

A comparison of centring parameterisations of Gaussian process-based models for Bayesian computation using MCMC

Mark R. Bass¹  · Sujit K. Sahu¹

Received: 18 February 2016 / Accepted: 30 August 2016
© The Author(s) 2016. This article is published with open access at Springerlink.com

Abstract Markov chain Monte Carlo (MCMC) algorithms for Bayesian computation for Gaussian process-based models under default parameterisations are slow to converge due to the presence of spatial- and other-induced dependence structures. The main focus of this paper is to study the effect of the assumed spatial correlation structure on the convergence properties of the Gibbs sampler under the default non-centred parameterisation and a rival centred parameterisation (CP), for the mean structure of a general multi-process Gaussian spatial model. Our investigation finds answers to many pertinent, but as yet unanswered, questions on the choice between the two. Assuming the covariance parameters to be known, we compare the exact rates of convergence of the two by varying the strength of the spatial correlation, the level of covariance tapering, the scale of the spatially varying covariates, the number of data points, the number and the structure of block updating of the spatial effects and the amount of smoothness assumed in a Matérn covariance function. We also study the effects of introducing differing levels of geometric anisotropy in the spatial model. The case of unknown variance parameters is investigated using well-known MCMC convergence diagnostics. A simulation study and a real-data example on modelling air pollution levels in London are used for illustrations. A generic pattern emerges that the CP is preferable in the presence of more spatial correlation or more information obtained through, for example, additional data points or by increased covariate variability.

Keywords Bayesian inference · Gibbs sampler · Hierarchical models · Rate of convergence · Spatial data

✉ Sujit K. Sahu
S.K.Sahu@soton.ac.uk

¹ Mathematical Sciences, University of Southampton,
Southampton SO17 1BJ, UK

1 Introduction

Spatially correlated data are prevalent in many of the physical, biological and environmental sciences. It is natural to model these processes in a Bayesian modelling framework, employing Markov chain Monte Carlo (MCMC) techniques for model fitting and prediction, in particular Gibbs sampling type algorithms (Gelfand and Smith 1990). There is a growing interest among researchers in regression models with spatially varying coefficients (Gelfand et al. 2003). Fitting these highly overparameterised and nonstationary models is challenging and computationally expensive. Latent processes correlated across space produce dense covariance matrices that require calculations of order $O(n^3)$ to invert, for n spatial locations (Cressie and Johannesson 2008).

For normal linear hierarchical models with independent random effects, it is known that the ratio of the variance parameters determines the convergence rates of the Gibbs samplers (Papaspiliopoulos et al. 2003; Gelfand et al. 1995). When the data precision is relatively high, the centred parameterisation (CP) will yield an efficient Gibbs sampler, and when the data precision is relatively low, the non-centred parameterisation (NCP) is more efficient. Papaspiliopoulos et al. (2003) find that the NCP outperforms the CP for a Cauchy data model with Gaussian latent variables. Papaspiliopoulos and Roberts (2008) further investigate how the model parameterisation and the tail behaviour of the distributions of the data and the latent process all interact to determine the stability of the Gibbs sampler. They look at combinations of Cauchy, double exponential, Gaussian and exponential power distributions for the CP and the NCP. The heuristic remark that follows from this comparison is that the convergence of the CP is quicker when the data model has lighter tails than that of the latent variables, with the opposite scenario favouring the NCP.

There has been little investigation into the influence of correlation across the random effects on the rate of convergence of the Gibbs sampler. Simulation studies conducted by Papaspiliopoulos et al. (2003) on the spatial Poisson-log-Normal model suggest that stronger spatial correlation improves the sampling efficiency of the CP relative to that of the NCP. However, there are several unresolved questions regarding the choice of the CP versus NCP for the mean structure of a general multi-process Gaussian spatial model. Which of the two parameterisations will converge faster when spatial correlation is increased? What happens to the rates of convergence when tapering (Furrer et al. 2006; Kaufman et al. 2008) is introduced? How does the smoothness parameter in an assumed Matérn covariance function influence the rates? In addition, there are other unexplored issues regarding the choice and the number of blocks for the random effects, the influence of the scale of spatially varying covariates and the introduction of different levels of geometric anisotropy.

In this paper, we cast the general spatial model with multiple spatially varying covariates as a three stage normal linear hierarchical model. This model formulation allows us to compute the exact rates of convergence for both CP and NCP for known prior covariance matrices by following (Roberts and Sahu 1997). These exact rates of convergence facilitate comparison between the two rival parameterisations, CP and NCP. For an exponential correlation function, convergence for the CP is hastened when spatial correlation is stronger, the opposite being true for the NCP. This is also demonstrated in the context of tapered covariance matrices, geometric anisotropic correlation functions and the regression process associated with a spatially varying covariate. The exponential correlation function is a member of the broader Matérn family (Matérn 1986). When we increase the smoothness parameter, the effect is to slow convergence for both the CP and NCP. In the limiting case of the Gaussian correlation, when the smoothness parameter tends to infinity, both CP and NCP may fail to converge when the sample size is large enough due to the associated singularity of the covariance matrices, and Sect. 3.3 investigates the issues.

When the prior covariance matrices are unknown, the exact convergence is intractable and so the CP and NCP are compared by statistics based on well-known convergence diagnostics on the potential scale reduction factor, see e.g., Gelman and Rubin (1992). We use a simulation and a real-data example to show that increasing the effective range for an exponential correlation function improves the sampling efficiency of the CP, whereas shortening the effective range helps the NCP.

The following remarks are in order. First, a related approach is to marginalise over the random effects, thus reducing the dimension of the posterior distribution. This approach can be employed when the error structures of the

data and the random effects are both assumed to be Gaussian. Marginalised likelihoods are used by Gelfand et al. (2003) for fitting spatially varying coefficient models and by Banerjee et al. (2008) to implement Gaussian predictive process models. However, marginalisation results in a loss of conditional conjugacy of the variance parameters and means that they have to be updated by using Metropolis-type steps, which require difficult and time consuming tuning. On the other hand, the Gibbs sampler for the full model can potentially be completely automated and run without the need for any tuning.

Secondly, it is possible to generate intermediate partially centred parameterisations by considering CP and the NCP as extremes of a family of parameterisations. Indeed, this has been followed up by Bass and Sahu (2016) in a companion paper. Interweaving of the CP and NCP as proposed by Yu and Meng (2011) is particularly useful when the practitioner has little knowledge of the convergence properties of either parameterisation. These authors obtain an upper bound on the convergence rate of the interweaving algorithm based on an intractable maximal correlation between the latent variables under the two parameterisations and to our knowledge the exact convergence rate for the interweaving algorithm has not yet been computed. Lastly, both CP- and NCP-based computation methods are similar in spirit to various data augmentation (DA) schemes (Liu and Wu 1999; van Dyk and Meng 2001; Imai and van Dyk 2005; Filippone et al. 2013). A direct theoretical comparison between the exact convergence rates of the centring parameterisations for Gaussian process (GP)-based models and the DA algorithms is desirable, as has been done by Sahu and Roberts (1999) for the EM algorithm and the Gibbs sampler. However, this requires further methodological development of DA algorithms for the GP models and evaluation of their exact rates of convergence.

The rest of this paper is organized as follows. In Sect. 2, we give details of a general spatial model and obtain expressions for the rates of convergence. A simple example here illustrates the rates and brings out the rivalry between the two parameterisations. Section 3 is devoted to comparison of the rates of convergence under different settings of correlation structures and introduction of geometric anisotropy. It also studies the effect of tapering and the scale of the spatially varying covariates on the rates of convergence. In Sect. 4, we drop the assumption of known precision matrices and use two convergence diagnostics to judge the sampling efficiencies of the CP and NCP methods. Analyses are carried out on simulated data and PM10 concentration data taken from Greater London in 2011. Section 5 contains some concluding remarks. Appendices A and B, respectively, contain the technical details for calculating the rates of convergence and the full conditional distributions needed for Gibbs sampling.

