}_u) locations is

$$\begin{pmatrix} \mathbf{z}_o \\ \mathbf{z}_u \end{pmatrix} = \begin{pmatrix} \mathbf{X}_o \\ \mathbf{X}_u \end{pmatrix} \beta + \begin{pmatrix} \eta_o \\ \eta_u \end{pmatrix}, \quad \eta \equiv \begin{pmatrix} \eta_o \\ \eta_u \end{pmatrix} \sim N(\mathbf{0}, \Sigma) \quad (8)$$

The most natural approach to RSR in this case would be to restrict the random effect $\eta \equiv (\eta_o' \eta_u')'$ to be orthogonal to $\mathbf{X} \equiv (\mathbf{X}_o' \mathbf{X}_u')'$ on the full spatial support of the system, so that

$$\mathbf{X}' \eta = \mathbf{0}. \quad (9)$$

Model fitting in the case of partially observed discrete spatial support could be accomplished using the approach of Hughes and Haran (2013) by computing the Moran eigenvectors of the graph of the full spatial support. However, the resulting regression parameters would no longer be analogous to the NS (3) parameter estimates, as the restriction in (9) no longer requires the spatial random effect η_o to be orthogonal to the fixed effects \mathbf{X}_o on the support of the observations.

If the spatial support is continuous rather than discrete, then both $X(\mathbf{s})$ and $\eta(\mathbf{s})$ are defined continuously on $\mathbf{s} \in \mathbb{S}$, although typically, observations of the process of interest are only available at a discrete set of locations. In this case, the RSR restriction could be defined by an inner product

$$\int_{\mathbb{S}} X(\mathbf{s}) \eta(\mathbf{s}) d\mathbf{s} = 0. \quad (10)$$

In practice, this restriction could be approximated by considering both $X(\mathbf{s})$ and $\eta(\mathbf{s})$ on a regular (or non-regular) grid of fine spatial resolution and approximating the integral with an appropriate summation using quadrature weights. If conditioning by Kriging (6)–(7) is used, the RSR restriction does not provide a significant computational hurdle. However, approximating the constraint (10) requires inference on η on a fine grid regardless of the dimensionality of the observations, which can be computationally intensive. Additionally, (10) requires knowledge of the spatial covariates $\mathbf{X}(\mathbf{s})$ continuously (or on a fine grid) across the study region. If MCMC is used for inference, computationally efficient approaches for conditional simulation of η (and \mathbf{X} , if needed) on a fine grid include Gaussian Markov random field (GMRF) approximations to the Gaussian field η (Lindgren *et al.*, 2011) reduced rank approximations (Cressie and Johannesson, 2008), such as predictive process approaches (Banerjee *et al.*, 2008), or fast spectral methods (Nychka *et al.*, 2002; Wikle, 2002; Hooten *et al.*, 2003; Royle and Wikle, 2005) for sampling random fields.

As an alternative to requiring orthogonality on the full support as in (10), we could consider restricting the spatial random effect to be orthogonal to the fixed effects only at the observation locations. Under this approach, the constraint is given by

$$\mathbf{X}_o' \eta_o = \mathbf{0} \quad (11)$$

thus relaxing the constraint on η_u in (8). This approach is appealing because it only requires observation of covariates \mathbf{X}_o at the spatial locations of the observations \mathbf{y}_o and because inference only requires inversion of matrices with dimension equal to the number of observed spatial locations. However, this approach ignores the full support of the spatial field, and the interpretation of the regression parameters δ is less clear. In the case of the SGLMM (1), the regression parameters β model the conditional linear relationship between \mathbf{X} and \mathbf{y} , conditioned on the full spatial random effect η . Under the restriction (9) on the full spatial support, the regression parameters δ model the unconditional linear relationship between \mathbf{X} and \mathbf{z} on the full spatial support. Under the restriction (11), the regression parameters model the unconditional linear relationship between \mathbf{X}_o and \mathbf{y}_o at the observed spatial locations.

3. COLLINEARITY OF FIXED AND RANDOM EFFECTS

The RSR (2) has been proposed as an approach to recover the point estimates of regression parameters in an NS (3) model. However, under model misspecification, inference for the unconditional regression parameters δ can lead to spurious conclusions about the importance of the fixed effects in the model. When the spatial random effect η is collinear with one or more covariates, the potential confounding between them can result in unidentifiability (or weak identifiability) of the relationship between the covariates and the response \mathbf{y} (see e.g., Paciorek, 2010, Section 2.1). For example, consider a Gaussian spatial linear mixed model (SLMM) with one fixed effect \mathbf{x} that has no relationship with the response \mathbf{y} , so $\beta = 0$ and

$$\mathbf{y} = \mu \mathbf{1} + 0\mathbf{x} + \eta + \epsilon, \quad \epsilon \sim N(\mathbf{0}, \tau^2 \mathbf{I}), \quad \eta \sim N(\mathbf{0}, \Sigma_\eta). \quad (12)$$

In (12), the response \mathbf{y} is independent of \mathbf{x} but is dependent on η , which is an unobserved spatially smooth covariate modeled by a Gaussian process. If \mathbf{x} and η are collinear, then it will be easy to misspecify the model. For example, if either an NS linear model

$$\mathbf{y} = \mu \mathbf{1} + \beta \mathbf{x} + \epsilon, \quad \epsilon \sim N(\mathbf{0}, \tau^2 \mathbf{I}) \quad (13)$$

or an RSR model

$$\mathbf{y} = \mu \mathbf{1} + \beta \mathbf{x} + (\mathbf{I} - \mathbf{P}_x)\eta + \epsilon, \quad \epsilon \sim N(\mathbf{0}, \tau^2 \mathbf{I}), \quad \eta \sim N(\mathbf{0}, \Sigma_\eta) \quad (14)$$

were used, parameter estimates could indicate a significant relationship between \mathbf{x} and \mathbf{y} through a non-zero estimate of β . We formalize this through an examination of the coefficient of determination (R^2).

