

Semiparametric Bayesian analysis of transformation spatial mixed models for large datasets*

YING WU, DAN CHEN[†], AND NIANSHENG TANG

In *spatial mixed models* (SMMs), it is commonly assumed that stationary spatial process and random errors independently follow the Gaussian distribution. However, in some applications, this assumption may be inappropriate. To this end, this paper proposes a *transformation spatial mixed models* (TSMMs) to accommodate large dataset that follows the non-Gaussian distribution. With the help of Gibbs sampler algorithm, a semiparametric Bayesian approach is developed to make inference on TSMMs by using Bayesian P-splines to approximate transformation function, and a fixed number of known but not necessarily orthogonal spatial basis functions with multi-resolution analysis method to approximate nonstationary spatial process. Instead of Wishart distribution assumption for the prior of precision matrix of random effects, we consider Cholesky decomposition of the precision matrix, and specify the priors for unknown components in low unit triangular matrix and diagonal matrix. Simulation studies and an example are used to illustrate the proposed methodologies.

AMS 2000 SUBJECT CLASSIFICATIONS: Primary 62F15; secondary 62J05.

KEYWORDS AND PHRASES: Gibbs sampler, Large dataset, Matrix decomposition, Semiparametric Bayesian analysis, Transformation spatial Mixed Models.

1. INTRODUCTION

In much economical, social, physical and epidemiological research, substantive theory and practice usually involve the analysis of spatial data (e.g., Cressie [3], Anselin and Bera [4], Anselin and Syabri [12]). Various statistical models, theories and methods have been developed to analyze spatial data over the past years (e.g., Bai [5], Cai [9], Chan [10]). In particular, SMMs are widely adopted to capture spatial structures of spatial and spatial-temporal datasets from very

fine to very large scales by incorporating random effects in the analysis of spatial data. For example, see Cressie and Johannesson [12] for the analysis of very large datasets, Kang and Cressie [23] for Bayesian inference, Nychka, Wikle and Royle [32], Nychka et al. [31] and Katzfuss [25] for multi-resolution analysis based on the W-wavelet theories (e.g., Kwong and Tang [26]).

At present, there are more than a dozen spatial data analysis software packages to adapt to strong demand in various fields, for example, see SpaceStat package [1], GeoDa [2], spdep and DCluster on the open source Comprehensive R Archive Network, and ArcGIS 10.7. However, these packages cannot deal with non-normal spatial datasets. In particular, they cannot deal with binary/poisson data with a binomial/poisson distribution, which is often encountered in economical, social and epidemiological studies. For example, a remote sensing dataset, which was retrieved by the *Multi-angle Imaging SpectroRadiometer* (MISR) instrument on NASA's Terra and Aqua satellite, was collected monthly at a spatial resolution, which was converted to 3-degrees dataset (e.g., 1.0×1.0 , 0.5×0.5 , 0.25×0.25). The MISR instrument is one of the key instruments on board that were widely used to collect the global atmosphere information such as the *aerosol optical depth* (AOD), *Water Vapor* (WV), *Net Radiation* (NR) and *Carbon Monoxide* (CM). The histograms of the 2-degrees AOD, WV, NR and CM datasets during May 2016 between longitudes -179.75 and $+179.75$, and between latitudes -89.75 and $+89.75$ (e.g., Figure 6) show that the AOD, WV, NR and CM data are non-normal. Hence, the large scale dataset motivates considering the transformation spatial mixed models.

Considerable work has been focused on approaches to relax the normality assumption in spatial models over the past years. For example, see De Oliveira et al. [15], Genton and Zhang [19], Zhang and EI-Shaarawi [41], and Zareifard and Khahedi [40]. The aforementioned literature has been developed under the skewed distributional assumption of spatial data by utilizing the Box-Cox family of power transformations or a scale mixing of a unified skew Gaussian process. However, for transformation linear mixed models rather than spatial mixed models, Gurka [21] demonstrated that estimating transformation parameter in the Box-Cox family of power transformations may yield biased estimation

*This research is supported by the National Natural Science Foundation of China (No. 11671349) and the Key Projects of the National Natural Science Foundation of China (No. 11731101), the National Social Science Fund of China (No. 17BTJ038), Yunnan Applied Basic Research Projects (No. 2016FD087).

[†]Corresponding author.

of variance components, and maximum likelihood estimator of transformation parameter is not consistent. To address these issues, various transformation models have been developed in a Bayesian framework. In particular, Bayesian P-splines approach has been used to handle transformation models. For example, Song and Lu [36] presented a semiparametric transformation model with P-splines; Song and Liu [35] considered a transformation structure equation model with highly non-normal and incomplete data with P-splines; Tang, Wu and Chen [37] extended to a transformation linear mixed model based on the truncated centered Dirichlet Process prior approximation to the distribution of random effects. However, to our knowledge, little work has been done on TSMMs based on the idea of Bayesian P-splines in a Bayesian framework.

The main contributions of this paper include as follows. First, the P-splines is employed to approximate transformation functions because of its good properties, such as allowing for simultaneous estimation of smooth functions and smoothing parameters, allowing the smoothing parameters to be locally adaptive, and the flexibility of implementing Bayesian analysis via the R program. Second, the multi-resolution analysis technique with a fixed number of known but not necessarily orthogonal spatial basis functions is adopted to approximate nonstationary spatial process. Third, a *Markov chain Monte Carlo* (MCMC) algorithm is developed to simultaneously obtain Bayesian estimations of unknown parameters, random effects and transformation functions by combining the Gibbs sampler and W-wavelet theories. Fourth, instead of traditional Wishart distribution assumption for the prior of precision matrix of random effects, we consider Cholesky decomposition of the precision matrix for random effects, and then specify the priors for unknown components in low unit triangular matrix and diagonal matrix, which largely reduce the number of unknown parameters and improve the convergence rate of MCMC algorithm. Fifth, we investigate the identifiability of the considered model. Sixth, we present a Bayesian model comparison approach by utilizing Bayes factor and Path sampling method, and discuss the goodness-of-fit assessment problem via the posterior predictive (PP) p -value.

The rest of this paper is organized as follows. Section 2 introduces TSMMs including the selection of the W-wavelet matrix, Bayesian P-splines approach to approximate the transformation function. Bayesian inference including parameter estimation, model comparison and goodness-of-fit statistic are given in Section 3. Simulation studies and an example are illustrated in Section 4. Technical details are presented in the Appendix.

2. TRANSFORMATION SPATIAL MIXED MODELS

2.1 Models and notation

Let $\{Y(\mathbf{s}) : \mathbf{s} \in \mathbb{S} \subset \mathbb{R}^d\}$ be a real-valued spatial process. Our main interest is to make statistical inference on

the Y -process on the basis of the actually observed non-normal dataset $\{Z(\mathbf{s}) : \mathbf{s} \in \mathbb{S}\}$. To this end, we consider the following measurement error model

$$(1) \quad f(Z(\mathbf{s})) = Y(\mathbf{s}) + \varepsilon(\mathbf{s}), \quad \mathbf{s} \in \mathbb{S},$$

where $\{\varepsilon(\mathbf{s}) : \mathbf{s} \in \mathbb{S}\}$ is an independent Gaussian process that is independent of $Y(\cdot)$ and has mean zero and $\text{var}\{\varepsilon(\mathbf{s})\} = \sigma_\varepsilon^2$. It is assumed that the process $Z(\cdot)$ is known only at a finite number of spatial locations, e.g., $\{\mathbf{s}_1, \dots, \mathbf{s}_N\}$. Let $\mathbf{Z} \equiv \{Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_N)\}$ be the available data, $f(\cdot)$ be an unknown smooth transformation function for $Z(\mathbf{s})$, and denote $Z^*(\mathbf{s}) = f(Z(\mathbf{s}))$. Suppose that $f(\cdot)$ is strictly monotone and differentiable. Following Cressie and Johannesson [12] and Kang and Cressie [23], we consider the following *spatial mixed models* (SMMs):

$$(2) \quad \begin{cases} Y(\mathbf{s}) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \nu(\mathbf{s}) + \xi(\mathbf{s}), \\ \nu(\mathbf{s}) = \mathbf{W}(\mathbf{s})^\top \boldsymbol{\eta}, \quad \mathbf{s} \in \mathbb{S}, \end{cases}$$

where $\mathbf{X}(\cdot) = (X_1(\cdot), \dots, X_p(\cdot))^\top$ is a $p \times 1$ vector of covariate processes, $\boldsymbol{\gamma}$ is a $p \times 1$ vector of coefficients corresponding to the fixed effects, $\mathbf{X}(\cdot)^\top \boldsymbol{\gamma}$ describes the large-scale spatial variation, $\nu(\cdot)$ represents the smooth small-scale spatial variation, $\xi(\cdot)$ is the residual term, $\boldsymbol{\eta}$ is a $q \times 1$ vector of random effects, $\boldsymbol{\eta}$ and $\xi(\cdot)$ define an SMM that can handle spatial variability for small and fine scales, $\mathbf{W}(\cdot)$ is a $q \times 1$ vector of q deterministic, known and multi-resolutional spatial basis functions that are not necessarily orthogonal of each other, i.e., $\mathbf{W}(\cdot) = (\mathbf{w}_1(\cdot), \dots, \mathbf{w}_q(\cdot))^\top$. It is assumed that $\xi(\cdot)$ is independent of $\boldsymbol{\eta}$, $\boldsymbol{\eta}$ is distributed as the multivariate Gaussian distribution with mean zero and $\text{cov}(\boldsymbol{\eta}) = \boldsymbol{\Sigma}$, i.e., $\boldsymbol{\eta} \sim \mathcal{N}_q(\mathbf{0}, \boldsymbol{\Sigma})$, and $\xi(\cdot)$ follows the Gaussian distribution $\mathcal{N}(0, \sigma_\xi^2)$. Thus, Equations (1) and (2) define a *transformation spatial mixed models* (TSMMs). From the definition of the TSMMs, it is easily seen that $Z(\mathbf{s})$ is non-normal dataset, but $Z^*(\mathbf{s})$ is normally distributed utilizing the transformation function $f(\cdot)$.