2 General spatial model

2.1 Model specification

For data observed at a set of locations s_1, \dots, s_n , we consider the following normal linear model with spatially varying regression coefficients (Gelfand et al. 2003):

$$Y(s_i) = \sum_{k=0}^{p-1} \{\theta_k + \beta_k(s_i)\} x_k(s_i) + \epsilon(s_i) \quad (i = 1, \dots, n). \quad (1)$$

We model (measurement or micro-scale) errors $\epsilon(s_i)$ as independent and normally distributed with mean zero and variance σ_ϵ^2 . Spatially indexed observations $\mathbf{Y} = \{Y(s_1), \dots, Y(s_n)\}^T$ are conditionally independent and normally distributed as

$$Y(s_i) \sim N(\mathbf{x}^T(s_i)\{\boldsymbol{\theta} + \boldsymbol{\beta}(s_i)\}, \sigma_\epsilon^2),$$

where $\mathbf{x}(s_i) = \{1, x_1(s_i), \dots, x_{p-1}(s_i)\}^T$ is a vector containing covariate information for site s_i and $\boldsymbol{\theta} = (\theta_0, \dots, \theta_{p-1})^T$ is a vector of global regression coefficients. The k th element of $\boldsymbol{\theta}$ is locally perturbed by a realisation of a zero mean independent Gaussian process, denoted $\beta_k(s_i)$, which are collected into a vector $\boldsymbol{\beta}(s_i) = \{\beta_0(s_i), \dots, \beta_{p-1}(s_i)\}^T$. The n realisations of the Gaussian process associated with the k th covariate are given by $\boldsymbol{\beta}_k = \{\beta_k(s_1), \dots, \beta_k(s_n)\}^T \sim N(0, \boldsymbol{\Sigma}_k)$ ($k = 0, \dots, p-1$), where $\boldsymbol{\Sigma}_k = \sigma_k^2 \mathbf{R}_k$, and $(\mathbf{R}_k)_{ij} = \text{corr}\{\beta_k(s_i), \beta_k(s_j)\}$. The form of the model given in (1) is known as the NCP. The CP is found by introducing the variables $\tilde{\beta}_k(s_i) = \theta_k + \beta_k(s_i)$. Therefore, $\tilde{\boldsymbol{\beta}}_k = \{\tilde{\beta}_k(s_1), \dots, \tilde{\beta}_k(s_n)\}^T \sim N(\theta_k \mathbf{1}, \boldsymbol{\Sigma}_k)$.

Global effects $\boldsymbol{\theta}$ are assumed to be multivariate normal *a priori* and so we write model (1) in its hierarchically centred form as

$$\mathbf{Y}|\tilde{\boldsymbol{\beta}} \sim N(\mathbf{X}_1 \tilde{\boldsymbol{\beta}}, \mathbf{C}_1), \quad \tilde{\boldsymbol{\beta}}|\boldsymbol{\theta} \sim N(\mathbf{X}_2 \boldsymbol{\theta}, \mathbf{C}_2), \quad \boldsymbol{\theta} \sim N(\mathbf{m}, \mathbf{C}_3),$$

where $\mathbf{C}_1 = \sigma_\epsilon^2 \mathbf{I}$ and $\mathbf{X}_1 = (\mathbf{I}, \mathbf{D}_1, \dots, \mathbf{D}_{p-1})$ is the $n \times np$ design matrix for the first stage where \mathbf{D}_k is a diagonal matrix with entries $\mathbf{x}_k = \{x_k(s_1), \dots, x_k(s_n)\}^T$. We denote by $\tilde{\boldsymbol{\beta}} = (\tilde{\boldsymbol{\beta}}_0^T, \dots, \tilde{\boldsymbol{\beta}}_{p-1}^T)^T$ the $np \times 1$ vector of centred, spatially correlated random effects.

The design matrix for the second stage, \mathbf{X}_2 , is a $np \times p$ block diagonal matrix, the blocks made of vectors of ones of length n . The p processes are assumed independent *a priori* and so \mathbf{C}_2 is block diagonal where the k th block is $\boldsymbol{\Sigma}_k$.

2.2 Prior distributions

The global effects $\boldsymbol{\theta} = (\theta_0, \theta_1, \dots, \theta_{p-1})^T$ are assumed to be independent *a priori* with the k th element assigned an

independent Gaussian prior distribution with mean m_k and variance $\sigma_k^2 v_k$, and hence, we write $\theta_k \sim N(m_k, \sigma_k^2 v_k)$ for $k = 0, \dots, p-1$. Therefore, $\mathbf{m} = (m_0, \dots, m_{p-1})^T$ and \mathbf{C}_3 is a diagonal matrix with diagonal entries $\sigma_k^2 v_k$.

The realisations of the k th non-centred Gaussian process, $\boldsymbol{\beta}_k$, have a prior covariance matrix $\boldsymbol{\Sigma}_k = \sigma_k^2 \mathbf{R}_k$. This prior covariance matrix is shared by the k th centred Gaussian process, $\tilde{\boldsymbol{\beta}}_k$. The prior distributions for the variance parameters are $\sigma_k^2 \sim IG(a_k, b_k)$ ($k = 0, \dots, p-1$), $\sigma_\epsilon^2 \sim IG(a_\epsilon, b_\epsilon)$, where we write $X \sim IG(a, b)$ if X has a density proportional to $x^{-(a+1)} e^{-b/x}$. The entries of the \mathbf{R}_k are $(\mathbf{R}_k)_{ij} = \text{corr}\{\beta_k(s_i), \beta_k(s_j)\} = \rho_k(d_{ij}; \phi_k, v_k)$ where $d_{ij} = \|s_i - s_j\|$ denotes the distance between s_i and s_j and ρ_k is a correlation function from the Matérn family (Handcock and Stein 1993; Matérn 1986).

The Matérn correlation function for a pair of random variables at sites s_i and s_j is

$$\rho(d_{ij}, \phi, \nu) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\sqrt{2\nu\phi} d_{ij} \right)^\nu K_\nu \left(\sqrt{2\nu\phi} d_{ij} \right), \quad \phi > 0, \nu > 0, \quad (2)$$

where $\Gamma(\cdot)$ is the gamma function and $K_\nu(\cdot)$ is the modified Bessel function of the second kind of order ν (Abramowitz and Stegun 1972, Sect. 9.6). The parameter ϕ controls the rate of decay of correlation between two points as their separation increases. The smoothness of the realised random field is controlled by ν , as the process realisations are $[\nu]$ -times mean-square differentiable. A number of parameterisations of the Matérn correlation function exist, for examples see Schabenberger and Gotway (2004, Sect. 4.7.2). The form given in (2) is taken from Rasmussen and Williams (2006, Sect. 4.2.1) and its special cases for different values of ν are discussed in Sect. 3.2.

2.3 Exact rates of convergence

For a Gibbs sampler with Gaussian target distribution with a known precision matrix, we can compute the exact rate of convergence (Roberts and Sahu 1997). The convergence rate λ is bounded in the interval $[0, 1]$, with $\lambda = 0$ indicating immediate convergence and $\lambda = 1$ indicating sub-geometric convergence (Meyn and Tweedie 1993).

Suppose we block update all random effects $\boldsymbol{\beta}$, or $\tilde{\boldsymbol{\beta}}$ in the case of the CP, and block update all global effects $\boldsymbol{\theta}$. Using the results given in Appendix 1 we can show that the respective rates of convergence for the CP and the NCP of model (1) are given by the maximum modulus eigenvalue of

$$\mathbf{F}^c = \left(\mathbf{X}_2' \mathbf{C}_2^{-1} \mathbf{X}_2 + \mathbf{C}_3^{-1} \right)^{-1} \mathbf{X}_2' \mathbf{C}_2^{-1} \times \left(\mathbf{X}_1' \mathbf{C}_1^{-1} \mathbf{X}_1 + \mathbf{C}_2^{-1} \right)^{-1} \mathbf{C}_2^{-1} \mathbf{X}_2, \quad (3)$$

and the maximum modulus eigenvalue of

$$\begin{aligned} F^{nc} = & \left(X_2' X_1 C_1^{-1} X_1 X_2 + C_3^{-1} \right)^{-1} \\ & \times X_2' X_1 C_1^{-1} X_1 \left(X_1' C_1^{-1} X_1 + C_2^{-1} \right)^{-1} \\ & \times X_1' C_1^{-1} X_1 X_2. \end{aligned} \quad (4)$$

In Sect. 3, we use this result to investigate how the entries of X_1 , C_1 and C_2 determine the rate of convergence for the CP and the NCP.

2.4 A simple example

To illustrate how parameterisation effects the posterior correlation of the model parameters and the rate of convergence of the Gibbs sampler, consider the following simple model, taken from Gelfand et al. (1996, Sect. 2). Let

$$Y_i = \theta + \beta_i + \epsilon_i, \quad (5)$$

with $\beta_i \sim N(0, \sigma_\beta^2)$ and $\epsilon_i \sim N(0, \sigma_\epsilon^2)$ independently distributed for all $i = 1, \dots, n$. The form of the model given by (5) is the NCP. The CP is found by replacing β_i with $\tilde{\beta}_i = \beta_i + \theta$, and so $Y_i = \tilde{\beta}_i + \epsilon_i$, and $\tilde{\beta}_i \sim N(\theta, \sigma_\beta^2)$.