Consider the case where \mathbf{x} and η are independent mean-zero Gaussian random fields with Matern covariance (5)

$$\mathbf{x} \sim N(\mathbf{0}, \Sigma_x), \quad \eta \sim N(\mathbf{0}, \Sigma_\eta) \quad (15)$$

$$\Sigma_{x;ij} \equiv \sigma_x^2 C_v(d_{ij}; \phi_x) \quad \Sigma_{\eta;ij} \equiv \sigma_\eta^2 C_\eta(d_{ij}; \phi_\eta)$$

The random fields \mathbf{x} and η are by definition independent, but we show in the supporting information, Section S1, that the mean of the coefficient of determination (R^2) between \mathbf{x} and η is approximately

$$E[R^2] \approx \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n [C_{v_x}(d_{ij}, \phi_x) \cdot C_{v_\eta}(d_{ij}, \phi_\eta)] \quad (16)$$

under a linear approximation and the assumption that both \mathbf{x} and η are zero-centered. For a fixed distance d , the Matern correlation function $C_v(d, \phi)$ approaches zero as $\phi \rightarrow \infty$ and approaches 1 as $\phi \rightarrow 0$, leading to the following properties of the linear approximation (16) to the coefficient of determination between two Matern random fields:

1. As $\phi_x \rightarrow 0$, $E[R^2] \rightarrow 0$.
2. As $\phi_\eta \rightarrow 0$, $E[R^2] \rightarrow 0$.
3. As $\phi_x \rightarrow \infty$ and $\phi_\eta \rightarrow \infty$, $E[R^2] \rightarrow 1$.

These results show that even when η and \mathbf{x} are independent, if both are spatially smooth, then η and \mathbf{x} are likely to be collinear and unconditional inference on β through the NS or RSR models could lead to spurious conclusions about the importance of \mathbf{x} in the model. In Section 5 and Section 6, we verify, through simulation, that posterior credible intervals for RSR regression parameters can be inappropriately narrow under model misspecification.

Reich *et al.* (2006) note that collinearity in fixed and random effects can lead to variance inflation and suggest that RSR alleviates this inflation because of confounding. However, inference on regression parameters under RSR is strongly related to inference on regression parameters in the NS model (which also assumes that all random effects are orthogonal to the fixed effects). In the supporting information, Section S2, we consider specifying a uniform prior for δ and taking an empirical Bayes approach to estimate the marginal variance σ^2 under the NS model and the non-spatial (nugget) variance σ_{nug}^2 under RSR. In this situation, the posterior mean of $\delta|\mathbf{y}$ under NS is the OLS estimate $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$, and the posterior variance is $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$. Under RSR, the posterior mean is also $(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$, and the posterior variance is $\sigma_{\text{nug}}^2(\mathbf{X}'\mathbf{X})^{-1}$, which is a scaled version of the NS estimate. Under an REML approach, the estimates of marginal (under NS) and nugget (under RSR) variances satisfy $\hat{\sigma}_{\text{nug}}^2 \leq \hat{\sigma}^2$, and the resulting marginal posterior variances of regression parameters under RSR will be smaller than marginal posterior variances under the NS model. Thus, when the data are Gaussian, posterior variances for regression parameters under

RSR can be smaller than corresponding posterior variances under a NS model, and RSR and NS models will often provide similar inference on regression parameters as both constrain all variation in the space spanned by \mathbf{X} to be explained by the fixed effects.

In the next section, we propose a posterior predictive approach to inference on regression parameters that maintains the RSR posterior mean but appropriately inflates the posterior variance to reflect uncertainty because of potential confounding with random effects.

For non-Gaussian data, a nonlinear link function will induce some dependence between the fixed effects and the RSR random effect, even though they are orthogonal. This dependence can result in differences in inference on the regression parameters under NS or RSR models. This discrepancy in inference is evident in our data example in Section 6.

4. POSTERIOR PREDICTIVE INFERENCE ON REGRESSION PARAMETERS

Paciorek (2010) noted that, without prior assumptions or constraints on the spatial random effect η , the SLMM model may be unidentifiable. We show that the RSR model (2) can be seen as a reparameterization of the SGLMM (1) when conditioning on the spatial random effect η . The linear predictor can be written as follows:

$$\begin{aligned} \mathbf{z} &= \mathbf{X}\boldsymbol{\beta} + \eta \\ &= \mathbf{X}\boldsymbol{\beta} + \mathbf{P}_x\eta + (\mathbf{I} - \mathbf{P}_x)\eta \\ &= \mathbf{X}\boldsymbol{\beta} + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\eta + (\mathbf{I} - \mathbf{P}_x)\eta \\ &= \mathbf{X}\left[\boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\eta\right] + (\mathbf{I} - \mathbf{P}_x)\eta \\ &= \mathbf{X}\boldsymbol{\delta} + (\mathbf{I} - \mathbf{P}_x)\eta. \end{aligned}$$

Thus, given data $\{\mathbf{y}, \mathbf{X}\}$, the conditional likelihood $L_{\text{SGLMM}}(\mathbf{y}; \boldsymbol{\beta}, \eta)$ of the data under (1) is identical to the likelihood $L_{\text{RSR}}(\mathbf{y}; \boldsymbol{\delta}, \eta)$ under (2), where

$$\boldsymbol{\delta} \equiv \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\eta. \quad (17)$$

It would thus be impossible to choose between the SLMM and RSLMM models based on in-sample data alone. This has important implications for studies in which understanding the fixed effects has a high priority, as $\boldsymbol{\beta}$ in (1) and $\boldsymbol{\delta}$ in (2) may lead to vastly different interpretations of the effect of covariates in \mathbf{X} on the response \mathbf{y} .

As SGLMM and RSR are reparameterizations of each other, it is possible to obtain inference on both the conditional ($\boldsymbol{\beta}$) and unconditional ($\boldsymbol{\delta}$) regression parameters jointly through posterior prediction (Gelman *et al.*, 2004, p.8) of derived quantities like (17). When fitting the SGLMM model (1) to the data using MCMC, at the k -th iteration of the MCMC sampler, after sampling $\boldsymbol{\beta}^{(k)}$ and $\eta^{(k)}$ from their respective full-conditional distributions, a sample from the posterior predictive distribution of the marginal regression parameters $\boldsymbol{\delta}$ can be directly obtained via evaluating (17) using the current state of the Markov chain

$$\boldsymbol{\delta}^{(k)} = \boldsymbol{\beta}^{(k)} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\eta^{(k)},$$

and posterior predictive inference on $\boldsymbol{\delta}|\mathbf{y}$ can be obtained using the resulting samples.