2.2 Generating multiresolutional spatial basis functions via W wavelet

In model (2), $\mathbf{W}(\cdot)$ is usually unknown. To this end, following Kwong and Tang [26], Nychka, Wikle and Royle [32] and Katzfuss [25], we generate a spatial basis using the repeated translations and scalings of several fixed functions, it is called discrete wavelet transform. To generate W wavelet, we first choose a B-splines as the fixed function, and then use the repeated W translations and scalings of B-splines to obtain two templates (e.g., mother and father), and finally determine the number of W wavelets at different resolutions in the same domain. As an illustration, we consider wavelets defined in the interval $[0, 1]$, and a coarsest level of resolution (e.g., say L). The W wavelets used here are plotted in Figure 1, and a basis of 32 functions for $L = 3$ is displayed in Figure 2. The first 3 resolution basis functions are similar

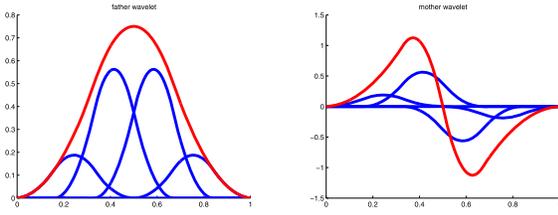


Figure 1. Continuous versions of the father and mother W -transform wavelets.

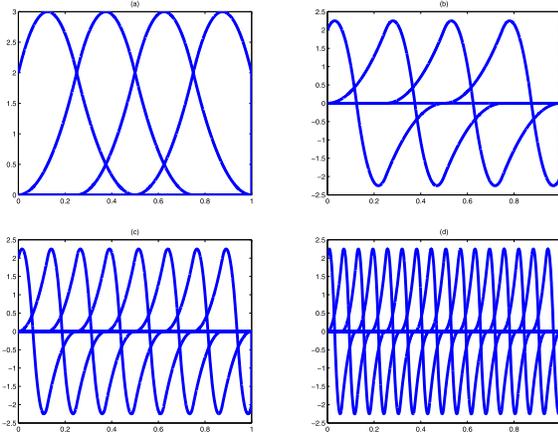


Figure 2. Family of 32 basis functions, (a) is the first 4 functions based on translations of the father wavelet; (b) is the same translated functions for the mother; (c) is 8 translations of the mother wavelet reduced by a factor of two; (d) is 16 translations of the reduced mother wavelet.

to the father wavelet translated in equally spaced locations (e.g., see Figure 2(a)). The first 3 resolution basis functions are the mother wavelet translated by the same manner, and are shown in Figure 2(b). The next generation of basis functions has twice as the resolution, and is similar to a scaling and translation of the mother wavelet. Figure 2(c) shows this generation. This cascade continues with the number of members in each subsequent generation, and the resolution increasing by a factor of two. Figure 2(d) completes the basis of size 32. The relationship between W -matrix and W -wavelet is given in Appendix.

2.3 P-splines approximation of transformation function

Due to the unknown form of function $f(\cdot)$, it is impossible to make statistical inference on γ and Σ via models (1) and (2). To this end, we here adopt the widely used Bayesian P-splines method to approximate $f(\cdot)$. Following Lang and Brezger [27], $f(Z(\mathbf{s}))$ can be approximated by the linear combination of the finite B-splines (e.g., De Boor [14]). That is, $f(Z(\mathbf{s})) \approx \sum_{k=1}^K \beta_k B_k(Z(\mathbf{s}))$, where K is the number of B-splines determined by the number of knots in

the support of $Z(\mathbf{s})$ and the degrees of B-splines, β_k 's are unknown coefficients, and $B_k(Z(\mathbf{s}))$ is the k th B-splines for $Z(\mathbf{s})$. Thus, based on the above approximation of $f(Z(\mathbf{s}))$, for any $\mathbf{s} \in \mathbb{S}$, (1) and (2) can be written as the following unified form

$$(3) \quad \sum_{k=1}^K \beta_k B_k(Z(\mathbf{s})) \approx \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \mathbf{W}(\mathbf{s})^\top \boldsymbol{\eta} + \xi(\mathbf{s}) + \varepsilon(\mathbf{s}).$$

To control unknown function $f(\cdot)$ via penalizing the coefficients of the adjacent B-splines, we consider a random walk prior to β_k . In particular, we consider the following first and second order random walks for β_k :

$$(4) \quad \begin{aligned} \beta_k &= \beta_{k-1} + u_k, & \text{for } k = 2, \dots, K, \\ \beta_k &= 2\beta_{k-1} - \beta_{k-2} + u_k, & \text{for } k = 3, \dots, K, \end{aligned}$$

respectively, where $u_k \sim \mathcal{N}(0, \tau^2/\varphi_k)$. It is assumed that β_1 follows a diffuse prior, i.e., $\beta_1 \propto \text{constant}$ for the first-order random walk, and $\beta_1 \propto \text{constant}$ and $\beta_2 \propto \text{constant}$ for the second-order random walk. Here τ^2 can be regarded as an inverse smoothing parameter that controls the smoothness of the resulting function $f(\cdot)$, and φ_k 's are treated as local smoothing parameters, which are employed to control the local smoothness of a function with significantly different curvatures at different $Z(\mathbf{s})$'s. The above defined prior distribution for $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)^\top$ can be written as

$$p(\boldsymbol{\beta}) \sim \left(\frac{1}{\sqrt{2\pi\tau}}\right)^{K-d} \exp\left\{-\frac{1}{2\tau^2} \boldsymbol{\beta}^\top \mathbf{M}(\boldsymbol{\varphi}) \boldsymbol{\beta}\right\} \times I(\beta_1 < \dots < \beta_K, \mathbf{Q}\boldsymbol{\beta} = 0),$$

where $\mathbf{M}(\boldsymbol{\varphi}) = \mathbb{D}^\top \text{diag}(\boldsymbol{\varphi}) \mathbb{D}$, $\mathbb{D} = \mathbf{D}_{d-1} \mathbf{D}_{d-2} \cdots \mathbf{D}_1 \mathbf{D}_0$, $\boldsymbol{\varphi} = (\varphi_{d+1}, \dots, \varphi_K)^\top$, \mathbf{D}_h is a $(K-h-1) \times (K-h)$ matrix and is defined as

$$\mathbf{D}_h = \begin{pmatrix} -1 & 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & -1 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & -1 & 1 \end{pmatrix}$$

for $h = 0, \dots, d-1$ with $d = 1$ and 2 , and d is the order of the random walk.

There are two kinds of approaches to select K . The first one is to fix K at some fixed value, and the other one is to regard K as a random variable that is selected by using the reversible jump MCMC algorithm or Birth-and-Death method. Lang and Brezger [27] compared the above mentioned two methods, and found that the former performed better for functions with moderate curvature, whilst the latter performed better for highly oscillating functions. When the transformation function is monotonic, it is more likely to be a function with the moderate curvature, so we fix K in some applications. More discussions on the selection of K can refer to Lang and Brezger [27].

2.4 Identification

For the above considered TSMMs, it is assumed that the transformation function $f(\mathbf{Z}(\mathbf{s})) = \sum_{k=1}^K \beta_k B_k(\mathbf{Z}(\mathbf{s}))$ is strictly monotone increasing, and β_k 's satisfy the following constraints $\beta_1 < \beta_2 < \dots < \beta_K$, which is widely adopted to make the estimation of the monotone transformation function easier even if this condition is sufficient but not necessary.

Note that the above considered TSMMs is unidentifiable because (β_k, μ) and $(\beta_k + c, \mu + c)$ have the same likelihood function for any constant c due to the following fact that $\sum_{k=1}^K \beta_k B_k(\mathbf{Z}(\mathbf{s})) - \mu = \sum_{k=1}^K (\beta_k + c) B_k(\mathbf{Z}(\mathbf{s})) - (\mu + c)$ for any constant c , which is obtained by the property of the B-spline: $\sum_{k=1}^K B_k(\cdot) = 1$.

For identifiability, we assume that $\sigma_\varepsilon^2 = 1$, the function $f(\cdot)$ has zero mean, i.e., $\sum_{\mathbf{s} \in \mathbb{S}} \sum_{k=1}^K \beta_k B_k(\mathbf{Z}(\mathbf{s})) = 0$, which is equivalent to $\mathbf{Q}\boldsymbol{\beta} = 0$, where $\mathbf{Q} = (\mathbf{B}_1^*, \dots, \mathbf{B}_K^*)$ in which $\mathbf{B}_k^* = \sum_{\mathbf{s} \in \mathbb{S}} B_k(\mathbf{Z}(\mathbf{s}))$ for $k = 1, \dots, K$, and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)^\top$.