Assuming a locally uniform prior distribution for θ , and that σ_β^2 and σ_ϵ^2 are known, Papaspiliopoulos et al. (2003) show that the exact convergence rates of the CP and the NCP of model (5) are

$$\lambda_c = \frac{\sigma_\epsilon^2}{\sigma_\epsilon^2 + \sigma_\beta^2},$$

and $\lambda_{nc} = 1 - \lambda_c$.

The rates of convergence highlight two important features of the sampling efficiency of the CP and the NCP. Firstly, that the ratio of the variance parameters is an important quantity in determining which parameterisation should be employed for model fitting, and secondly, that a change in variance ratio has opposing effects on each of the parameterisations.

3 CP versus NCP

In this section, we investigate how the rates of convergence are affected by the variance parameters and the correlation structure of the spatial processes. To focus on these relationships, we let $p = 1$ in model (1), giving us the following hierarchically centred model:

$$\begin{aligned} Y | \tilde{\beta}_0 & \sim N(\tilde{\beta}_0, \sigma_\epsilon^2 \mathbf{I}) \\ \tilde{\beta}_0 | \theta_0 & \sim N(\theta_0 \mathbf{1}, \sigma_0^2 \mathbf{R}_0) \\ \theta_0 & \sim N(m_0, v_0). \end{aligned} \quad (6)$$

It follows from Eqs. (3) and (4) that the respective rates of convergence for the CP and the NCP of model (6) are

$$\begin{aligned} \lambda_c = & \left(1/\sigma_0^2 \mathbf{1}^T \mathbf{R}_0^{-1} \mathbf{1} + 1/(\sigma_0^2 v_0) \right)^{-1} 1/\sigma_0^2 \mathbf{1}^T \mathbf{R}_0^{-1} \\ & \times \left(1/\sigma_\epsilon^2 \mathbf{I} + 1/\sigma_0^2 \mathbf{R}_0^{-1} \right)^{-1} 1/\sigma_0^2 \mathbf{R}_0^{-1} \mathbf{1}, \end{aligned} \quad (7)$$

and

$$\begin{aligned} \lambda_{nc} = & \left(n/\sigma_\epsilon^2 + 1/(\sigma_0^2 v_0) \right)^{-1} 1/\sigma_\epsilon^2 \mathbf{1}^T \\ & \times \left(1/\sigma_\epsilon^2 \mathbf{I} + 1/\sigma_0^2 \mathbf{R}_0^{-1} \right)^{-1} 1/\sigma_\epsilon^2 \mathbf{1}. \end{aligned} \quad (8)$$

For independent random effects, the ratio of variance parameters is important in determining the rates of convergence and so we introduce the quantity: $\delta_0 = \sigma_0^2/\sigma_\epsilon^2$. In Sects. 3.1–3.5, we use expressions (7) and (8) to compare the convergence rates for the CP and the NCP for different values of $\delta_0 = \sigma_0^2/\sigma_\epsilon^2$ and forms of \mathbf{R}_0 . In Sect. 3.6, we alter the model to include a covariate.

In Sects. 3.2–3.5, we confine ourselves to the case $1/v_0 = 0$, such that we have an improper prior distribution for θ_0 . This serves to clarify the effects of the other parameters on the convergence rates. To see the effect that the prior precision of θ_0 has on the convergence rate consider two different values for v_0 ; $v_{0,1}$ and $v_{0,2}$, with corresponding rates of convergence $\lambda_{c,1}$, $\lambda_{nc,1}$, $\lambda_{c,2}$ and $\lambda_{nc,2}$. Comparing the ratio of the convergence rates for the two different priors, we have

$$\frac{\lambda_{c,1}}{\lambda_{c,2}} = \frac{\mathbf{1}^T \mathbf{R}_0^{-1} \mathbf{1} + 1/v_{0,2}}{\mathbf{1}^T \mathbf{R}_0^{-1} \mathbf{1} + 1/v_{0,1}},$$

and clearly if $v_{0,1} < v_{0,2}$ then $\lambda_{c,1} < \lambda_{c,2}$. The same result can be seen for the NCP where

$$\frac{\lambda_{nc,1}}{\lambda_{nc,2}} = \frac{\sigma_0^2 n + \sigma_\epsilon^2/v_{0,2}}{\sigma_0^2 n + \sigma_\epsilon^2/v_{0,1}}.$$

Therefore, a more precise prior distribution for θ_0 will hasten convergence for both the CP and the NCP.

3.1 Convergence rates for equi-correlated random effects

To illustrate how changing the strength of correlation between the random effects influences the convergence rates of the different parameterisations, we begin by assuming a equi-correlation model. We suppose that

$$(\mathbf{R}_0)_{ij} = \begin{cases} \rho & \text{if } i \neq j \\ 1 & \text{if } i = j, \end{cases} \quad (9)$$

for $0 \leq \rho < 1$. We restrict ρ to take only non-negative values, as is usual in spatial data modelling. Roberts and Sahu (1997) consider a similar structure for the dispersion matrix of a Gaussian target distribution but do not include a global mean parameter as we do here.

To assist in the computation of convergence rates λ_c and λ_{nc} we make use of the following two matrix inversion identities. The first is the Sherman-Morrison-Woodbury formula, see for example Harville (1997, p. 423). Let N be an $n \times n$ matrix, U be an $n \times m$ matrix, M be an $m \times m$ matrix and V be an $m \times n$ matrix, then

$$(N + U M V)^{-1} = N^{-1} - N^{-1} U (M^{-1} + V N^{-1} U)^{-1} V N^{-1}. \quad (10)$$

The second result takes I to be the $n \times n$ identity matrix and J the $n \times n$ matrix of ones, then

$$(aI + bJ)^{-1} = \frac{1}{a}I - \frac{b}{a(a + nb)}J, \quad (11)$$

for constants $a > 0$, $b \neq -(a/n)$. This can be easily checked by direct multiplication and noting that

$$JJ = \mathbf{1}\mathbf{1}'\mathbf{1}' = \mathbf{1}n\mathbf{1}' = nJ.$$

Note also that identity (11) follows from (10) if we set $N = aI$, $U = \mathbf{1}$, $M = bI$ and $V = \mathbf{1}'$.

To compute the convergence rates given in (7) and (8) we must invert matrices $\sigma_0^2 R_0$ and $(1/\sigma_\epsilon^2 I + 1/\sigma_0^2 R_0^{-1})$. Using Eq. (10) we see that

$$\begin{aligned} (1/\sigma_\epsilon^2 I + 1/\sigma_0^2 R_0^{-1})^{-1} \\ = \sigma_\epsilon^2 I - \sigma_\epsilon^2 I (\sigma_\epsilon^2 I + \sigma_0^2 R_0)^{-1} \sigma_\epsilon^2 I. \end{aligned}$$

For R_0 defined by (9), we write

$$\sigma_0^2 R_0 = \sigma_0^2(1 - \rho)I + \sigma_0^2 \rho J, \quad (12)$$

and

$$\sigma_\epsilon^2 I + \sigma_0^2 R_0 = (\sigma_\epsilon^2 + \sigma_0^2(1 - \rho))I + \sigma_0^2 \rho J, \quad (13)$$

and we can invert the matrices given in (12) and (13) using Eq. (11).

After some cancellation we find the convergence rates for the CP and the NCP to be

$$\lambda_c = \frac{nv_0}{\sigma_0^2(1 - \rho) + n\sigma_0^2 \rho + nv_0} \left(\frac{\sigma_\epsilon^2}{\sigma_\epsilon^2 + \sigma_0^2(1 - \rho) + n\sigma_0^2 \rho} \right), \quad (14)$$

and

$$\lambda_{nc} = \frac{n\sigma_0^2 v_0}{\sigma_\epsilon^2 + n\sigma_0^2 v_0} \left(\frac{\sigma_0^2(1 - \rho) + n\sigma_0^2 \rho}{\sigma_\epsilon^2 + \sigma_0^2(1 - \rho) + n\sigma_0^2 \rho} \right). \quad (15)$$

If we assume an improper prior distribution, achieved by letting $1/v_0 = 0$, then $\lambda_{nc} = 1 - \lambda_c$, but otherwise equality does not hold. Note also that when $\rho = 0$, we recover the rates for the independent random effects model, see Sect. 2.4.