When fitting the RSR model to the data, samples of η are not obtained in the MCMC algorithm; rather, samples of $(\mathbf{I} - \mathbf{P}_x)\eta$ are obtained through conditioning by Kriging (6)–(7). As a result, we cannot use (17) to directly obtain posterior predictive samples of $\boldsymbol{\beta}$ under RSR. As an alternative, we propose posterior predictive inference on the derived quantity

$$\tilde{\boldsymbol{\beta}} = \boldsymbol{\delta} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\tilde{\eta}, \quad \tilde{\eta} \sim N(\mathbf{0}, \boldsymbol{\Sigma}), \quad (18)$$

which retains the mean value of the unconditional regression parameter but adjusts the variance to reflect the possible collinearity between fixed and random effects. Samples from the posterior predictive distribution of $\tilde{\boldsymbol{\beta}}$ can be obtained by conditioning on $\boldsymbol{\delta}^{(k)}$ and $\boldsymbol{\Sigma}^{(k)}$;

$$\tilde{\boldsymbol{\beta}}^{(k)} \sim N\left(\boldsymbol{\delta}^{(k)}, (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Sigma}^{(k)}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\right). \quad (19)$$

Our goal in proposing the use of posterior predictive inference on $\tilde{\boldsymbol{\beta}}$ in RSR is to retain the mean value of the unconditional regression parameter $\boldsymbol{\delta}$ from an RSR analysis and also retain the computational benefits of RSR, while adjusting the posterior variance to reflect the possible collinearity between fixed and random effects.

5. SIMULATION STUDY

In comparing SGLMM and RSR approaches with spatial modeling, we have described situations where posterior distributions of the unconditional regression parameters $\boldsymbol{\delta}|\mathbf{y}$ under RSR have the potential to be inappropriately narrow. Under a frequentist framework, this would result in an elevated rate of Type-I errors in which $\boldsymbol{\delta}$ was incorrectly found to be significantly different from zero. A Bayesian analogue to the Type-I error is the Type-S error of Gelman and Tuerlinckx (2000). If a regression parameter (β_k or δ_k) is truly zero, then we will consider a Type-S error to have occurred if a 95% equal-tailed posterior credible interval for the regression parameter does not overlap zero. We conducted a simulation study to examine Type-S errors for RSR and SGLMM under a variety of scenarios, including model misspecification and different combinations of variance parameters.

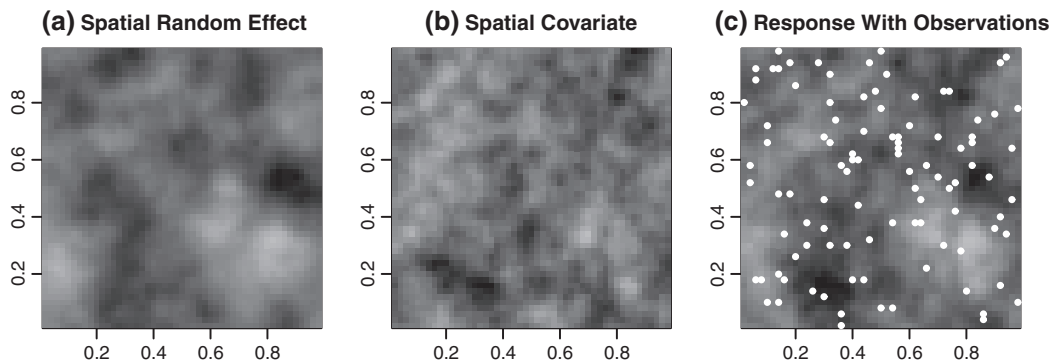


Figure 1. Simulation study to examine Type-S error rate in restricted spatial regression model. A spatial random effect η (a) and a covariate \mathbf{x} (b) were simulated on a fine grid. The resulting response \mathbf{y} (c) was simulated. To simulate partial observation of the spatial random field \mathbf{y} , 100 locations (white circles) in the unit square were randomly chosen. Varying levels of Matern smoothness parameter ν , spatial range ρ_η of η , and spatial range ρ_x of \mathbf{x} were used in different simulations (this example realization was generated with $\nu = 3/2$, $\rho_\eta = 0.4$, and $\rho_x = 0.1$)

5.1. Simulation study 1: effect of covariance parameters on Type-S error rates

In simulation, we first examined the effects of spatial structure on Type-S error rates in continuous space models. In Section 3 and in the supporting information, Section S1, we showed that the potential for collinearity between independent fixed and random effects increases as the range of spatial autocorrelation increases. To examine this behavior, we randomly chose 100 point locations from $(0, 1) \times (0, 1)$ (Figure 1) and simulated both \mathbf{x} and η as independent mean-zero Gaussian random fields with Matern covariance as in (5).

We then simulated the response variable \mathbf{y} using the SLMM (12) as the true model with model parameters given by $\beta = 0$, $\mu = -1.5$, $\sigma_\eta^2 = 1$, and $\tau^2 = 0.1$. We considered three values of the Matern smoothness parameter for both \mathbf{x} and η ($\nu = \nu_x = \nu_\eta \in \{0.5, 1.5, 2.5\}$), and four values of the range parameters ϕ_x and ϕ_η chosen so that the range ρ of spatial dependence (defined as the distance ρ at which $C_\nu(\rho, \phi) = 0.1$) varied from $\rho = 0.02$ (little spatial dependence) to $\rho = 0.40$ (moderate spatial dependence). These values of ρ were calculated using the empirical relationship $\rho = \phi\sqrt{8\nu}$ derived by Lindgren *et al.* (2011).

For each combination of smoothness parameter ν (shared by both \mathbf{x} and η), spatial range ρ_x of \mathbf{x} and spatial range ρ_η of η , we repeated this simulation process 100 times. For each resulting data set, both SGLMM (1) and RSR (2) models were fit using MCMC, with Gibbs updates for all parameters except for the range parameter ϕ_η , for which we used a Metropolis–Hastings step tuned so that acceptance was roughly 40%. Each MCMC sampler was run until the Monte Carlo standard error (e.g., Flegal *et al.*, 2008) was smaller than 0.02 for all model parameters. We considered a Type-S error to have occurred if a 95% equal-tailed posterior credible interval for the marginal posterior of $\beta|\mathbf{y}$ (for SLM) or $\delta|\mathbf{y}$ (for RSR) did not overlap zero. For this simulation study, we only considered RSR with the constraint that the spatial random effect be orthogonal to the fixed effects on the support of the data (11).

We note that, while (17) could be used to estimate the conditional regression parameters while fitting the SGLMM, the results we report for RSR are from a separate MCMC sampler. We have also estimated the conditional (RSR) regression parameters using (17) and found that the results agree with the RSR results in our simulation study (results not shown).

Figure 2 shows the estimated Type-S error rates for RSR and SGLMM in this simulation where the true model is (12). Additional numerical results are shown in the supporting information, Section S3. Over all simulations, the Type-S error rate for the SLM model was 0.051, while the Type-S error rate for the RSR model was 0.218. The Type-S error rate for RSR increases in general as the range ρ_η of the spatial random effect increases but is not affected strongly by the range ρ_x of the covariate \mathbf{x} or by the Matern smoothness parameter ν .