3. BAYESIAN INFERENCE

3.1 Prior distribution

Following Kang and Cressie [23] and Song and Lu [36], we consider the following priors for γ , τ^{-2} , φ_k , σ_ε and $\boldsymbol{\beta}$:

$$\begin{aligned} p(\gamma) &\sim \mathcal{N}_p(\boldsymbol{\mu}_{\gamma 0}, \sigma_{\gamma 0}^2 \mathbf{I}_p), & p(\tau^{-2}) &\sim \Gamma(e_1, e_2), \\ p(\varphi_k) &\sim \Gamma(b/2, b/2), & p(\sigma_\varepsilon) &\sim U(0, \kappa_\varepsilon), \end{aligned}$$

where $\Gamma(\cdot, \cdot)$ represents the gamma distribution, $U(\cdot, \cdot)$ is the uniform distribution; $\boldsymbol{\mu}_{\gamma 0}$, $\sigma_{\gamma 0}^2$, e_1 , e_2 , b and κ_ε are the pre-given hyperparameters.

To present the prior of $\boldsymbol{\Sigma}$ associated with $\boldsymbol{\eta}$, following Chan and Jeliaskov [10], we consider the following Cholesky decomposition of $\boldsymbol{\Sigma}^{-1}$:

$$(5) \quad \boldsymbol{\Sigma}^{-1} = \mathbf{A}^\top \boldsymbol{\Lambda}^{-1} \mathbf{A},$$

where \mathbf{A} is a lower unit triangular matrix, and $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_q)$ in which $\lambda_j > 0$ for $j = 1, \dots, q$. For simplicity, let a_{ij} ($1 \leq j < i \leq q$) denote the free elements of lower unitriangular matrix \mathbf{A} , i.e.

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ a_{21} & 1 & 0 & \dots & 0 & 0 \\ a_{31} & a_{32} & 1 & \dots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ a_{q1} & a_{q2} & \dots & \dots & a_{q,q-1} & 1 \end{pmatrix}.$$

Also, denote $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_q)^\top$, $\mathbf{a}_i = (a_{i1}, \dots, a_{i,i-1})^\top$ for $i = 2, \dots, q$, and $\mathbf{a} = (\mathbf{a}_2^\top, \dots, \mathbf{a}_q^\top)^\top$. Based on the above parameterization, we consider the following priors for λ_i and \mathbf{a}_j ($i = 1, \dots, q$, $j = 2, \dots, q$):

$$\lambda_i^{-1} \sim \Gamma\left(\frac{v_{0i}}{2}, \frac{u_{0i}}{2}\right), \quad \mathbf{a}_j \sim \mathcal{N}_{j-1}(\mathbf{a}_{j0}, \mathbf{A}_{j0})$$

where v_{0i} , u_{0i} , \mathbf{a}_{j0} and \mathbf{A}_{j0} are the pre-given hyperparameters.

3.2 Gibbs sampler

To avoid the high-dimensional integral involved in making statistical inference via the marginal probability density function $p(\mathbf{Z})$, the Gibbs sampler is adopted to draw a sequence of random observations from the joint conditional distribution $p(\gamma, \boldsymbol{\eta}, \boldsymbol{\beta}, \sigma_\varepsilon^{-2}, \tau^{-2}, \boldsymbol{\varphi}, \mathbf{a}, \boldsymbol{\lambda} | \mathcal{D})$, where $\mathcal{D} = \{\mathbf{Z}, \mathbf{X}, \mathbf{W}\}$. The Gibbs sampler is implemented by iteratively drawing observations from the following conditional distributions: (a) $p(\gamma | \boldsymbol{\eta}, \sigma_\varepsilon^{-2}, \boldsymbol{\beta}, \mathcal{D})$, (b) $p(\boldsymbol{\eta} | \xi, \gamma, \boldsymbol{\beta}, \mathbf{a}, \boldsymbol{\lambda}, \mathcal{D})$, (c) $p(\xi | \gamma, \boldsymbol{\eta}, \sigma_\varepsilon^{-2}, \boldsymbol{\beta}, \mathcal{D})$, (d) $p(\boldsymbol{\beta} | \gamma, \boldsymbol{\eta}, \sigma_\varepsilon^{-2}, \boldsymbol{\varphi}, \mathcal{D})$, (e) $p(\sigma_\varepsilon^{-2} | \gamma, \boldsymbol{\eta}, \boldsymbol{\beta}, \mathcal{D})$, (f) $p(\tau^{-2} | \boldsymbol{\beta}, \boldsymbol{\varphi})$, (g) $p(\boldsymbol{\varphi} | \boldsymbol{\beta}, \tau^{-2})$ (h) $p(\mathbf{a} | \boldsymbol{\eta}, \boldsymbol{\lambda})$, (i) $p(\boldsymbol{\lambda} | \boldsymbol{\eta}, \mathbf{a})$.

The conditional distributions required in implementing the aforementioned Gibbs sampler are presented in the Appendix. Convergence of the above presented Gibbs sampler algorithm can be monitored by the *estimated potential scale reduction* (EPSR) value evaluated from several parallel sequences of observations as the runs proceed [17]. Convergence is claimed if all the EPSR values of unknown parameters are less than 1.2.

3.3 Bayesian estimates

Let $\{(\boldsymbol{\gamma}^{(m)}, \boldsymbol{\eta}^{(m)}, \boldsymbol{\beta}^{(m)}, \boldsymbol{\theta}_\eta^{(m)}) : m = 1, \dots, \mathcal{M}\}$ be the observations of $(\boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\beta}, \boldsymbol{\theta}_\eta)$ simulated from the joint conditional distribution $p(\boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\beta}, \sigma_\varepsilon^{-2}, \tau^{-2}, \boldsymbol{\varphi}, \boldsymbol{\theta}_\eta | \mathcal{D})$ via the preceding presented MCMC algorithm after the MCMC algorithm converges, where $\boldsymbol{\theta}_\eta$ is a set of unknown parameters associated with the distribution of random effects. Thus, Bayesian estimates of $\boldsymbol{\theta}_\eta, \boldsymbol{\beta}, \boldsymbol{\gamma}$ and $\boldsymbol{\eta}$ can be evaluated by

$$\begin{aligned} \widehat{\boldsymbol{\theta}}_\eta &= \frac{1}{\mathcal{M}} \sum_{m=1}^{\mathcal{M}} \boldsymbol{\theta}_{\eta\xi}^{(m)}, & \widehat{\boldsymbol{\beta}} &= \frac{1}{\mathcal{M}} \sum_{m=1}^{\mathcal{M}} \boldsymbol{\beta}^{(m)}, \\ \widehat{\boldsymbol{\gamma}} &= \frac{1}{\mathcal{M}} \sum_{m=1}^{\mathcal{M}} \boldsymbol{\gamma}^{(m)}, & \widehat{\boldsymbol{\eta}} &= \frac{1}{\mathcal{M}} \sum_{m=1}^{\mathcal{M}} \boldsymbol{\eta}^{(m)}, \end{aligned}$$

respectively. The function $f(\mathbf{Z})$ can be estimated by $\widehat{f}(\mathbf{Z}) = \sum_{k=1}^K \widehat{\beta}_k B_k(\mathbf{Z})$, where $\widehat{\beta}_k$ is the k th component of $\widehat{\boldsymbol{\beta}}$. Similarly, the consistent estimates of covariance matrices of $\widehat{\boldsymbol{\theta}}_\eta$, $\widehat{\boldsymbol{\beta}}$, $\widehat{\boldsymbol{\gamma}}$ and $\widehat{\boldsymbol{\eta}}$ can be obtained by using their corresponding sample covariance matrices of the simulated observations $\{(\boldsymbol{\theta}_\eta^{(m)}, \boldsymbol{\beta}^{(m)}, \boldsymbol{\gamma}^{(m)}, \boldsymbol{\eta}^{(m)}) : m = 1, \dots, \mathcal{M}\}$. For example, $\widehat{\text{var}}(\widehat{\boldsymbol{\beta}}) = \{\mathcal{M} - 1\}^{-1} \sum_{m=1}^{\mathcal{M}} (\boldsymbol{\beta}^{(m)} - \widehat{\boldsymbol{\beta}})(\boldsymbol{\beta}^{(m)} - \widehat{\boldsymbol{\beta}})^\top$. Then, the standard errors for the components of $\widehat{\boldsymbol{\beta}}$ can be obtained by using the square roots of their corresponding diagonal elements of $\widehat{\text{var}}(\widehat{\boldsymbol{\beta}})$.

3.4 Bayesian model comparison

Bayes factor is an important statistic for comparing several competing models in a Bayesian framework, and is

widely used to make model comparison in various statistical models. But it is rather difficult to calculate Bayes factor because of the high-dimensional integral involved. To address the issue, some methods such as bridge sampling (e.g., Meng and Wong [30]) and path sampling (e.g., Gelman and Meng [18]) have been developed. Here, the path sampling method is used to calculate Bayes factor for comparing two competing models \mathcal{H}_1 and \mathcal{H}_0 :

$$B_{10} = \frac{p(\mathbf{Z}^*, \mathbf{X}, \mathbf{W} | \mathcal{H}_1)}{p(\mathbf{Z}^*, \mathbf{X}, \mathbf{W} | \mathcal{H}_0)},$$

where $p(\mathbf{Z}^*, \mathbf{X}, \mathbf{W} | \mathcal{H}_t)$ is the marginal probability density of $(\mathbf{Z}^*, \mathbf{X}, \mathbf{W})$ under \mathcal{H}_t , $\mathbf{Z}^* = f(\mathbf{Z})$, $p(\mathbf{Z}^*, \mathbf{X}, \mathbf{W} | \mathcal{H}_t) = \int p(\mathbf{Z}^*, \boldsymbol{\eta}, \mathbf{X}, \mathbf{W} | \boldsymbol{\vartheta}_t, \mathcal{H}_t) p(\boldsymbol{\vartheta}_t | \mathcal{H}_t) d\boldsymbol{\eta} d\boldsymbol{\vartheta}_t$, $\boldsymbol{\vartheta}_t = \{\boldsymbol{\gamma}, \boldsymbol{\beta}, \sigma_\zeta^{-2}, \mathbf{a}, \boldsymbol{\lambda}, \boldsymbol{\varphi}, \tau^{-2}\}$ is the set of parameters under \mathcal{H}_t for $t = 0$ and 1. Similar to Lee and Song [28] and Lee and Tang [29], for a continuous parameter $t \in [0, 1]$, we consider the following class of probability density functions