It is useful to re-write Eqs. (14) and (15) as

$$\lambda_c^{-1} = \left\{ 1 + v_0^{-1} \left[\sigma_0^2(1 - \rho)/n + \sigma_0^2 \rho \right] \right\} \{ 1 + \delta_0[1 + (n - 1)\rho] \}, \quad (16)$$

and

$$\lambda_{nc}^{-1} = [1 + (n\delta_0 v_0)^{-1}](1 + \{\delta_0[1 + (n - 1)\rho]\}^{-1}), \quad (17)$$

respectively. Consider first the case $1/v_0 = 0$ and recall that a lower rate indicates faster convergence. For the CP, increasing either δ_0 , ρ or n speeds up convergence. Increasing any one of these quantities has the opposing effect on the NCP. When $1/v_0 \neq 0$, λ_{nc} behaves as in the improper case. This is true of λ_c with respect to δ_0 and ρ , but it is no longer monotonic in n .

3.2 Effect of spatial correlation

In spatial modelling, the correlation between two realisations of a latent process is usually assumed to be a function of their separation. Here, we consider exponential correlation functions, which are used widely in applications (Sahu et al. 2010; Berrocal et al. 2010; Sahu et al. 2007; Huerta et al. 2004). We have that

$$(R_0)_{ij} = \exp(-\phi_0 d_{ij}),$$

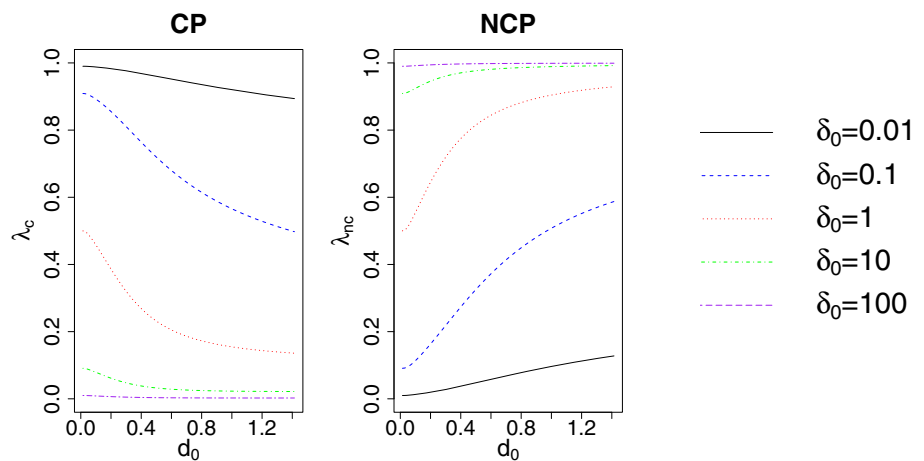
where ϕ_0 is the spatial decay parameter. The exponential correlation function belongs to the Matérn family and is found by letting $\nu = 0.5$ in Eq. (2). To see this we can use the following results taken from Schabenberger and Gotway (2004, Sect. 4.3.2)

$$\Gamma(0.5) = \sqrt{\pi}, \quad K_{0.5}(t) = \sqrt{\frac{\pi}{2t}} e^{-t}.$$

We characterise the strength of correlation in terms of the effective range, which we define as the distance, d_0 , such that $\text{corr}\{\beta_0(s_i), \beta_0(s_j)\} = 0.05$. For an exponential correlation function, we have that

$$d_0 = -\log(0.05)/\phi_0 \approx 3/\phi_0.$$

Fig. 1 Convergence rate against effective range for the CP and the NCP at different levels of δ_0



We cannot compute explicit expressions for the entries of \mathbf{R}_0^{-1} , and hence, we cannot find expressions for the convergence rate in terms of ϕ_0 . Therefore, we use a simulation approach. We take the unit square to be the spatial domain, randomly selecting $n = 40$ locations which will be used throughout the rest of this section.

We consider five values of the variance ratio δ_0 and vary the strength of spatial correlation by controlling the effective range $d_0 = -\log(0.05)/\phi_0$. For each value of δ_0 , we compute the convergence rates given in Eqs. (7) and (8) for effective ranges between zero, implying no spatial correlation, and $\sqrt{2}$, the maximum possible separation of two points in the domain.

Convergence rates are plotted against the effective range for the CP and the NCP in Fig. 1, where again a lower rate indicates faster convergence. For a fixed d_0 , we can see that increasing the δ_0 decreases the convergence rate for the CP, but increases it for the NCP, as we might expect given the results for independent random effects. We also observe that for a fixed level of δ_0 increasing d_0 , thus increasing the strength of correlation between the random effects, decreases the convergence rate for the CP and increases it for the NCP. Hence, the complimentary nature of the CP and NCP is maintained when we vary the strength of exponential correlation across the random effects. However, the two rates do not add to 1 unlike in the simpler case of just below-mentioned Eq. (15).

The convergence rates are dependent on the set of sampling locations. For a different set of locations, the convergence rates are changed, but for our given set of locations, the overall picture is not; increasing δ_0 or d_0 quickens convergence for the CP and slows convergence for the NCP.

3.3 Effect of the smoothness parameter in the Matérn correlation function

In this section, we consider different correlation functions from the Matérn family, see Eq. (2) for the general form. ν

is a half integer, such that $\nu = b + 0.5$ for $b = 0, 1, 2, \dots$, the correlation function takes on a simpler form. Taken from [Rasmussen and Williams \(2006, Sect. 4.2\)](#), we have that

$$\rho(d_{ij}, \phi, \nu) = \exp\left(-\sqrt{2\nu}\phi d_{ij}\right) \frac{\Gamma(b+1)}{\Gamma(2b+1)} \sum_{r=0}^b \frac{(b+r)!}{r!(b-r)!} \left(\sqrt{8\nu}\phi d_{ij}\right)^{b-r}.$$

In particular, when $\nu = 1.5$, the correlation function is

$$\rho(d_{ij}, \phi) = \left(1 + \sqrt{3}\phi d_{ij}\right) \exp\left(-\sqrt{3}\phi d_{ij}\right), \quad (18)$$

and $\nu = 2.5$, it becomes

$$\rho(d_{ij}, \phi) = \left(1 + \sqrt{5}\phi d_{ij} + \frac{5\phi^2 d_{ij}^2}{3}\right) \exp(-\sqrt{5}\phi d_{ij}). \quad (19)$$

As $\nu \rightarrow \infty$, the correlation function goes to

$$\rho(d_{ij}, \phi) = \exp\left(-\frac{\phi^2 d_{ij}^2}{2}\right),$$

which is sometimes known as the squared exponential or Gaussian correlation function.

We consider model (6) and compare the convergence rates for the CP and the NCP for the exponential, $\nu = 1.5$, $\nu = 2.5$ and Gaussian correlation functions. In Sect. 3.2, the strength of correlation is considered in terms of the effective range, which for the exponential correlation function is $-\log(0.05)/\phi_0$. In terms of ϕ_0 , the effective range for the Gaussian correlation function is given by $\sqrt{-2\log(0.05)}/\phi_0$. For other members of the Matérn class, there is no closed form expression for the effective range. Therefore, for the cases ν is equal to 1.5 and 2.5, we take an effective range d_0

Fig. 2 Convergence rates for the CP of model (6) for different values of ν . **a** $\delta_0 = 0.01$, **b** $\delta_0 = 0.1$, **c** $\delta_0 = 1$, **d** $\delta_0 = 10$, **e** $\delta_0 = 100$