5.2. Simulation study 2: Type-S error rates under model Misspecification

We next considered the effects of model misspecification on Type-S error rates in continuous space spatial models. As our first simulation study showed that the Matern smoothness parameter ν had little to no impact on Type-S error rates, we consider only $\nu = 0.5$ for this simulation study. As in our first simulation study, we fixed $\beta = \delta = 0$, $\mu = -1.5$, $\sigma_\eta^2 = 1$, and considered four values of the range parameters ϕ_x and ϕ_η chosen so that the range ρ varied from $\rho = 0.02$ to $\rho = 0.40$. We fixed $\tau^2 = 0.4$ (we specified $\tau^2 = 0.1$ in our first simulation study) to examine the effect of increased noise on Type-S error rates. We simulated 100 observation locations again randomly from the unit square and simulated η , \mathbf{x} , and \mathbf{y} for each of three models: the SGLMM (12), an analogous RSR model (14) with η constrained to be orthogonal to \mathbf{x} on the support of the observations (11), and an RSR model with the constraint that η be orthogonal to \mathbf{x} on the entire unit square (10). In this last case, we approximated the integral constraint (10) with the constraint that \mathbf{x} and η be orthogonal on a 50×50 regular grid on the unit square. We then fit both SGLMM and RSR models to the resulting three data sets.

This process was repeated 100 times, and the overall Type-S error rates are shown in Figure 3 (additional numerical results are shown in the supporting information, Section S3). RSR shows inflated Type-S error rates under model misspecification when the true model is either the SGLMM or RSR with constraint on the full support. This holds true whenever $\rho_\eta \geq 0.1$. Type-S error rates for SGLMM are always near or below the nominal rate of 0.05 under model misspecification. The increased marginal variance in the second simulation study reduced the Type-S error rate in general, but the patterns of Type-S error rates for RSR under model misspecification are similar with less noise (Figure 2) and more noise (Figure 3).

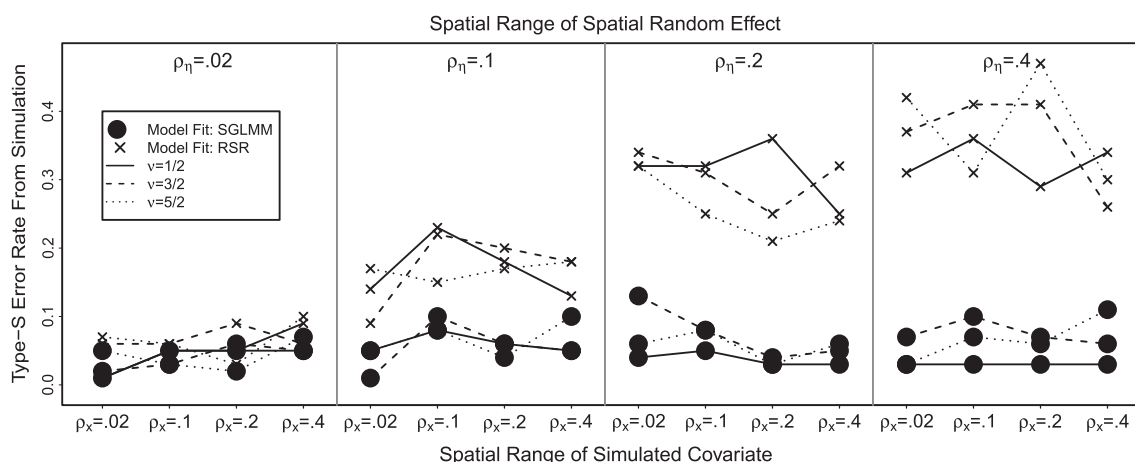


Figure 2. Simulation study to examine Type-S error rates of spatial generalized linear mixed model (SGLMM) and restricted spatial regression approaches for continuous space models when the true model is an SGLMM. Both a spatial random effect η and a covariate \mathbf{x} were simulated at 100 randomly chosen locations in the unit square. Varying levels of Matern smoothness parameter ν , spatial range ρ_η of η and spatial range ρ_x of \mathbf{x} were used. Collinearity between independent \mathbf{x} and η can lead to increased Type-S error rates of restricted spatial regression under model misspecification, with increased range of the spatial random effect η leading to increased Type-S error rates. This potential collinearity does not increase the Type-S error rate under the SGLMM approach for the simulations we considered. Additional numerical results are shown in the supporting information, Section S3

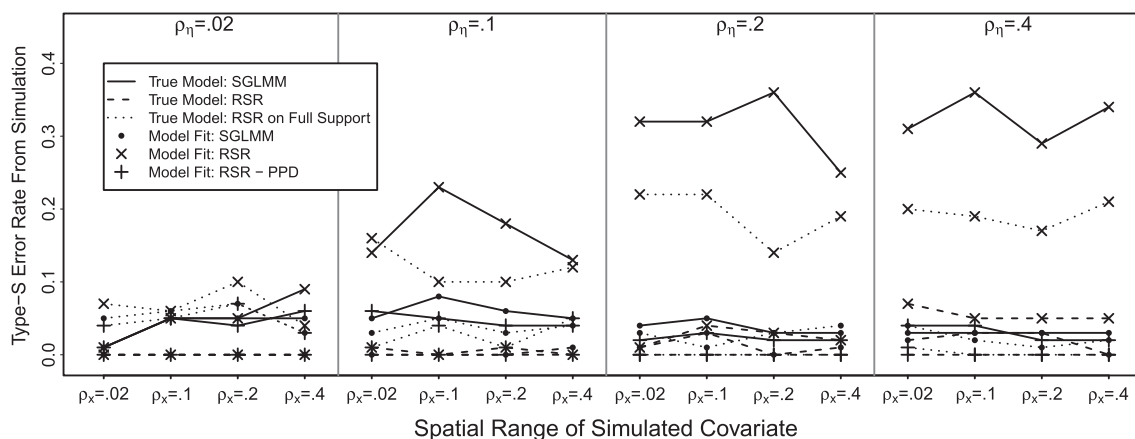


Figure 3. Simulation study to examine Type-S error rates of the spatial generalized linear mixed model (SGLMM) and restricted spatial regression (RSR) approaches for continuous spatial support models under a variety of true models used for simulation. RSR shows inflated Type-S error rates under model misspecification (when the true model is the SGLMM or RSR with restriction on the full spatial support), while the SGLMM and our proposed posterior predictive approach (RSR-PPD) are conservative under model misspecification. Additional numerical results are shown in the supporting information, Section S3

5.3. Simulation study discussion

Our simulation studies indicate that RSR has an increased risk of Type-S error for regression parameters under model misspecification for models with continuous spatial support. Conversely, Type-S error rates under the SGLMM when RSR is the generating model can be below then nominal rate of 0.05. Additional research is needed to examine whether these low Type-S error rates correspond to a lack of power under model misspecification. Type-S error rates for SGLMM and for the PPD approach (19) are always near or below the nominal rate. An additional application specific simulation study in Section 6.2 confirms these results in a binary regression setting with multiple predictor variables. We also conducted a similar simulation study for areal spatial data with ICAR covariance with analogous results (results not shown).