$$F(t) = \int p(\mathbf{Z}^*, \boldsymbol{\eta}, \mathbf{X}, \mathbf{W}, t | \boldsymbol{\vartheta}_t) p(\boldsymbol{\vartheta}_t) d\boldsymbol{\eta} d\boldsymbol{\vartheta}_t,$$

where $p(\mathbf{Z}^*, \boldsymbol{\eta}, \mathbf{X}, \mathbf{W}, t | \boldsymbol{\vartheta}_t)$ is the joint probability density function of $(\mathbf{Z}^*, \boldsymbol{\eta}, \mathbf{X}, \mathbf{W})$ under \mathcal{H}_t that links \mathcal{H}_0 and \mathcal{H}_1 with the continuous parameter t such that $\mathcal{H}_t = \mathcal{H}_0$ if $t = 0$, and $\mathcal{H}_t = \mathcal{H}_1$ if $t = 1$. Following Gelman and Meng [18], we have

$$\log B_{10} = \log \frac{F(1)}{F(0)} = \int_0^1 \mathbf{E}\{H(\mathbf{Z}^*, \boldsymbol{\eta}, \mathbf{X}, \mathbf{W}, \boldsymbol{\vartheta}_t, t)\} dt,$$

where \mathbf{E} represents the expectation taken with respect to the joint conditional probability density $p(\mathbf{Z}^*, \boldsymbol{\eta}, \boldsymbol{\vartheta}_t | \mathbf{X}, \mathbf{W}, t)$, $H(\mathbf{Z}^*, \boldsymbol{\eta}, \mathbf{X}, \mathbf{W}, \boldsymbol{\vartheta}_t, t) = d \log p(\mathbf{Z}^*, \boldsymbol{\eta}, \mathbf{X}, \mathbf{W}, t | \boldsymbol{\vartheta}_t) / dt$. Then, $\log B_{10}$ can be estimated by

$$\widehat{\log B_{10}} = \frac{1}{2} \sum_{\ell=0}^h (t_{(\ell+1)} - t_{(\ell)}) (\bar{H}_{(\ell+1)} + \bar{H}_{(\ell)}),$$

where $0 = t_{(0)} < t_{(1)} < \dots < t_{(h)} < t_{(h+1)} = 1$, $\bar{H}_{(\ell)} = \mathcal{M}^{-1} \sum_{m=1}^{\mathcal{M}} H(\mathbf{Z}^*, \boldsymbol{\eta}^{(m)}, \boldsymbol{\vartheta}_{t_{(\ell)}}^{(m)}, t_{(\ell)}, \mathbf{X}, \mathbf{W})$ and $\{(\boldsymbol{\eta}^{(m)}, \boldsymbol{\vartheta}_{t_{(\ell)}}^{(m)}) : m = 1, \dots, \mathcal{M}\}$ are the observations generated from the joint conditional probability density $p(\mathbf{Z}^*, \boldsymbol{\eta}, \boldsymbol{\vartheta}_t | \mathbf{X}, \mathbf{W}, t_{(\ell)})$ via the preceding presented Gibbs sampler.

3.5 Goodness-of-fit statistic

To assess the plausibility of the posited model, the posterior predictive (PP) p -value [43] is here adopted. Following Gelman, Meng and Stern [43], the PP p -value is defined as

$$p_B = P\{D(\mathbf{Z}_{\text{rep}} | \boldsymbol{\Theta}) \geq D(\mathbf{Z} | \boldsymbol{\Theta}) | \mathbf{X}, \mathbf{W}, H_0\},$$

where \mathbf{Z}_{rep} denotes a replication of \mathbf{Z} , H_0 represents the plausibility of the above considered model, $D(\cdot | \cdot)$ is a discrepancy variable, and $\boldsymbol{\Theta} = \{\boldsymbol{\gamma}, \boldsymbol{\eta}, \sigma_\zeta^{-2}, \tau^{-2}, \boldsymbol{\varphi}, \boldsymbol{\beta}, \mathbf{a}, \boldsymbol{\lambda}\}$. For

our considered model H_0 , we take the discrepancy variable as

$$D(\mathbf{Z}_{\text{rep}} | \boldsymbol{\Theta}) = \sum_{s \in \mathbb{S}} \{f(\mathbf{Z}_{\text{rep}}(\mathbf{s})) - \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} - \mathbf{W}(\mathbf{s})^\top \boldsymbol{\eta} - \xi(\mathbf{s})\}^2 / \sigma_\varepsilon^2,$$

which is asymptotically distributed as the chi-squared distribution with N degrees of freedom. Thus, the PP p -value can be rewritten as

$$p_B = \int P(\chi^2(N) \geq D(\mathbf{Z} | \boldsymbol{\Theta})) p(\boldsymbol{\Theta} | \mathbf{X}, \mathbf{W}) d\boldsymbol{\Theta}.$$

It is rather challenging to evaluate the above integral. To solve the issue, the commonly used Monte Carlo method is adopted to approximate p_B , i.e.,

$$\hat{p}_B = \frac{1}{\mathcal{M}} \sum_{m=1}^{\mathcal{M}} P(\chi^2(N) \geq D(\mathbf{Z} | \boldsymbol{\Theta}^{(m)})),$$

where the observations $\{\boldsymbol{\Theta}^{(m)} : m = 1, \dots, \mathcal{M}\}$ are generated from the joint probability density $p(\boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\beta}, \sigma_\zeta^{-2}, \tau^{-2}, \boldsymbol{\varphi}, \mathbf{a}, \boldsymbol{\lambda} | \mathbf{Z}, \mathbf{X}, \mathbf{W})$ via the preceding presented Gibbs sampler. Following Lee and Tang [29], the model is plausible if \hat{p}_B is not far from 0.5 (e.g., within the interval (0.3, 0.7)).

4. NUMERICAL EXAMPLE

4.1 Simulation studies

To investigate the performance of the proposed Bayesian estimation procedure, we conducted the first simulation study. In this simulation study, the dataset $\{\mathbf{Z}^*(\mathbf{s}), \mathbf{X}(\mathbf{s}), \mathbf{W}(\mathbf{s}) : \mathbf{s} \in \mathbb{S}\}$ was generated from the TSMMs defined in Equations (3). Here, covariate $\mathbf{X} = (X_1, X_2, X_3)^\top$ was generated by sampling X_1 from the uniform distribution Uniform(-2,2), X_2 and X_3 from the standard normal distribution; we took $L = 2$ resolution W-wavelets $\mathbf{W}(\cdot)_{q=12}$ as the spatial basis functions; $\boldsymbol{\eta}$ was generated from the multivariate normal distribution $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$, ξ was sampled from the normal distribution $\mathcal{N}(0, \sigma_\xi^2)$; ε was sampled from the normal distribution $\mathcal{N}(0, \sigma_\varepsilon^2)$; and we took $\mathbb{S} = \{s : s = 1, \dots, N\}$ with $N = \text{longitude} \times \text{latitude}$, where longitude and latitude were taken as 20, respectively. The true values of parameters $\boldsymbol{\gamma}$, σ_ζ^2 , σ_ε^2 and $\boldsymbol{\Sigma}$ were taken as $\boldsymbol{\gamma} = (-0.5, 0.5, 0.5)^\top$, $\sigma_\zeta^2 = 0.001$, $\sigma_\varepsilon^2 = 1$ and $\boldsymbol{\Sigma} = g\boldsymbol{\Sigma}^0$, respectively, where $g = N(1 - \sigma_\xi^2) / \text{tr}(\mathbf{W}^\top \boldsymbol{\Sigma}^0 \mathbf{W})$, and $\boldsymbol{\Sigma}^0 = \mathbf{R}^{-1} \mathbf{Q}^\top \boldsymbol{\Delta}^0 \mathbf{Q} \mathbf{R}^{-\top}$ in which $\boldsymbol{\Delta}^0 = (\delta_{ij}^0)$ was a $N \times N$ stationary covariance matrix with $\delta_{ij}^0 = (1 - \sigma_\xi^2) \exp(-|i - j|/25)$, and \mathbf{Q} and \mathbf{R} were the QR decompositions of matrix \mathbf{W} (i.e., $\mathbf{W} = \mathbf{Q}\mathbf{R}$).

To investigate the effect of different transformation functions on Bayesian estimations of parameters of interest, we considered the following three types of functions for $f(\mathbf{Z}(\mathbf{s}))$:

Type A. $f^{-1}(\mathbf{Z}^*(\mathbf{s})) = \exp(0.5\mathbf{Z}^*(\mathbf{s}) - 1.0) / (1.0 + \exp(0.5\mathbf{Z}^*(\mathbf{s}) - 1.0))$ leading to highly skewed data;

Table 1. Performance of Bayesian estimates for different transformation functions (TFs) in the first simulation study

TF	γ_1		γ_2		γ_3		
	Bias	RMS	Bias	RMS	Bias	RMS	
A	PS	-0.017	0.100	0.013	0.118	0.024	0.102
	BCO	0.367	0.370	-0.365	0.368	-0.365	0.369
	BC1	0.451	0.452	-0.451	0.452	-0.453	0.454
B	PS	0.002	0.074	0.000	0.084	0.002	0.082
	BCO	0.302	0.308	-0.303	0.311	-0.302	0.310
	BC1	0.403	0.403	-0.404	0.404	-0.407	0.405
C	PS	-0.017	0.078	0.019	0.095	0.013	0.080
	BCO	0.387	0.391	-0.386	0.392	-0.384	0.390
	BC1	0.387	0.391	-0.386	0.392	-0.383	0.390

Note: ‘PS’ represents estimation obtained with P-splines method, ‘BCO’ denotes estimation obtained with Box-Cox transformation at the optimal value of ρ , and ‘BC1’ is estimation obtained with Box-Cox transformation at $\rho = 1.0$.