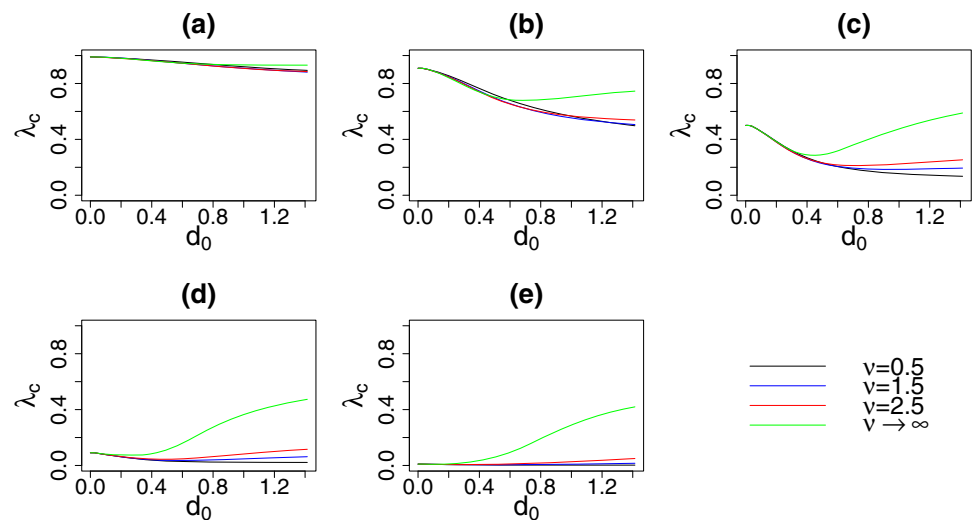
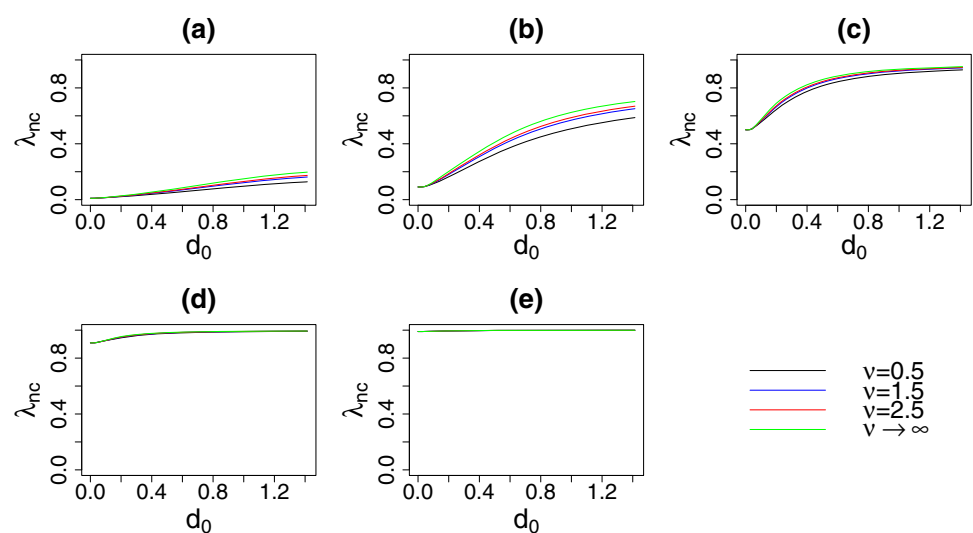


Fig. 3 Convergence rates for the NCP of model (6) for different values of ν . **a** $\delta_0 = 0.01$, **b** $\delta_0 = 0.1$, **c** $\delta_0 = 1$, **d** $\delta_0 = 10$, **e** $\delta_0 = 100$



and search for the value of ϕ_0 that solves

$$\rho(d_0, \phi_0) - 0.05 = 0, \quad (20)$$

where $\rho(d_0, \phi_0)$ is given by functions (18) and (19) respectively.

Convergence rates are computed for each parameterisation for effective ranges between 0 and $\sqrt{2}$ and for five values of $\delta_0 = 0.01, 0.1, 1, 10, 100$. Recall that as previously in Sect. 3.2 we take the unit square to be the spatial domain and randomly select 40 locations. The results for the CP are given in Fig. 2. We see that for fixed ν and ϕ_0 , increasing the δ_0 reduces the convergence rate. Also we see that for fixed ϕ_0 and δ_0 , the convergence rate is increased when ν is increased, except for the $\delta_0 = 0.1$ case where the ordering only becomes apparent as the effective range is increased. Unlike in the case of $\nu = 0.5$, increasing the effective range does not reduce the convergence rate for other values of ν .

The equivalent plot for the NCP is given in Fig. 3. For fixed ν and ϕ_0 , the increasing δ_0 increases the convergence

rate. For fixed ϕ_0 and δ_0 , the increasing ν slows convergence as it does for the CP. The convergence rate is monotonically increasing with the increasing effective range for all four correlation functions. We also note that convergence rates for the NCP are not as sensitive to changes in ν as they are for the CP.

Figures 2 and 3 show that increasing ν , the smoothness parameter, leads to slower convergence for both the CP and NCP, which contradicts the complimentary behaviour of the rates of convergence seen in the previous two subsections. In order to understand this further, we investigate as follows. The rates of convergence, as obtained in (7) and (8) depend on R_0^{-1} which is guaranteed to be positive definite for any n when a member of the Matérn class of correlation functions is adopted. However, in practical numerical calculations with a large value of n (much greater than 40 used in Figs. 2 and 3) R_0 becomes increasingly singular in the presence of high level of spatial correlation and smoothness. The minimum value of n for which this near singularity is observed depends on the minimum distance between the " C_2 pairs of the n

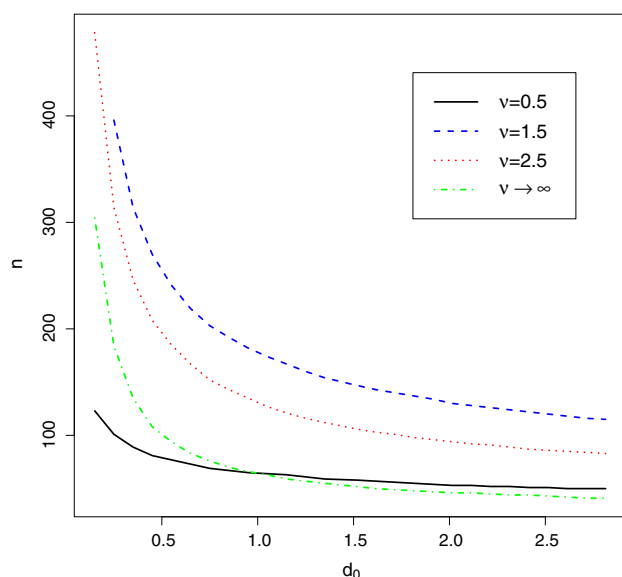


Fig. 4 The minimum value of n for which R_0 is nearly singular for different values of d_0 and ν

locations and the values of the parameters ϕ and ν in the Matérn correlation function. To investigate this minimum value of n , we randomly sample 500 points in the unit square where the minimum Euclidean distance between any two points is greater than a threshold value, which is taken to be 0.01. This ensures that the singularity is not caused by very closely located points in the sample.

Next we find the minimum value of n for which R_0 is approximately singular, viz., the value of the determinant less than 10^{-7} , for a given value of ν and a given value of d_0 . The value of ϕ_0 is chosen by (20) $\nu = 1.5$ and 2.5 and $\phi_0 = -\log(0.05)/d_0$ when $\nu = 0.5$, i.e. the case of exponential correlation function and $\phi_0 = \sqrt{-2 \log(0.05)}/d_0$ in the case of the Gaussian correlation function, when $\nu \rightarrow \infty$. The minimum value is plotted in Fig. 4 against d_0 for all four correlation functions considered here. As expected, for smaller values of d_0 , tending to the case of zero spatial correlation, the minimum value of n becomes very large. Increasing spatial correlation, as effected by increasing d_0 , leads to a smaller value of n at which near singularity of R_0 is reached. Interestingly, increasing smoothness, through values of ν hasten this except for the case of the exponential correlation function. This is due to the nonlinear effect of the ν and ϕ on the Matérn correlation function.

The near singularity in the limiting Gaussian case is related to the near predictability of the associated spatial processes. In the limiting Gaussian case, it is possible to predict $Y(s')$ for any s' in the same spatial domain upon observing the same spatial process $Y(s)$ at any location s , see, e.g. page 62 of Banerjee et al. (2015). Such deterministic behaviour leads to the near singularity of covariance matrices which in turn leads

to non-convergence. Indeed, Stein (1999) (page 70) explicitly recommends not to use the Gaussian correlation function to model physical processes. The investigation here confirms this view by pointing to non-convergence of the MCMC fitting algorithms for smooth and highly correlated processes for large values of n . This also allows us to conclude that increasing n in an infill asymptotic sense (Zhang 2004), in the presence of high spatial correlation, will lead to near singularity of the covariance matrices, which in turn will lead to non-convergence of either of the two parameterisations. This asymptotic result, however, does not spell disaster for the CP and NCP in practical problems where the strength of the spatial correlation is such that the resulting effective range (d_0 in the exponential case) is significantly less than the maximum distance ($\sqrt{2}$ in the unit square) between any two points in a closed spatial domain. Such a situation will allow modelling of a reasonably large number of spatial observations (e.g. low 100s) using the centring parameterisations.