These results indicate that posterior credible intervals obtained under RSR are likely to be inappropriately narrow under model misspecification. Care should be taken in making statements about the importance or statistical significance of variables based on RSR alone. The PPD approach (19) retains the posterior mean of the RSR regression parameters, but in simulation, it increases the variance of the posterior distribution of regression parameters and results in low Type-S error rates for all true models.

Our simulation study examines the effects of model misspecification when the true spatial generating model is known. In practice, the true generating process that gives rise to observed spatially correlated data is almost never known. Our intent is to illustrate differences in the inference provided by RSR and SGLMM approaches. Our simulation studies illustrate cases where spatial confounding is likely to be

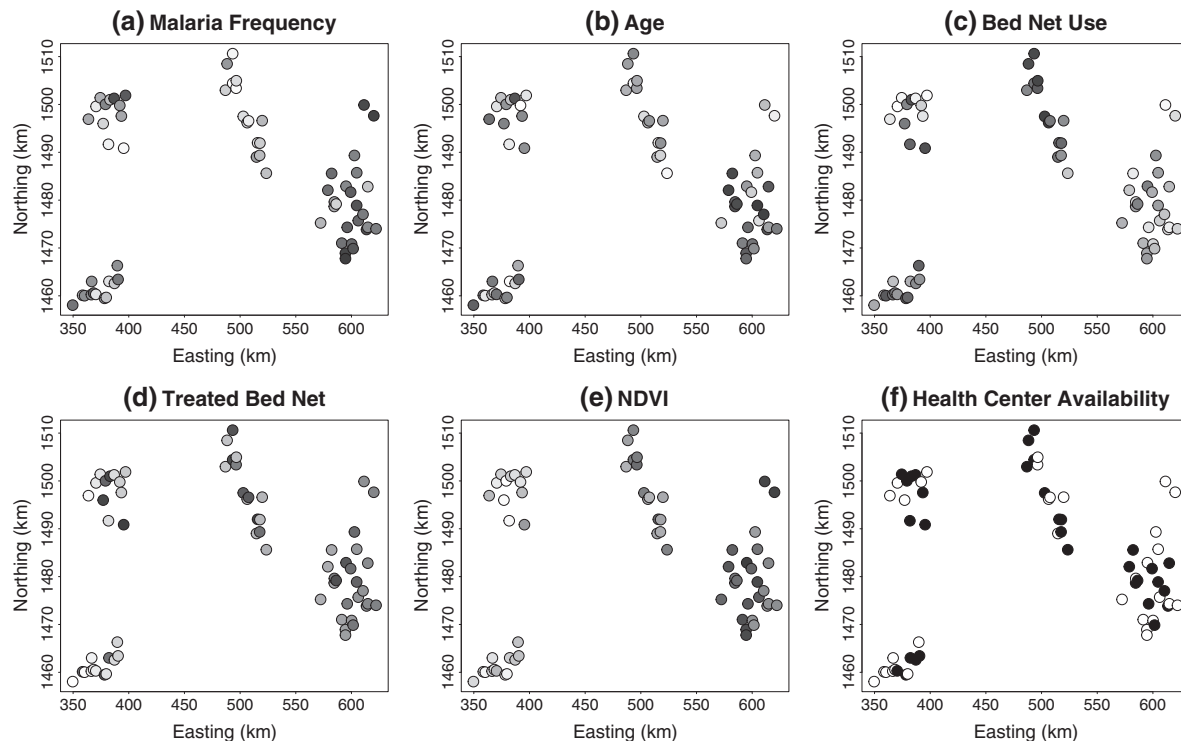


Figure 4. Spatial locations of 65 villages in The Gambia, Africa where children were tested for malaria. Shading indicates percent of children with malaria (a), average age of individuals in the study (b), percent of individuals using bed nets (c), percent of those bed nets treated with insecticide (d), the normalized difference vegetation index (e), and an indicator of whether a health center is present in the village (f). In all cases, darker shading indicates increased values of the variable shown

present (when spatial covariates and the spatial random effect have a large range parameter) and indicate that the conditional regression parameter estimates from RSR are likely to be less conservative than the unconditional regression parameter estimates from SGLMM.

6. EXAMPLE: MALARIA IN GAMBIA

We illustrate modeling in the presence of spatial confounding using binary presence and absence data of malaria in The Gambia, Africa (Figure 4). Malarial infection status (present or not) was recorded for $n = 2035$ children in $m = 65$ villages in The Gambia. Additional information recorded for each child includes the child's age, an indicator of whether (1) or not (0) the child sleeps under a bed net, and if so, an indicator for whether or not his or her bed net is treated with insecticide, the normalized difference vegetation index (NDVI) at each village derived from satellite data, as well as an indicator of presence or absence of a health center in the village. The NDVI measure and indicator for health center availability are common across all children in the i -th village. Further details about the data can be found in Thomson *et al.* (1999), who first analyzed this data by fitting a logistic regression model adjusted for spatially correlated errors using generalized estimating equations. Diggle *et al.* (2002) fit a geostatistical model to the data using an SGLMM.

Our goal is not to replicate these previous analyses. Rather, we use this data set that is publicly available in the 'geoR' package (Diggle and Ribeiro, 2007) in the R statistical computing environment (R Core Team, 2013) to illustrate the benefits and potential pitfalls of modeling SGLMMs using RSR. We begin with a data analysis and conclude with a simulation study.