Table 2. MSPEs of Y for different transformation functions (TFs) in the first simulation study

TF	P-splines	Box-Cox
A	0.113	1.056
B	0.066	0.812
C	0.079	1.221

Type B. $f^{-1}(Z^*(\mathbf{s})) = \text{arctan}\{\exp(Z^*(\mathbf{s})) + 0.5\}$ yielding non-symmetrically U-shaped data;

Type C. $f^{-1}(Z^*(\mathbf{s})) = Z^*(\mathbf{s})/5 + \sin((Z^*(\mathbf{s}) - 5)/20) + 2$ leading to bimodal data.

Based on the above generated dataset, the preceding presented MCMC algorithm was adopted to calculate Bayesian estimates of parameters and $\mathcal{M} = 5000$ observations after 5000 burn-in iterations. For the P-splines approximation of $f(Z(\mathbf{s}))$, we fixed the total number of knots to be $n = n_0 + 2n_1 + 1 = 30$, where n_0 is the number of intervals for dividing the domain of x_{\min} and x_{\max} into n_0 intervals, $n_1 = 3$ is the degree of B-splines. Thus, the number of B-splines K was taken to be $K = n_0 + n_1 = 26$. We selected the second order random walk (i.e., $d = 2$) for specifying the prior of β_k 's. The selection of the hyperparameters was given in Appendix. For comparison, we also calculated Bayesian estimates of parameters via the Box-Cox transformation, i.e., $f(Z) = (Z^\rho - 1)/\rho$ for $\rho \neq 0$ and $f(Z) = \log(Z)$ for $\rho = 0$, under the following two cases: (i) $\rho = 1.0$ and (ii) the optimal value of ρ , which is 0.303, 0.343 and 0.788 for Types A, B and C, respectively. Results for 100 replications were given in Table 1, where ‘Bias’ was the difference between the true value and the mean of the estimates based on 100 replications, and ‘RMS’ was the root mean square between the estimates based on 100 replications and its true value. The average mean squared prediction errors (MSPEs) of Y for 100 replications via P-splines and Box-Cox transformation with the optimal value of ρ were presented in Table 2.

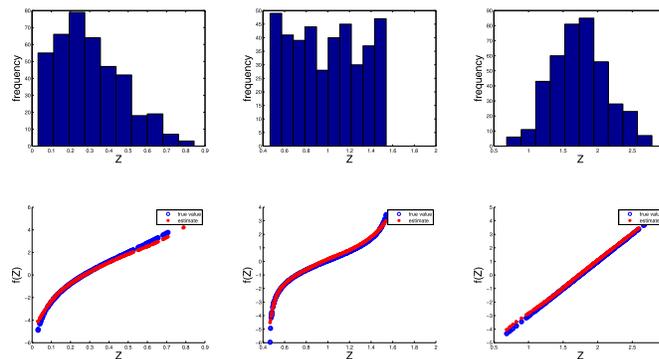


Figure 3. The histogram of Z (upper panel), estimated transformation function via P-spline of $f(Z)$ (lower panel) for type A (left), type B (middle) and type C (right) for a randomly selected replication in the first simulation study.

Examination of Tables 1 and 2 indicated that the above proposed Bayesian estimation procedure performed better than the Box-Cox transformation method, and Box-Cox transformation method with the optimal value of ρ was better than that with a particular value for ρ .

Figure 3 plotted the true and estimated curves of $f(Z)$, Figure 4 plotted the true process of Y and its predicted processes via the P-splines, and Figure 5 plotted the true and estimated covariance matrices of $\boldsymbol{\eta}$ for a randomly selected replication. Inspection of Figures 3 and 4 implied that the preceding proposed Bayesian P-splines method was an effective non-parameter method for estimating unknown transformation function. Examination of Figure 5 showed that the above presented Bayesian method behaved well in estimating covariance matrix of $\boldsymbol{\eta}$.

To illustrate the presented goodness-of-fit statistic, we calculated PP p -value for Types A, B and C, which were 0.435, 0.448 and 0.358 for a randomly selected replication, respectively, which showed that the considered model can fit the generated data well as expected.

To illustrate Bayes factor for comparing two competing spatial mixed models associated with random effects, we conducted the second simulation study. In this simulation study, we considered the following two competing models:

$$\begin{aligned} \mathcal{H}_0 &: f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \nu(\mathbf{s}) + \varepsilon(\mathbf{s}), \\ \mathcal{H}_1 &: f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \nu(\mathbf{s}) + \xi(\mathbf{s}) + \varepsilon(\mathbf{s}). \end{aligned}$$

We simulated the dataset on the basis of the same setting as given in the first simulation study. Thus, \mathcal{H}_1 represented the true model. Similar to Lee and Tang [29], \mathcal{H}_0 and \mathcal{H}_1 was linked by $\mathcal{H}_{t01}: f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \nu(\mathbf{s}) + t\xi(\mathbf{s}) + \varepsilon(\mathbf{s})$, where $t \in [0, 1]$. Clearly, \mathcal{H}_{t01} was equal to \mathcal{H}_0 when $t = 0$, and \mathcal{H}_{t01} was equal to \mathcal{H}_1 when $t = 1$. To calculate logarithm Bayes factor via the path sampling procedure, we took $h = 10$ and $\mathcal{M} = 5000$ after 5000 burn-in iterations in evaluating $\widehat{T}_{(\ell)}$ for $\ell = 0, 1, \dots, h + 1$. The estimated logarithm Bayes factors corresponding to three transformation functions (i.e., Types A, B and C) were $\log \widehat{B}_{10} = 10.186, 4.949, \text{ and } 3.939$,

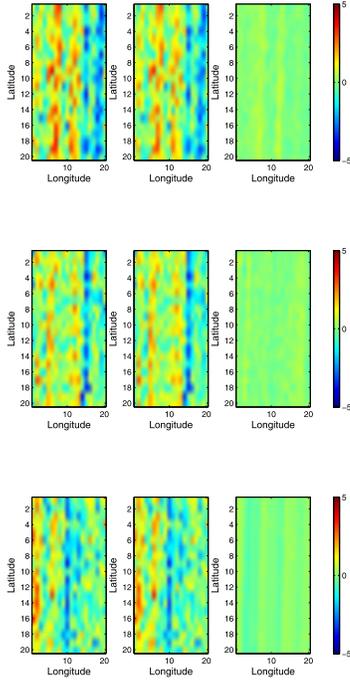


Figure 4. True process of Y (left), predicted process of \hat{Y} with P-spline (middle) and predicted bias of Y (right) for type A (upper), type B (middle) and type C (lower panel) for a randomly selected replication in the first simulation study.

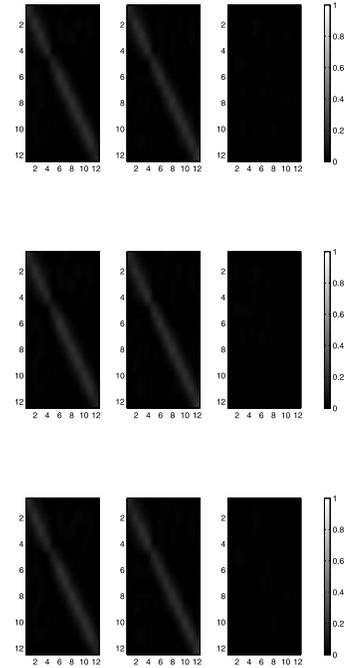


Figure 5. True covariance matrix K of η (left), estimated \hat{K} (middle) and predicted bias of K (right) for type A (upper), type B (middle) and type C (lower panel) in the first simulation study.

respectively. Following Kass and Raftery [24], the true model \mathcal{H}_1 was selected regardless of transformation functions as expected, which implied that the preceding proposed model comparison procedure is quite effective.