3.4 Effect of introducing geometric anisotropy

The class of Matérn correlation functions is isotropic. This means that the correlation between the random variables at any two points, s_i and s_j , depends on the distance between them $d_{ij} = \|s_i - s_j\|$ (and parameters ϕ and ν), and hence, the contours of iso-correlation are circular. The assumption that spatial dependence is the same in all directions is not always appropriate and therefore we may seek an anisotropic specification for the correlation structure.

Anisotropic correlation functions are widely used and have been employed to model, for example, scallop abundance in the North Atlantic (Ecker and Gelfand 1999), extreme precipitation in Western Australia (Apputhurai and Stephenson 2013) and the phenotypic traits of trees in northern Sweden (Banerjee et al. 2010).

Different forms of anisotropy exist, see Zimmerman (1993), but we consider only geometric anisotropy. Geometric anisotropic correlation functions can be constructed from isotropic correlation functions by taking a linear transformation of the lag vector $s_i - s_j$. Let

$$d_{ij}^* = \|\mathbf{G}(s_i - s_j)\|, \quad (21)$$

where \mathbf{G} is a 2×2 transformation matrix. In Euclidean space, (21) is equivalent to

$$d_{ij}^* = \{(s_i - s_j)^T \mathbf{H}(s_i - s_j)\}^{1/2},$$

where $\mathbf{H} = \mathbf{G}^T \mathbf{G}$. The matrix \mathbf{H} must be positive definite, i.e. $d_{ij}^* > 0$ for $s_i \neq s_j$, which is ensured if \mathbf{G} is non-singular, see Harville (1997, Corollary 14.2.14). By replacing d_{ij} with d_{ij}^* in (2), we have a geometric anisotropic Matérn correlation function with elliptical contours of iso-correlation.

$$\rho(d_{ij}^*, \phi) = \exp(-\phi d_{ij}^*). \quad (23)$$

We take the point $\mathbf{s}^* = (0.5, 0.5)^T$ in the unit square and fix decay parameter $\phi = 1$. We then compute the correlation between \mathbf{s}^* and all points on a 20×20 grid, according to the correlation function given in (23). The values are then smoothed to produce a correlation surface. This is repeated for each of the nine combinations of α and ψ and displayed in Fig. 5.

We can see that setting $\alpha = 0.5$ strengthens correlation in the x-direction. This is because for the purposes of computing correlation, the separation of two points in the x-direction is halved. When $\alpha = 1$, the angle of rotation ψ does not effect the contours as they are circular. Clearly, setting $\alpha = 2$ has the effect of weakening correlation in the x-direction.

To assess the impact of anisotropy on the convergence rates for the CP and the NCP we return to model (6). We consider an anisotropic exponential correlation function for the spatial process and so

$$(\mathbf{R}_0)_{ij} = \exp \left(-\phi_0 d_{ij}^* \right),$$

where d_{ij}^* is given by Eq. (21). We begin by fixing $\psi = 0$ and letting $\alpha = 0.5, 1, 2$. This corresponds to pan-

 Springer

Fig. 6 Convergence rates for the CP of model (6) with an anisotropic exponential correlation function for different values of α . **a** $\delta_0 = 0.01$, **b** $\delta_0 = 0.1$, **c** $\delta_0 = 1$, **d** $\delta_0 = 10$, **e** $\delta_0 = 100$

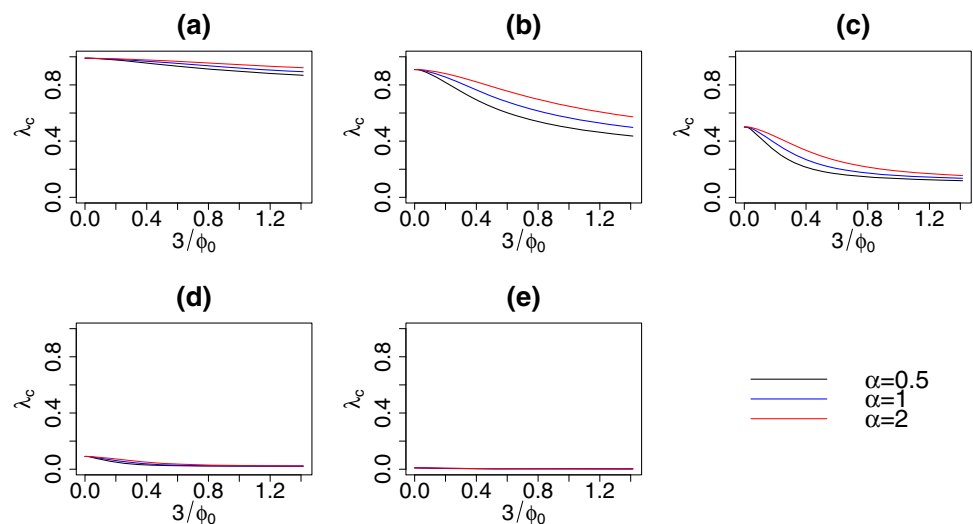
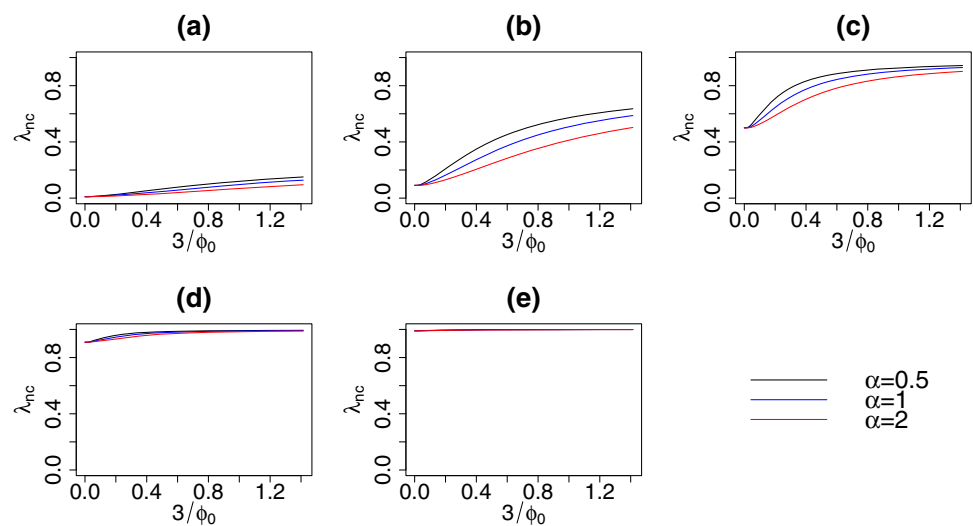


Fig. 7 Convergence rates for the NCP of model (6) with an anisotropic exponential correlation function for different values of α . **a** $\delta_0 = 0.01$, **b** $\delta_0 = 0.1$, **c** $\delta_0 = 1$, **d** $\delta_0 = 10$, **e** $\delta_0 = 100$



els 1(a), 2(a), and 3(a), in Fig. 5. We use five values for $\delta_0 = \sigma_0^2/\sigma_\epsilon^2 = 0.01, 0.1, 1, 10, 100$ and vary ϕ_0 such that $3/\phi_0 \in (0, \sqrt{2}]$. Here, the effective range is direction dependent so we no longer refer to $3/\phi_0$ as the effective range.

We compute convergence rates for the CP and the NCP and plot results in Figs. 6 and 7 respectively. As α is reduced, we increase the strength of correlation in the x-direction. This result is faster convergence for the CP and slower convergence for the NCP. This is consistent with the results of Sect. 3.2 which shows that increasing the effective range of an isotropic exponential correlation function, thus strengthening the correlation in all directions, helps the CP and hinders the NCP.

We now look at the effect of rotating the axis. If $\alpha = 1$ then a rotation has no impact on the correlation function as \mathbf{H} is the identity. We consider four combinations of $\alpha = 0.5, 2$ and $\psi = \pi/4, \pi/2$. These values correspond to panels 1(b)

and 1(c) for $\alpha = 0.5$, and 3(b) and 3(c) for $\alpha = 2$ in Fig. 5. Again, we let $\delta_0 = 0.01, 0.1, 1, 10, 100$ and vary ϕ_0 such that $3/\phi_0 \in (0, \sqrt{2}]$.