6.1. Analysis of malaria data

Let y_{ij} be an indicator for the presence ($y_{ij} = 1$) or absence ($y_{ij} = 0$) of malaria in the j th child in the i th village. The covariate vector, \mathbf{x}_{ij} contains the covariates corresponding to the i -th child's age, net use, and net treatment with insecticide and the j -th village's NDVI and health center availability. We consider binary SGLMM models with the probit (Gaussian CDF) link function and analyze this data using probit regression models without spatial correlation (NS), with an unrestricted spatial random effect (SGLMM), and with a spatial random effect constrained to be orthogonal to the covariates (RSR). Let

$$y_{ij} \sim \text{Binom}(1, \Phi(z_{ij})), \quad (20)$$

where $\Phi(\cdot)$ is the probit link function ($\Phi(w) = P(W \leq w)$ if $w \sim N(0, 1)$). To match notation with previous sections, let $\mathbf{z} = (z_{11} \dots z_{n_m m})'$ be a vector of the z_{ij} for all villages ($j = 1, 2, \dots, m$) and children ($i = 1, 2, \dots, n_j$) in the study. Then the NS probit regression model is specified by

$$\mathbf{z} = \mathbf{X}\boldsymbol{\beta} \quad (21)$$

Diggle *et al.* (2002) found that, after fitting a similar non-spatial binary regression, a residual analysis indicated spatial heterogeneity in the residuals. Our residual analyses yielded similar results (results not shown). To specify an appropriate model in this case, we follow Diggle *et al.* (2002) and add a spatial random effect to the conditional mean of \mathbf{z} in (21). The latent unrestricted model (SGLMM) we specify is

$$\mathbf{z} = \mathbf{X}\boldsymbol{\beta} + \mathbf{K}\boldsymbol{\eta} \quad , \quad \boldsymbol{\eta} \sim N(\mathbf{0}, \boldsymbol{\Sigma}), \quad (22)$$

where $\boldsymbol{\Sigma}$ is a covariance matrix for a Matern (5) random field with smoothness parameter $\nu = 1/2$ (exponential covariance), σ^2 is the partial sill parameter, and ϕ is the Matern range parameter. The design matrix \mathbf{K} links observations to spatial locations and has $N = 2350$ rows, equal to the total number of children in the study and $m = 65$ columns equal to the number of unique spatial locations (villages) in the study. All entries in the l -th row of \mathbf{K} are zero except for $K_{lj} = 1$, where the l -th individual in the study comes from the j -th village.

To illustrate an RSR approach to this analysis, we replace the spatial random effect $\mathbf{K}\boldsymbol{\eta}$ in (22) with a spatial random effect constrained to be orthogonal to the fixed effects \mathbf{X} . The resulting latent model is

$$\mathbf{z} = \mathbf{X}\boldsymbol{\delta} + (\mathbf{I} - \mathbf{P}_X)\mathbf{K}\boldsymbol{\eta} \quad , \quad \boldsymbol{\eta} \sim N(\mathbf{0}, \boldsymbol{\Sigma}). \quad (23)$$

We note that, while we are modeling the latent spatial random effect $\boldsymbol{\eta}$ as continuous in space, we only have covariate values at the locations of the 65 villages in the study. Without continuous spatial covariate information, we cannot constrain the spatial random effect to be orthogonal to the fixed effects on the spatial domain of the study (10). Instead, we are constraining the random effect to be orthogonal to the fixed effects on the support of the observations (11).

We have chosen in (23) to constrain the spatial random effect $\mathbf{K}\boldsymbol{\eta}$ to be orthogonal to \mathbf{X} . This forces orthogonality on the support of \mathbf{X} . An alternative approach would be to constrain orthogonality between $\boldsymbol{\eta}$ and the mean of \mathbf{X} at each spatial location. This could be accomplished by constraining $\boldsymbol{\eta}$ to be orthogonal to $(\mathbf{K}'\mathbf{K})^{-1}\mathbf{K}'\mathbf{X}$ using conditioning by Kriging (7).

We specified diffuse priors for the conditional ($\boldsymbol{\beta}$) and marginal ($\boldsymbol{\delta}$) regression parameters;

$$\boldsymbol{\beta} \sim N(\mathbf{0}, 10^6 \mathbf{I}) \quad (24)$$

$$\boldsymbol{\delta} \sim N(\mathbf{0}, 10^6 \mathbf{I}). \quad (25)$$

We specified a conjugate inverse gamma prior for σ^2 with mean and variance of 10 and a truncated half Cauchy prior for the range parameter ϕ ;

$$[\phi] \propto \frac{1}{1 + \phi^2/100^2} 1_{\{\phi \in (0, 75000)\}}, \quad (26)$$

where the upper bound of the prior distribution ($\phi < 75000$ m) was chosen so that the effective range of spatial correlation was less than half of the maximum distance between two villages in the study region.

We constructed an MCMC algorithm to obtain samples from the posterior distributions of parameters in the NS, SLM, and RSR models. In each case, the MCMC sampler was run until the Monte Carlo standard error (e.g., Flegal *et al.*, 2008) was smaller than 0.01 for all model parameters. Conjugate updates were used for all parameters except for ϕ , for which we used random walk Metropolis–Hastings updates. Using methods described in Section 4, we also drew a sample from the posterior predictive distribution (RSR-PPD) of the conditional regression parameters (19) in the RSR model at each iteration of the MCMC sampler.

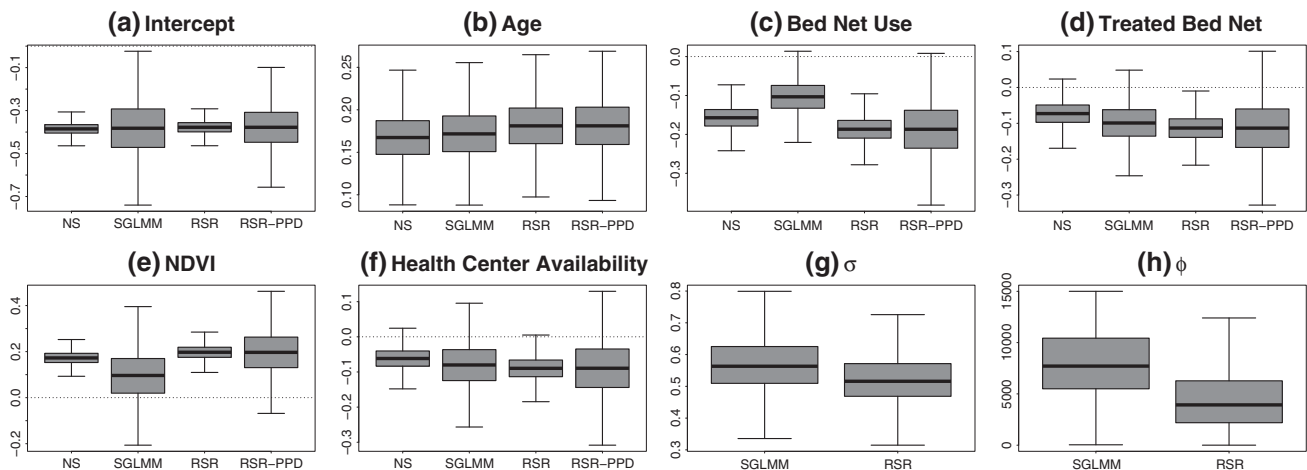


Figure 5. Boxplots of posterior distributions of model parameters given the Gambia malaria incidence data are shown for a non-spatial (NS) probit regression model, a spatial generalized linear mixed model, a restricted spatial regression model (RSR), and a posterior predictive approach (RSR-PPD)

Table 1. Type-S error rates under simulation for the regression parameter associated with normalized difference vegetation index in the Gambia malaria data

True Model	Model Fit		
	SLM	RSR	RSR-PPD
SLM	0.069	0.700	0.069
RSR	0.004	0.083	0.001

Boxplots of posterior distributions of model parameters are shown in Figure 5. To illustrate modeling choices and consequences in the presence of spatial confounding, consider the results for the regression parameters related to bed net use (Figure 5(c)) and NDVI (Figure 5(e)). Under both the NS and RSR models, the 95% credible intervals of these parameters do not overlap zero, while under the SLM model, the 95% credible intervals do overlap zero, leading to potentially conflicting results between the conditional and unconditional relationships between bed net use, NDVI, and malaria prevalence.