To illustrate Bayes factor in selecting the number of spatial basis functions or the number of resolution levels, we conducted the third simulation study. In this simulation study, we considered the following two competing models: $\mathcal{H}_{w0} : f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \mathbf{W}^{(1)}(\mathbf{s})^\top \boldsymbol{\eta}^{(1)} + \mathbf{W}^{(2)}(\mathbf{s})^\top \boldsymbol{\eta}^{(2)} + \xi(\mathbf{s}) + \varepsilon(\mathbf{s})$, $\mathcal{H}_{w1} : f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \mathbf{W}^{(1)}(\mathbf{s})^\top \boldsymbol{\eta}^{(1)} + \xi(\mathbf{s}) + \varepsilon(\mathbf{s})$, where $\mathbf{W}^{(1)}(\mathbf{s}) = (\mathbf{w}_1(\mathbf{s}), \dots, \mathbf{w}_4(\mathbf{s}))^\top$, $\mathbf{W}^{(2)}(\mathbf{s}) = (\mathbf{w}_5(\mathbf{s}), \dots, \mathbf{w}_{12}(\mathbf{s}))^\top$, $\boldsymbol{\eta}^{(1)} = (\eta_1, \dots, \eta_4)^\top$, and $\boldsymbol{\eta}^{(2)} = (\eta_5, \dots, \eta_{12})^\top$. Thus, \mathcal{H}_{w0} was the true model. Similarly, \mathcal{H}_{w0} and \mathcal{H}_{w1} was linked by $\mathcal{H}_{tw01} : f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \mathbf{W}^{(1)}(\mathbf{s})^\top \boldsymbol{\eta}^{(1)} + (1-t)\mathbf{W}^{(2)}(\mathbf{s})^\top \boldsymbol{\eta}^{(2)} + \xi(\mathbf{s}) + \varepsilon(\mathbf{s})$, where $t \in [0, 1]$. Clearly, \mathcal{H}_{tw01} was equal to \mathcal{H}_{w0} when $t = 0$, and \mathcal{H}_{tw01} was equal to \mathcal{H}_{w1} when $t = 1$. We generated the dataset from the model \mathcal{H}_{w0} on the basis of the same setting as given in the first simulation study. The preceding presented path sampling approach with $h = 10$ and $\mathcal{M} = 5000$ after 5000 burn-in iterations was employed to calculate Bayes factors for each of three transformation functions: Types A, B and C. The estimated logarithm Bayes factors corresponding to Types A, B and C transformation functions were $\log \hat{B}_{10} = -202.892, -180.995$ and -152.184 , respectively, which showed that the true model \mathcal{H}_{w0} was selected regardless of transformation functions as expected.

4.2 An example

In this subsection, the remote sensing dataset described in Introduction was used to illustrate the above presented TSMMs together with Bayesian method. As an illustration, we only fitted a small study region \mathbb{S} between longitudes $+30.25$ and $+15.75$, and between latitudes -30.25 and -15.75 to the considered TSMMs. There were some missing data in \mathbb{S} , whose missing proportion was 2.56%. For simplicity, we deleted these missing data. We selected three covariates WV, NR and CM to construct a TSMMs for AOD, i.e., in the considered TSMMs, $Z(\mathbf{s}) = \text{AOD}(\mathbf{s})$, and $\mathbf{X}(\mathbf{s}) = (\text{CM}(\mathbf{s}), \text{NR}(\mathbf{s}), \text{WV}(\mathbf{s}))$. Because the units of covariates and response were inconsistent, the raw data were standardized on the basis of the fully observed data. In this case, $N = 30 \times 30 = 900$ in \mathbb{S} .

The above developed Bayes factor together with the path sampling method was used to select the number of the spatial basis functions. Thus, we considered the following two competing models: $\mathcal{H}_{w0} : f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \mathbf{W}^{(1)}(\mathbf{s})^\top \boldsymbol{\eta}^{(1)} + \mathbf{W}^{(2)}(\mathbf{s})^\top \boldsymbol{\eta}^{(2)} + \xi(\mathbf{s}) + \varepsilon(\mathbf{s})$, $\mathcal{H}_{w1} : f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \mathbf{W}^{(1)}(\mathbf{s})^\top \boldsymbol{\eta}^{(1)} + \xi(\mathbf{s}) + \varepsilon(\mathbf{s})$, where $\boldsymbol{\eta}^{(1)}$ and $\boldsymbol{\eta}^{(2)}$ are 4×1 and 8×1 vectors of random effects corresponding to multi-resolutional spatial basis functions $\mathbf{W}^{(1)}(\mathbf{s})$ and $\mathbf{W}^{(2)}(\mathbf{s})$, respectively. Similarly, \mathcal{H}_{w0} and \mathcal{H}_{w1} was linked by $\mathcal{H}_{tw01} : f(Z(\mathbf{s})) = \mathbf{X}(\mathbf{s})^\top \boldsymbol{\gamma} + \mathbf{W}^{(1)}(\mathbf{s})^\top \boldsymbol{\eta}^{(1)} + (1-t)\mathbf{W}^{(2)}(\mathbf{s})^\top \boldsymbol{\eta}^{(2)} + \xi(\mathbf{s}) + \varepsilon(\mathbf{s})$, where $t \in [0, 1]$. In this case,

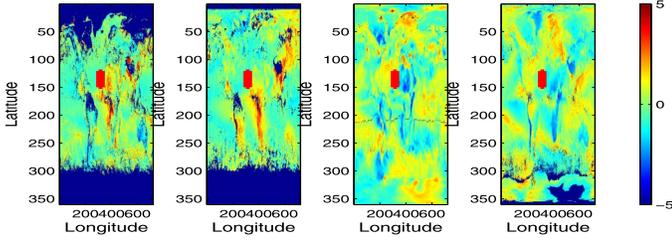


Figure 6. Images of AOD, CM, NR and WV datasets in which the studied datasets are showed in red region.

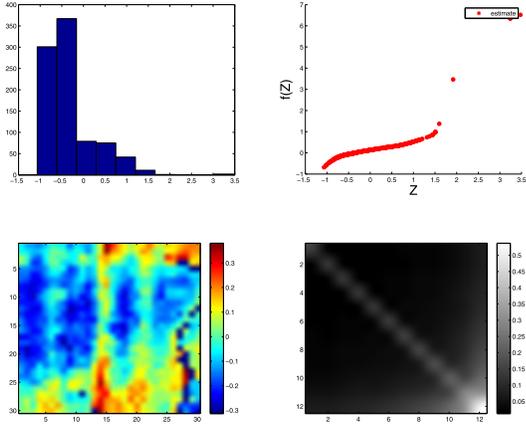


Figure 7. The hist of Z (first), estimated transformation function via P -spline of f (second), predicted process of \hat{Y} with P -spline (third), estimated covariance matrix K of η (fourth).

\mathcal{H}_{tw01} reduced to \mathcal{H}_{w0} when $t = 0$, and \mathcal{H}_{tw01} was \mathcal{H}_{w1} when $t = 1$. Similar to the third simulation study, we took $L = 2$, $h = 10$ and $\mathcal{M} = 5000$ observations after 5000 burn-in iterations in computing $\bar{T}_{(\ell)}$ for $\ell = 0, 1, \dots, h + 1$. The estimated logarithm Bayes factor was $\log \widehat{B}_{10} = -12.741$, which showed that we should select the model \mathcal{H}_{w0} to fit the considered dataset, i.e., $q = 12$.

Based on the selected model \mathcal{H}_{w0} , we obtained Bayesian estimates of γ : $\hat{\gamma} = (-0.2006, -0.183, 0.248)$, which showed that CM and NR had a negative effect on AOD, while WV had a positive effect on AOD. Also, the PP p -values was 0.312, which showed that the selected model fitted the data well. Figure 7 showed that the histogram of the considered dataset, the estimated transformation function of $f(Z)$, the predicted process of \hat{Y} , and the estimated covariance matrix K of η .

APPENDIX A

A.1 Specification of hyperparameters

Generally, one takes $e_2 = 1$ and a small value of e_1 such as $e_1 = 0.005$ leading to an almost diffuse prior for τ^2 . Also, we

set $b = 1$. Following Kang and Cressie [23], we take $\mu_{\gamma_0} = \mathbf{0}$ and $\sigma_{\gamma_0}^2 = 0.25$, $\hat{\sigma}_{\zeta}^2$ as an initial estimate of σ_{ζ}^2 , $\kappa_{\zeta} = \zeta \hat{\sigma}_{\zeta}^2$ in which $\zeta = 8$ in simulation studies and an example analysis. For the prior of Σ , we choose $v_{0i} = 8$, $u_{0i} = 1$, $\mathbf{a}_{j0} = 1$, $\mathbf{A}_{j0} = 0.25\mathbf{I}$ in which \mathbf{I} is an identity matrix.

A.2 MCMC algorithm

Step (a). The conditional distribution for γ is given by

$$p(\gamma|\xi, \eta, \sigma_{\varepsilon}^{-2}, \beta, \mathcal{D}) \sim \mathcal{N}_{\gamma}(\mu_{\gamma}, \Sigma_{\gamma}),$$

where $\mu_{\gamma} = \Sigma_{\gamma}[\sigma_{\varepsilon}^{-2}\mathbf{X}^{\top}\{\sum_{k=1}^K \beta_k \mathbf{B}_k(\mathbf{Z}) - \mathbf{W}\eta - \xi\} + \sigma_{\gamma_0}^{-2}\mu_{\gamma_0}]$, $\Sigma_{\gamma} = (\mathbf{X}^{\top}\mathbf{X}/\sigma_{\varepsilon}^2 + \sigma_{\gamma_0}^{-2}\mathbf{I}_p)^{-1}$, $\mathbf{X} = (\mathbf{X}(s_1), \dots, \mathbf{X}(s_N))^{\top}$, $\mathbf{B}_k(\mathbf{Z}) = (B_k(Z(s_1)), \dots, B_k(Z(s_N)))^{\top}$, and $\mathbf{W} = (\mathbf{W}(s_1), \dots, \mathbf{W}(s_N))^{\top}$.

Step (b). The conditional distribution for η is given by

$$p(\eta|\gamma, \xi, \beta, \mathbf{a}, \lambda, \mathcal{D}) \sim \mathcal{N}_{\eta}(\mu_{\eta}, \Sigma_{\eta}),$$

where $\mu_{\eta} = \Sigma_{\eta}\mathbf{W}\{\sum_{k=1}^K \beta_k \mathbf{B}_k(\mathbf{Z}) - \mathbf{X}\gamma - \xi\}/\sigma_{\varepsilon}^2$, and $\Sigma_{\eta} = (\sigma_{\varepsilon}^{-2}\mathbf{W}\mathbf{W}^{\top} + \Sigma^{-1})^{-1}$.