The results for the CP and the NCP are given in Figs. 8 and 9 respectively. We can see that changing ψ has very little effect on the convergence rates of either parameterisation as expected since d_{ij}^* is free of ψ . However, because we apply the rotation matrix first the subsequent stretch effectively acts in a different direction, that direction depending on ψ , and the resulting values for d_{ij}^* may be different. Take the example we used here. Let $\mathbf{s} = (s_1, s_2)^T$ and $\mathbf{s}_i - \mathbf{s}_j = (s_{i1} - s_{j1}, s_{i2} - s_{j2})^T = (l_1, l_2)^T$. For $\psi = \pi/4$, $d_{ij}^* = \sqrt{0.5[(l_1 - l_2)^2 + \alpha^2(l_1 + l_2)^2]}$, whereas if $\psi = \pi/2$ then $d_{ij}^* = \sqrt{l_1^2 + \alpha^2 l_2^2}$. The different values of d_{ij}^* may lead to different rates of convergence. Further investigation is needed to determine whether similar results hold for patterned sampling locations.

Fig. 8 Convergence rates for the CP of model (6) with an anisotropic exponential correlation function for different values of α and ψ . **a** $\delta_0 = 0.01$, **b** $\delta_0 = 0.1$, **c** $\delta_0 = 1$, **d** $\delta_0 = 10$, **e** $\delta_0 = 100$

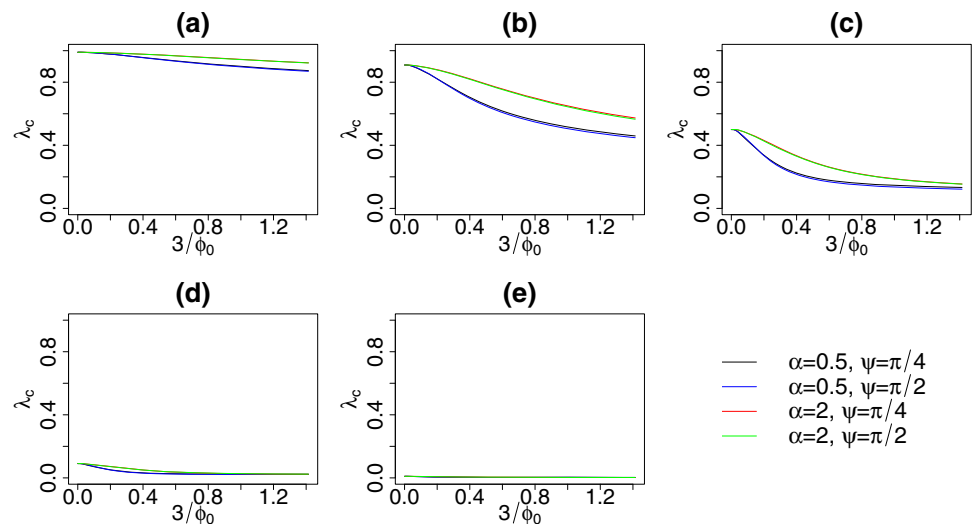
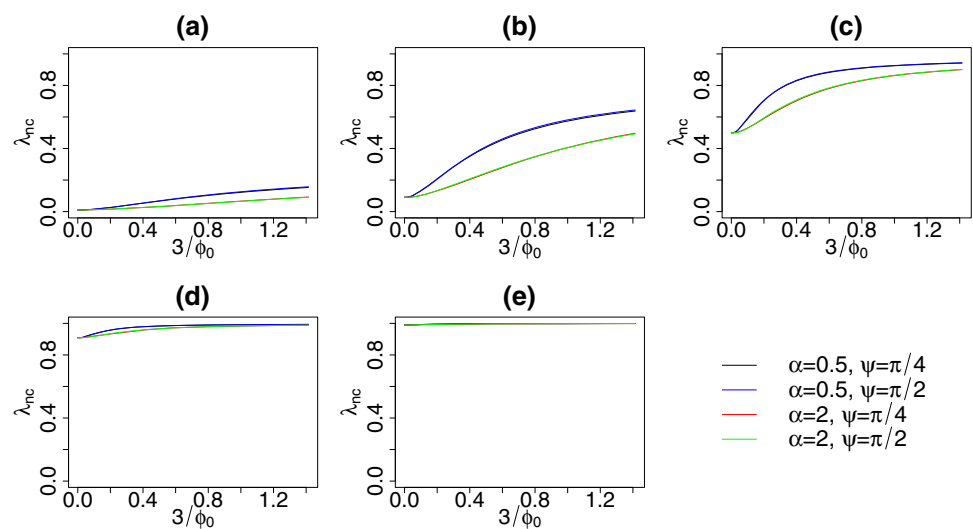


Fig. 9 Convergence rates for the NCP of model (6) with an anisotropic exponential correlation function for different values of α and ψ . **a** $\delta_0 = 0.01$, **b** $\delta_0 = 0.1$, **c** $\delta_0 = 1$, **d** $\delta_0 = 10$, **e** $\delta_0 = 100$



3.5 Effect of introducing tapered covariance matrices

When spatial association is modelled as a Gaussian process, the resulting covariances matrices are dense and inverting them can be slow or even infeasible for large n . One strategy to deal with this is covariance tapering (Furrer et al. 2006; Kaufman et al. 2008). The idea is to set to zero the entries in the covariance matrix that correspond to pairs of locations that are separated by a distance greater than a predetermined range. This results in sparse matrices that can be inverted more quickly than the original. In this section, we investigate the effect covariance tapering on the convergence rates for the CP and the NCP. We take model (6) with an exponential correlation function for \mathbf{R}_0 and compare the convergence rates found in Sect. 3.2 with those computed when we use a tapered covariance matrix.

The tapered correlation matrix, \mathbf{R}_{Tap} , is the element-wise product of the original correlation matrix \mathbf{R}_0 and the tapering correlation matrix \mathbf{T} , where \mathbf{T} is a sparse matrix with ij th entry equal to zero if d_{ij} is greater than some threshold dis-

tance. Positive definiteness of \mathbf{R}_{Tap} is assured if \mathbf{T} is positive definite (Horn and Johnson 2012, Theorem 7.5.3).

Given that our original correlation function is an exponential one, we follow Furrer et al. (2006) and use a spherical tapering function such that

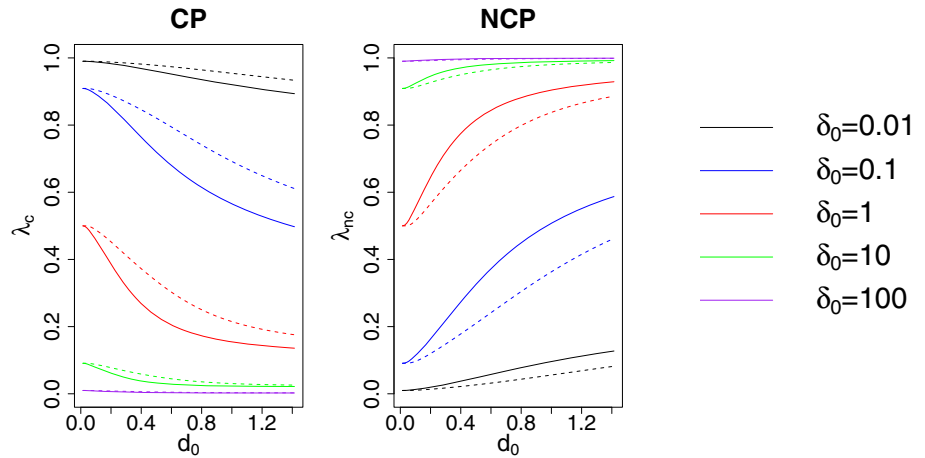
$$T_{ij} = \begin{cases} 1 - \frac{3d_{ij}\chi}{2} + \frac{d_{ij}^3\chi^3}{2} & \text{if } d_{ij} < 1/\chi, \chi > 0 \\ 0 & \text{otherwise,} \end{cases}$$

with decay parameter χ , where $1/\chi$ is equal to the effective range, so that here we have $\chi = -\phi_0/\log(0.05)$. Therefore,

$$(\mathbf{R}_{Tap})_{ij} = \begin{cases} \exp(-\phi_0 d_{ij}) \left(1 - \frac{3d_{ij}\chi}{2} + \frac{d_{ij}^3\chi^3}{2} \right) & \text{if } d_{ij} < d_0, \phi_0 > 0, \chi > 0 \\ 0 & \text{otherwise,} \end{cases}$$

where $d_0 = -\log(0.05)/\phi_0$ is the effective range.