The RSR decorrelates the fixed and random effects and results in faster mixing of model parameters, leading to quicker convergence of the MCMC algorithms most often used for fitting spatial models. To illustrate this, we compared the effective sample size (ESS) (e.g., Givens and Hoeting, 2012) of the MCMC samples of the regression parameter corresponding to NDVI for both SGLMM and RSR approaches with the malaria data. In both cases, the MCMC algorithms were tuned so that the ϕ parameter was accepted in the Metropolis–Hastings step approximately 40% of the time. All other parameters were sampled using Gibbs updates from their respective full-conditional distributions. As the MCMC samplers were terminated based on a fixed width stopping rule (Flegal *et al.*, 2008), the MCMC chains under RSR, and SGLMM approaches have different lengths. For a fair comparison, we considered the first 10,000 samples obtained from each MCMC sampler. The ESS for the samples obtained under the SGLMM approach (ESS = 197.2) is an order of magnitude smaller than under RSR (ESS = 3054.7). Similar results for ESS are obtained in our simulation studies and in the full analysis of the malaria data. In our experience, the computational benefits of RSR over SGLMMs are especially pronounced when using ICAR priors on spatial random effects for areal data and constitute one of the main advantages of RSR.

6.2. Simulation study

We conducted a simulation study motivated by the Gambia malaria data to compare SGLMM and RSR approaches in a binary regression setting with multiple covariates. We specified the true model to be simulated from as the SGLMM model (20 and 22) with model parameters (β, ϕ, σ) set equal to the posterior mean values from the SGLMM fit (Figure 5) except for the regression parameter corresponding to NDVI, which we set equal to zero to allow us to examine Type-S error rates of various models in binary regression. Under this model, we simulated malaria presence or absence (y_{ij}) for each individual in the study and then fit the SGLMM and RSR models to the simulated data. We repeated this simulation study for the case where the true model to be simulated from is the RSR model (20 and 23). Estimated Type-S error rates from 1000 simulations from each true model are shown in Table 1. As in our previous simulation study, the RSR model shows a high Type-S error rate when the true model is the SGLMM and an acceptable Type-S error rate when the true model is RSR. The SGLMM and the posterior predictive distribution of unconditional regression parameters under the RSR model (RSR-PPD) are more conservative and exhibit Type-S error rates closer to the nominal rate of 0.05 when the true model is the SLM and Type-S error rates below 0.05 when the true model is RSR.

7. DISCUSSION

It is critical in any modeling endeavor to consider the aims of the study when specifying a statistical model, and the same holds in the choice between SGLMM and RSR. Common goals of studies employing spatial random effects include inference on regression parameters related to fixed effects, inference on variance parameters governing spatial autocorrelation, prediction at unobserved locations, mean function estimation, and exploratory analysis of spatial variation not modeled in fixed effects. We consider the respective benefits of RSR and SGLMM approaches for each of these in turn.

If inference on regression parameters is of highest importance in the study, the advantages and disadvantages of RSR and SGLMM approaches are clear. RSR offers increased computational efficiency but at the cost of significantly increased risk of Type-S errors if the true model is the unconstrained SGLMM or RSR on the full spatial support. If RSR is employed, the PPD approach (19) offers an approach that retains the unconditional regression parameter mean but does not suffer from RSR's increased risk of Type-S errors. Conversely, if the spatial random effects are truly orthogonal to the fixed effects, then incorrectly fitting an unconstrained SGLMM could lead to underestimating the size of the effect of the spatial covariates on the response, as correlation between the spatial covariates and the spatial random effect can result in variance inflation. As SGLMM and RSR approaches can be represented as reparameterizations of each other (Section 4) it will be difficult to choose between SGLMM and RSR models using standard model selection approaches, such as Bayes factors or information criteria comparison (e.g., Spiegelhalter *et al.*, 2002; Hooten and Hobbs, 2015). Instead, choosing among RSR and SGLMM approaches should be based on an understanding of the system under study. In the absence of a compelling reason to assume that the random effects

should be orthogonal to the fixed effects, we recommend erring on the side of caution and employing the conservative SGLMM or RSR-PPD approaches over the RSR approach that is less conservative under model misspecification.

Consider choosing among RSR and SGLMM approaches when the generating mechanism for spatially autocorrelated observations is assumed to be a spatially smooth missing covariate. Employing an RSR approach in this setting assumes that this missing covariate is orthogonal to the measured covariates. In general, this is a strong assumption as smooth covariates are likely to be collinear (See the supporting information, Section S1).

One common situation where a version of RSR would be almost universally acceptable is in the case where an intercept is included in the fixed effects. The intercept is often treated as a nuisance parameter necessary in the model but not typically interpreted. Constraining the spatial random effect η to be orthogonal to the intercept (but with no constraint relative to other fixed effects) can often improve mixing, especially in the case where the range of spatial structure (ϕ in the Matern (5) covariance function) is large. This constraint that $\eta' \mathbf{1} = 0$ forces the random effect to be centered at zero and considers estimating spatial variability in the mean as a linear combination of the intercept, fixed effects, and constrained random effect ($\mathbf{X}\beta + (\mathbf{I} - \mathbf{1}\mathbf{1}'/n)\eta$). This constraint is commonly employed when dependence is modeled using ICAR covariance (e.g., Besag and Kooperberg, 1995; Besag *et al.*, 1995; Rue and Held, 2005).

If inference on the spatial mean function ($\mathbf{X}\beta + \eta$) is of primary importance, RSR offers increased computational efficiency, and the duality between RSR and SGLMM approaches described in Section 4 ensures that the mean function can be estimated equally well with either approach.