Step (c). The conditional distribution for ξ is given by

$$p(\xi|\gamma, \eta, \sigma_{\varepsilon}^{-2}, \beta, \mathcal{D}) \sim \mathcal{N}(\mu_{\xi}, \Sigma_{\xi}),$$

where $\mu_{\xi} = \Sigma_{\xi}\{\sum_{k=1}^K \beta_k \mathbf{B}_k(\mathbf{Z}) - \mathbf{X}\gamma - \mathbf{W}\eta\}/\sigma_{\varepsilon}^2$, and $\Sigma_{\xi} = (\sigma_{\varepsilon}^{-2} + \sigma_{\xi}^{-2})^{-1}\mathbf{I}$.

Step (d). It is easily shown from the prior distribution of β and the distribution of $\mathbf{Z}(s)$ that the full conditional distribution β can be expressed as

$$p(\beta|\gamma, \eta, \sigma_{\varepsilon}^{-2}, \varphi, \mathcal{D}) \propto \exp\left\{-\frac{1}{2}(\beta - \mathbf{m})^{\top}\mathbf{P}(\beta - \mathbf{m})\right\} \times \prod_{k=2}^K I(\beta_k \geq \beta_{k-1})I(\mathbf{Q}\beta = 0),$$

where $\mathbf{P} = \mathbf{B}\mathbf{B}^{\top}/\sigma_{\zeta}^2 + \mathbf{M}(\varphi)/\tau^2$ with $\sigma_{\zeta}^2 = \sigma_{\varepsilon}^2 + \sigma_{\xi}^2$, $\mathbf{B} = (B_1(\mathbf{Z}), \dots, B_K(\mathbf{Z}))^{\top}$, and $\mathbf{m} = \mathbf{P}^{-1}\mathbf{B}(\mathbf{X}\gamma + \mathbf{W}\eta)$. Clearly, the conditional distribution $p(\beta|\cdot)$ is not a familiar distribution. Thus, it is impossible to directly sample observation from the conditional distribution $p(\beta|\cdot)$. To this end, the Gibbs sampler is again adopted to sample each component β_k of β as follows.

For $t = 0, 1, 2, \dots$, at the current value $\beta^{(t)} = (\beta_1^{(t)}, \dots, \beta_K^{(t)})^{\top}$ of β with the constrained conditions $\beta_k^{(t)} \geq \beta_{k-1}^{(t)}$ for $k = 2, \dots, K$, the observations $\beta_1^{(t+1)}, \dots, \beta_K^{(t+1)}$ are iteratively drawn by

- (1) sampling $\beta_1^{(t+1)}$ from $TN(\mu_{\beta_1}, \varphi_{\beta_1}^2, -\infty, \beta_2^{(t)})$,
- (2) sampling $\beta_2^{(t+1)}$ from $TN(\mu_{\beta_2}, \varphi_{\beta_2}^2, \beta_1^{(t+1)}, \beta_3^{(t)})$,
- (3) sampling $\beta_3^{(t+1)}$ from $TN(\mu_{\beta_3}, \varphi_{\beta_3}^2, \beta_2^{(t+1)}, \beta_4^{(t)})$,

⋮

(K) sampling $\beta_K^{(t+1)}$ from $TN(\mu_{\beta_K}, \varphi_{\beta_K}^2, \beta_{K-1}^{(t+1)}, \infty)$, where $TN(\mu_d, \varphi_d^2, d_1, d_2)$ represents the one-dimensional

interval truncated Gaussian distribution in the interval (d_1, d_2) with mean μ_d and variance φ_d^2 , $\mu_{\beta_k} = m_k + \mathbf{P}_{[k]}^\top \mathbf{P}_{(k)}^{-1} (\boldsymbol{\beta}_{[k]}^{(t)} - \mathbf{m}_{[k]})$ and $\varphi_{\beta_k}^2 = P_{kk} - \mathbf{P}_{[k]}^\top \mathbf{P}_{(k)}^{-1} \mathbf{P}_{[k]}$ for $k = 1, \dots, K$, where P_{kk} is the k th diagonal element of matrix \mathbf{P} , $\mathbf{P}_{(k)}$ is a $(K-1) \times (K-1)$ submatrix of matrix \mathbf{P} with the components corresponding to the k th row and the k th column of matrix \mathbf{P} deleted, and $\mathbf{P}_{[k]}$ is a $(K-1) \times 1$ vector with the k th component in the k th column of matrix \mathbf{P} deleted, $\mathbf{m}_{[k]}$ and $\boldsymbol{\beta}_{[k]}$ are the $(K-1) \times 1$ vectors with the k th components of vectors \mathbf{m} and $\boldsymbol{\beta}$ deleted, respectively. For identifiability, the function $f(\cdot)$ satisfies the following restriction condition: $\mathbf{Q}\boldsymbol{\beta} = 0$, which can be implemented by centering the function $f(\cdot)$ on its mean at each of iterations of the Gibbs sampler when evaluating estimates of parameters via the above presented MCMC.

Step (e). The conditional distribution for σ_ξ^2 is given by

$$p(\sigma_\xi^2 | \boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\beta}, \mathcal{D}) \sim \text{IG}(a_{\sigma_\xi^2}, b_{\sigma_\xi^2}) I\{\sigma_\xi \in (0, \kappa_\xi)\},$$

where $a_{\sigma_\xi^2} = (N-1)/2$, $b_{\sigma_\xi^2} = (\mathbf{B}^\top \boldsymbol{\beta} - \mathbf{X}\boldsymbol{\gamma} - \mathbf{W}\boldsymbol{\eta})^\top (\mathbf{B}^\top \boldsymbol{\beta} - \mathbf{X}\boldsymbol{\gamma} - \mathbf{W}\boldsymbol{\eta})/2$, and $\text{IG}(\cdot, \cdot)$ represents the inverse-gamma distribution.

Step(f). The conditional distribution $p(\tau^{-2} | \boldsymbol{\beta}, \boldsymbol{\varphi})$ is given by

$$p(\tau^{-2} | \boldsymbol{\beta}, \boldsymbol{\varphi}) \sim \Gamma\left(e_1 + \frac{K-d}{2}, e_2 + \frac{1}{2} \boldsymbol{\beta}^\top \mathbf{M}(\boldsymbol{\varphi}) \boldsymbol{\beta}\right),$$

where d is the order of random walk.

Step(g). The conditional distribution $p(\varphi_k | \boldsymbol{\beta}, \tau^{-2})$ is given by

$$p(\varphi_k | \boldsymbol{\beta}, \tau^{-2}) \sim \Gamma\left(\frac{b+1}{2}, \frac{b + u_k^2/\tau^2}{2}\right),$$

where u_k is equal to $\beta_k - \beta_{k-1}$ and $\beta_k - 2\beta_{k-1} + \beta_{k-2}$ for the first and second order random walk, respectively.

Step (h). To obtain conditional distribution of $\boldsymbol{\lambda}$, it follows from $\boldsymbol{\omega} = \mathbf{A}\boldsymbol{\eta}$ and $|\boldsymbol{\Sigma}|^{-1} = |\mathbf{A}|^\top |\boldsymbol{\Lambda}|^{-1} |\mathbf{A}| = |\boldsymbol{\Lambda}^{-1}| = \prod_{i=1}^q \lambda_i^{-1}$ that the likelihood can be written as

$$\begin{aligned} \ell(\boldsymbol{\eta} | \boldsymbol{\Sigma}) &\propto |\boldsymbol{\Sigma}|^{-1/2} \exp\left\{-\frac{1}{2} \boldsymbol{\eta}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\eta}\right\} \\ &= \prod_{i=1}^q \lambda_i^{-1/2} \exp\left\{-\frac{1}{2} \boldsymbol{\omega}^\top \boldsymbol{\Lambda}^{-1} \boldsymbol{\omega}\right\} \\ &= \prod_{i=1}^q \lambda_i^{-1/2} \exp\left\{-\frac{1}{2} \text{tr}(\boldsymbol{\Lambda}^{-1} \boldsymbol{\omega} \boldsymbol{\omega}^\top)\right\} \\ &= \prod_{i=1}^q \lambda_i^{-1/2} \exp\left\{-\frac{s_i}{2\lambda_i}\right\} \end{aligned}$$

The conditional distribution for the i th component of $\boldsymbol{\lambda}$ is given by

$$p(\lambda_i^{-1} | \boldsymbol{\eta}, \mathbf{a}) \sim \Gamma\left(\frac{v_{0i} + 1}{2}, \frac{u_{0i} + s_i}{2}\right),$$

where s_i is the i th diagonal element of matrix $\boldsymbol{\omega} \boldsymbol{\omega}^\top$ for $i = 1, \dots, q$.

Step (i). Consider the conditional distribution of \mathbf{a} . Note that

$$\begin{aligned} \mathbf{A}\boldsymbol{\eta} &= \begin{pmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ a_{21} & 1 & 0 & \dots & 0 & 0 \\ a_{31} & a_{32} & 1 & \dots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ a_{q1} & a_{q2} & \dots & \dots & a_{q,q-1} & 1 \end{pmatrix} \begin{pmatrix} \eta_1 \\ \eta_2 \\ \eta_3 \\ \vdots \\ \eta_q \end{pmatrix} \\ (6) \quad &= \begin{pmatrix} \eta_1 \\ \eta_2 + a_{21}\eta_1 \\ \eta_3 + a_{31}\eta_1 + a_{32}\eta_2 \\ \vdots \\ \eta_q + \sum_{i=1}^{q-1} a_{qi}\eta_i \end{pmatrix}. \end{aligned}$$

Thus, we can rewrite the likelihood function as

$$\begin{aligned} \ell(\boldsymbol{\eta} | \boldsymbol{\Sigma}) &\propto \prod_{i=1}^q \lambda_i^{-1/2} \exp\left\{-\frac{1}{2} (\mathbf{A}\boldsymbol{\eta})^\top \boldsymbol{\Lambda}^{-1} (\mathbf{A}\boldsymbol{\eta})\right\} \\ &= \prod_{i=1}^q \lambda_i^{-1/2} \exp\left\{-\frac{\eta_i^2}{2\lambda_i}\right\} \exp\left\{-\frac{(\eta_2 + a_{21}\eta_1)^2}{2\lambda_2}\right\} \\ &\quad \times \exp\left\{-\frac{(\eta_q + \sum_{i=1}^{q-1} a_{qi}\eta_i)^2}{2\lambda_q}\right\}, \end{aligned}$$

which leads to

$$\ell(\boldsymbol{\eta} | \boldsymbol{\Sigma}) \propto \prod_{i=1}^q \lambda_i^{-1/2} \exp\left\{-\frac{1}{2} (\boldsymbol{\eta} - \nabla \mathbf{a})^\top \boldsymbol{\Lambda}^{-1} (\boldsymbol{\eta} - \nabla \mathbf{a})\right\},$$

where

$$\nabla = - \begin{pmatrix} 0 & 0 & 0 & \dots & 0 & 0 \\ \eta_1 & 0 & 0 & \dots & 0 & 0 \\ 0 & \eta_1 & \eta_2 & \dots & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \eta_1 & \dots & \eta_{q-1} \end{pmatrix}.$$

Therefore, the conditional distribution for \mathbf{a} is given by

$$p(\mathbf{a} | \boldsymbol{\eta}, \boldsymbol{\Lambda}) \sim \mathcal{N}(\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_a),$$

where $\boldsymbol{\Sigma}_a = (\mathbf{A}_0^{-1} + \nabla \boldsymbol{\Lambda}^{-1} \nabla)^{-1}$, and $\boldsymbol{\mu}_a = \boldsymbol{\Sigma}_a (\mathbf{A}_0^{-1} \mathbf{a}_0 + \nabla \boldsymbol{\Lambda}^{-1} \boldsymbol{\eta})$. Because $\boldsymbol{\Lambda}$ is a diagonal matrix, when \mathbf{A}_0 is also a diagonal or block-diagonal matrix corresponding to rows of \mathbf{A} , the derivations are simplified and the elements of \mathbf{a} can be updated in a series of independent steps

$$p(\mathbf{a}_i | \boldsymbol{\eta}, \boldsymbol{\Lambda}_i) \sim \mathcal{N}(\boldsymbol{\mu}_{ai}, \boldsymbol{\Sigma}_{ai}),$$

where $\boldsymbol{\Sigma}_{ai} = (\mathbf{A}_{i0}^{-1} + \lambda_i^{-1} \boldsymbol{\Omega}_i^\top \boldsymbol{\Omega}_i)^{-1}$, $\boldsymbol{\mu}_{ai} = \boldsymbol{\Sigma}_{ai} (\mathbf{A}_{i0}^{-1} \mathbf{a}_{i0} - \lambda_i^{-1} \boldsymbol{\Omega}_i^\top \boldsymbol{\eta}_i)$, $\boldsymbol{\Omega}_i = (\eta_1, \dots, \eta_{i-1})$ for $i = 2, \dots, q$. Note that \mathbf{A}_{i0} may depend on λ_i . Given the posterior draws \mathbf{a} and $\boldsymbol{\lambda}$, a posterior draw of $\boldsymbol{\Sigma}$ can be obtained by $\boldsymbol{\Sigma} = \mathbf{A}^{-1} \boldsymbol{\Lambda} \mathbf{A}^{-\top}$.

The conditional distributions corresponding to steps (a)-(i) are some familiar distributions, thus sampling observations from these distributions is straightforward.

- [20] GUINNESS, J. (2017). Spectral density estimation for random fields via periodic embeddings. *Journal of Computational and Graphical Statistics* **24**, 579–599.
- [21] GURKA, M. J. (2006). Selecting the best linear mixed model under REML. *American Statistician* **60**, 19–26. [MR2224133](#)
- [22] KANG, E. L. (2009). *Reduced-Dimension Hierarchical Statistical Models for Spatial and Spatio-Temporal Data*. Ph.D. thesis, The Ohio State University, Columbus, OH.
- [23] KANG, E. L. and CRESSIE, N. (2011). Bayesian Inference for the Spatial Random Effects Model. *Journal of the American Statistical Association* **106**, 972–983. [MR2894757](#)
- [24] KASS, R. E. and RAFTERY, A. E. (1995). Bayes factors. *Journal of the American Statistical Association* **90**, 773–795. [MR3363402](#)
- [25] KATZFUSS, M. (2015). A multi-resolution approximation for massive spatial datasets. *Journal of the American Statistical Association* **112**, 201–204. [MR3646566](#)
- [26] KWONG, M. K. and TANG, P. T. P. (1994). W-matrices, nonorthogonal multiresolution analysis, and finite signals of arbitrary length. Technical report, Argonne National Laboratory, 449-0794.
- [27] LANG, S. and BREZGER, S. (2004). Bayesian P-splines. *Journal of Computational and Graphical Statistics* **13**, 183–212. [MR2044877](#)
- [28] LEE, S. Y. and SONG, X. Y. (2003). Model comparison of nonlinear structural equation models with fixed covariates. *Psychometrika* **68**, 27–47. [MR2272368](#)
- [29] LEE, S. Y. and TANG, N. S. (2006). Bayesian analysis of nonlinear structural equation models with nonignorable missing data. *Psychometrika* **71**, 541–564. [MR2272542](#)
- [30] MENG, X. L. and WONG, W. H. (1996). Simulating ratios of normalizing constants via a simple identity: a theoretical exploration. *Statistica Sinica* **6**, 831–860. [MR1422406](#)
- [31] NYCHKA, D., BANDYOPADHYAY, S., HAMMERLING, D., LINDGREN, F. and SAIN, S. (2015). A multiresolution Gaussian process model for the analysis of large spatial datasets. *Journal of Computational and Graphical Statistics* **24**, 579–599. [MR3357396](#)
- [32] NYCHKA, D., WIKLE, C. and ROYLE, J. A. (2002). Multiresolution Models for Nonstationary Spatial Covariance Functions. *Statistical Modeling* **2**, 315–331. [MR1951588](#)
- [33] ROBERT, C. P., (1995). Simulation of truncated normal variables. *Statistics and Computing* **5**, 121–125.
- [34] SHI, T. and CRESSIE, N. (2007). Global statistical analysis of MISR aerosol data: a massive data product from NASA's terra satellite. *Environmetrics* **18**, 665–680. [MR2408937](#)
- [35] SONG, X. Y. and LIU, P. F. (2015). Transformation structural equation models with highly non-normal and incomplete Data. *Structural Equation Modeling: A Multidisciplinary Journal* **22**, 401–415. [MR3357141](#)
- [36] SONG, X. Y. and LU, Z. H. (2012). Semi-parametric transformation models with Bayesian P-splines. *Statistics and Computing* **22**, 1085–1098. [MR2950087](#)
- [37] TANG, N., WU, Y. and CHEN, D. (2018). Semiparametric Bayesian analysis of transformation linear mixed models. *Journal of Multivariate Analysis* **166**, 225–240. [MR3799645](#)
- [38] WIKLE, C. K., MILLIFF, R., NYCHKA, D. W. and BERLINER, L. M. (2001). Spatiotemporal hierarchical Bayesian modeling: Tropical ocean surface winds. Technical report, Argonne National Laboratory, 449-0794. [MR1939342](#)
- [39] XU, B., WIKLE, C. K. and FOX, N. (2005). A kernel-based spatio-temporal dynamical model for nowcasting radar precipitation. *Journal of the American Statistical Association* **100**, 1133–1144. [MR2236929](#)
- [40] ZAREIFARD, H. and KHALEDI, M. J. (2013). Non-Gaussian modeling of spatial data using scale mixing of a unified skew Gaussian process. *Journal of Multivariate Analysis* **114**, 16–28. [MR2993870](#)
- [41] ZHANG, H. and EI-SHAARAWI, A. (2010). On spatial skew-Gaussian processes and applications. *Environmetrics* **21**, 33–47. [MR2842222](#)
- [42] ZHU, H. T., IBRAHIM, J. G., CHI, Y. Y. and TANG, N. S. (2012). Bayesian influence measures for joint models for longitudinal and survival data. *Biometrics* **68**, 954–964. [MR3055200](#)
- [43] GELMAN, A., MENG, X.L., STERN, H. (1996). Posterior predictive assessment of model fitness via realized discrepancies. *Statistical Sinical* **6**, 733–807. [MR1422404](#)

Ying Wu
Yunnan Key Laboratory of Statistical Modeling
and Data Analysis
Yunnan University, Kunming 650091
P. R. of China
E-mail address: wy@cxtc.edu.cn

Dan Chen
Yunnan Key Laboratory of Statistical Modeling
and Data Analysis
Yunnan University, Kunming 650091
P. R. of China
E-mail address: danchen@ynu.edu.cn

Niansheng Tang
Yunnan Key Laboratory of Statistical Modeling
and Data Analysis
Yunnan University, Kunming 650091
P. R. of China
E-mail address: nstang@ynu.edu.cn