If inference on variance parameters is of prime importance in the study, the RSR approach can be viewed as a Bayesian approach to restricted maximum likelihood (REML) estimation. When the parameters related to fixed effects are considered nuisance parameters, the RSR approach provides increased computational efficiency and is recommended.

If the main goal of the study is prediction at unobserved locations, the chosen model needs to accommodate joint modeling of observed and unobserved locations. This is straightforward in the SGLMM approach and in RSR if the constraint is on the full support of the observed and unobserved locations (9). In simulation (results not shown), there is little difference in the mean-squared prediction error at unobserved locations between these two approaches. Under RSR with the constraint only specified on the support of the observations (11), there is no straightforward joint model for observed and unobserved responses. In this case, prediction is unclear.

In some studies, inference on the random effect η is an important exploratory goal. Examination of the posterior distribution of $\eta|y$ may suggest new studies or provide insights into unmodeled mechanisms driving variation in the response. In this case, the RSR approach affords an analysis of the residual spatial structure after accounting for the included spatial covariates. Plotting the posterior mean of $(\mathbf{I} - \mathbf{P}_x)\eta$ would give a spatial residual plot that may provide intuition into missing covariates or the spatial locations of anomalies in the data not accounted for by the observed spatial covariates.

In summary, we have shown that restricting spatial random effects to be orthogonal to fixed effects provides computational benefits and can be accomplished for models with continuous spatial support but should be used with care because of the assumption that all variation in the direction of the fixed effects be ascribed to them. If the parameters related to the fixed effects are considered nuisance parameters, then RSR is often appropriate, but if inference on parameters related to fixed effects is a key goal, then the SGLMM and RSR-PPD approaches are more robust to model misspecification.

Example code to carry out the SGLMM, RSR, and RSR-PPD approaches for geostatistical models as described in this manuscript is available at <http://sites.psu.edu/hanks/>. Supplementary material for this article is available online at the journal's website.

Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. EF-0914489 (Hoeting). Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the US government.

REFERENCES

- Banerjee S, Gelfand AE, Finley AO, Sang H. 2008. Gaussian predictive process models for large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **70**(4): 825–848.
- Besag J, Green P, Higdon D, Mengersen K. 1995. Bayesian computation and stochastic systems. *Statistical Science* **10**(1): 3–41.
- Besag J, Kooperberg C. 1995. On conditional and intrinsic autoregressions. *Biometrika* **82**(4): 733.
- Boots B, Tiefelsdorf M. 2000. Global and local spatial autocorrelation in bounded regular tessellations. *Journal of Geographical Systems* **2**(4): 319–348.
- Cressie N. 1993. *Statistics for Spatial Data*. Wiley-Interscience: New York City, NY.
- Cressie N, Johannesson G. 2008. Fixed rank kriging for very large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **70**(1): 209–226.
- Diggle P, Moyeed R, Rowlingson B, Thomson M. 2002. Childhood malaria in the Gambia: a case-study in model-based geostatistics. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* **51**(4): 493–506.
- Diggle P, Ribeiro PJ. 2007. *Model-Based Geostatistics*. Springer: New York City, NY.
- Diggle PJ, Ribeiro PJ, Jr. 2002. Bayesian inference in Gaussian model-based geostatistics. *Geographical and Environmental Modelling* **6**(2): 129–146.
- Flegal JM, Haran M, Jones GL. 2008. Markov chain Monte Carlo: can we trust the third significant figure? *Statistical Science* **23**(2): 250–260.
- Gelman A, Carlin BP, Stern H, Rubin DB. 2004. *Bayesian Data Analysis* 2nd ed. Chapman and Hall/CRC: Princeton, New Jersey, USA.
- Gelman A, Tuerlinckx F. 2000. Type S error rates for classical and Bayesian single and multiple comparison procedures. *Computational Statistics* **15**(3): 373–390.
- Givens GH, Hoeting JA. 2012. *Computational Statistics*, Vol. 708. John Wiley & Sons: New York City, NY.
- Griffith DA. 2000. A linear regression solution to the spatial autocorrelation problem. *Journal of Geographical Systems* **2**(2): 141–156.
- Griffith DA, Peres-Neto PR. 2006. Spatial modeling in ecology: the flexibility of eigenfunction spatial analyses. *Ecology* **87**(10): 2603–2613.

- Hodges JS, Reich BJ. 2010. Adding spatially-correlated errors can mess up the fixed effect you love. *The American Statistician* **64**(4): 325–334.
- Hooten MB, Hobbs NT. 2015. A guide to Bayesian model selection for ecologists. *Ecological Monographs*. DOI: 10.1890/14-0661.1.
- Hooten MB, Larsen DR, Wikle CK. 2003. Predicting the spatial distribution of ground flora on large domains using a hierarchical Bayesian model. *Landscape Ecology* **18**(5): 487–502.
- Hughes J, Haran M. 2013. Dimension reduction and alleviation of confounding for spatial generalized linear mixed models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **75**(1): 139–159.
- Lindgren F, Rue H, Lindström J. 2011. An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **73**(4): 423–498.
- Nychka D, Wikle C, Royle JA. 2002. Multiresolution models for nonstationary spatial covariance functions. *Statistical Modelling* **2**(4): 315–331.
- Paciorek CJ. 2010. The importance of scale for spatial-confounding bias and precision of spatial regression estimators. *Statistical Science* **25**(1): 107.
- R Core Team. 2013. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing: Vienna, Austria. ISBN 3-900051-07-0.
- 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.
- Royle JA, Wikle CK. 2005. Efficient statistical mapping of avian count data. *Environmental and Ecological Statistics* **12**(2): 225–243.
- Rue H, Held L. 2005. *Gaussian Markov Random Fields: Theory and Applications*, Monographs on Statistics and Applied Probability, vol. 104. Chapman & Hall: Boca Raton, FL.
- 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.
- Thomson MC, Connor SJ, D'Alessandro U, Rowlingson B, Diggle P, Cresswell M, Greenwood B. 1999. Predicting malaria infection in Gambian children from satellite data and bed net use surveys: the importance of spatial correlation in the interpretation of results. *American Journal of Tropical Medicine and Hygiene* **61**: 2–8.
- Tiefelsdorf M, Griffith DA. 2007. Semiparametric filtering of spatial autocorrelation: the eigenvector approach. *Environment and Planning A* **39**(5): 1193.
- Wikle. 2002. *Spatial Modeling of Count Data: A Case Study in Modeling Breeding Bird Survey Data on Large Spatial Domains*. Chapman and Hall/CRC: Boca Raton, FL, 199–209.

Copyright of Environmetrics is the property of John Wiley & Sons, Inc. and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.