Fig. 10 Convergence rates with tapered covariance matrices for the CP and the NCP at different levels of δ_0



We let $\delta_0 = 0.01, 0.1, 1, 10$ and 100 and vary d_0 between 0 and $\sqrt{2}$. The convergence rates for the CP and the NCP are given in Fig. 10. The dashed line represents the use of the tapered correlation matrix. The solid line for comparison are the rates achieved using the original correlation matrix \mathbf{R}_0 and are identical to those given in Fig. 1. Convergence rates are slowed by tapering for the CP and hastened for the NCP. Intuitively we can say that the under the CP stronger correlation is desirable and tapering reduces that, with the opposite being true for the NCP.

We can illustrate this effect by considering a spatial model with just two locations s_1 and s_2 such that $s_1 \neq s_2$. Let $0 \leq \text{corr}(\beta(s_1), \beta(s_2)) = \rho < 1$. Suppose that we use a tapering function that takes values ρ^* if $d_{12} < d_0$ and zero otherwise, where $0 \leq \rho^* < 1$. The tapered correlation is

$$\rho_{Tap} = \begin{cases} \rho\rho^* & \text{if } d_{12} < d_0 \\ 0 & \text{otherwise.} \end{cases}$$

Therefore, $\rho_{Tap} \leq \rho$, with equality attained only when $\rho = 0$. We know from Eqs. (16) and (17) that for equi-correlated random effects if ρ decreases, λ_c is increased and λ_{nc} is decreased. In other words, when $n = 2$, tapering can only worsen the performance of CP and improve the performance of the NCP.

3.6 Effect of covariates

In this section, we investigate the effect of the covariates upon the rate of convergence. We consider the following model

$$Y(s_i) = \{\theta_1 + \beta_1(s_i)\}x_1(s_i) + \epsilon(s_i) \quad (i = 1, \dots, n), \quad (24)$$

which may be found by letting $k = 1, \dots, p-1$, and $p = 2$ in model (1). Recalling that $\tilde{\beta}_1 = \{\tilde{\beta}_1(s_1), \dots, \tilde{\beta}_1(s_n)\}^T$, where $\tilde{\beta}_1(s_i) = \beta_1(s_i) + \theta_1$, and $\mathbf{x}_1 = \{x_1(s_1), \dots, x_1(s_n)\}^T$ and $\mathbf{D}_1 = \text{diag}(\mathbf{x}_1)$, we can write model (24) in the following form

$$\begin{aligned} Y | \tilde{\beta}_1 &\sim N(\mathbf{D}_1 \tilde{\beta}_1, \sigma_\epsilon^2 \mathbf{I}) \\ \tilde{\beta}_1 | \theta_1 &\sim N(\theta_1 \mathbf{1}, \sigma_1^2 \mathbf{R}_1) \\ \theta_1 &\sim N(m_1, \sigma_1^2 v_1). \end{aligned} \quad (25)$$

We consider only one covariate and so in the rest of this section we drop the subscript from \mathbf{D}_1 and \mathbf{x}_1 .

First suppose that random effects are independent. This can be considered the limiting case for weakening spatial correlation. For the sake of notational clarity, under the assumption of spatial independence we write $x(s_i) = x_i$ ($i = 1, \dots, n$). The convergence rate for the CP is

$$\lambda_c = \frac{1}{n + 1/v_1} \sum_{i=1}^n \frac{\sigma_\epsilon^2}{\sigma_\epsilon^2 + \sigma_1^2 x_i^2}.$$

If we let $1/v_1 = 0$, thus implying and improper prior for θ_1 , we can write λ_c as

$$\lambda_c = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + (\sigma_1^2/\sigma_\epsilon^2) x_i^2}. \quad (26)$$

We introduce the variable $\delta_1 = \sigma_1^2/\sigma_\epsilon^2$. For fixed \mathbf{x} , we can see that as δ_1 tends to zero, the convergence rate for the CP of model (24) tends to one. As δ_1 gets larger, the convergence rate goes to zero. To see the effect of the scale of \mathbf{x} , we introduce variables u_i , where

$$u_i = \frac{x_i - \bar{x}}{sd_x}, \quad i = 1, \dots, n, \quad (27)$$

and \bar{x} and sd_x are the sample mean and sample standard deviation of \mathbf{x} respectively. Substituting Eq. (27) into Eq. (26) we have

$$\lambda_c = \frac{1}{n} \sum_{i=1}^n \frac{1}{1 + (\sigma_1^2/\sigma_\epsilon^2) (u_i sd_x + \bar{x})^2}.$$

We suppose that the x_i 's have already been centred on zero and so $\bar{\mathbf{x}} = 0$. For fixed variance parameters, the effect of the scale of \mathbf{x} is clear; an increase in sd_x results in a decrease in the convergence rate and vice versa.

For the NCP, the convergence rate is

$$\lambda_{nc} = \frac{1}{\sum_{i=1}^n x_i^2 + \sigma_\epsilon^2 / (\sigma_1^2 v_1)} \sum_{i=1}^n \frac{\sigma_1^2 x_i^4}{\sigma_\epsilon^2 + \sigma_1^2 x_i^2}.$$

Letting $1/v_1 = 0$, we can write λ_{nc} as

$$\lambda_{nc} = \frac{1}{\sum_{i=1}^n x_i^2} \sum_{i=1}^n \frac{x_i^4}{(\sigma_\epsilon^2 / \sigma_1^2) + x_i^2}. \quad (28)$$

For fixed \mathbf{x} , if $\sigma_\epsilon^2 / \sigma_1^2$ goes to zero then λ_{nc} goes to one. In contrast, as the data variance dominates that of the random effects, the convergence rate falls. Note that in general $\lambda_c + \lambda_{nc} \neq 1$.

To see the effect of the scale of \mathbf{x} upon λ_{nc} we substitute Eq. (27) into equation (28). Then we have

$$\lambda_{nc} = \frac{1}{\sum_{i=1}^n (u_i sd_x + \bar{\mathbf{x}})^2} \sum_{i=1}^n \frac{(u_i sd_x + \bar{\mathbf{x}})^4}{(\sigma_\epsilon^2 / \sigma_1^2) + (u_i sd_x + \bar{\mathbf{x}})^2}.$$

Again, assuming $\bar{\mathbf{x}} = 0$, we get

$$\begin{aligned} \lambda_{nc} &= \frac{1}{\sum_{i=1}^n (u_i sd_x)^2} \sum_{i=1}^n \frac{(u_i sd_x)^4}{(\sigma_\epsilon^2 / \sigma_1^2) + (u_i sd_x)^2} \\ &= \frac{1}{\sum_{i=1}^n u_i^2} \sum_{i=1}^n \frac{u_i^4}{(\sigma_\epsilon^2 / \sigma_1^2 sd_x^2) + u_i^2}. \end{aligned}$$

Fixing σ_ϵ^2 and σ_1^2 , as sd_x tends to infinity, λ_{nc} tends to 1, as sd_x tends to zero, λ_{nc} tends to 0.

Now we investigate the effect that increasing the strength of correlation between realisations of the slope surface has upon the performance of the CP and the NCP. We let $(\mathbf{R}_1)_{ij} = \exp(-\phi_1 d_{ij})$ and so the effective range $d_1 = -\log(0.05)/\phi_1$. To generate the values of \mathbf{x} we select a point s_x , which we may imagine to be the site of a source of pollution. We assume that the value for the observed covariate at site s decays exponentially at a rate ϕ_x with increasing separation from s_x , so that

$$x(s_i) = \exp(-\phi_x \|s_i - s_x\|) \quad (i = 1, \dots, n).$$

The spatial decay parameter ϕ_x is chosen such that there is an effective spatial range of $\sqrt{2}/2$, i.e. if $\|s - s_x\| = \sqrt{2}/2$ then $x(s) = 0.05$. The values of \mathbf{x} are standardised by subtracting their sample mean and dividing by their sample standard deviation.

We compute the convergence rate for the CP and the NCP for model (25) for five values of $\delta_1 = 0.01, 0.1, 1, 10, 100$, and for an effective range d_1 between 0 and $\sqrt{2}$. Results are given in Fig. 11. We see that for the CP for a fixed d_1 , increasing δ_1 achieves faster convergence. If we fix δ_1 the performance of the CP is improved as the effective range is increased. The opposite is seen for the NCP, whose performance is improved by decreasing δ_1 or shortening the effective range. Therefore, the variance ratio δ_1 and the decay parameter ϕ_1 have same influence on the convergence rates of the CP and NCP as δ_0 and ϕ_0 .

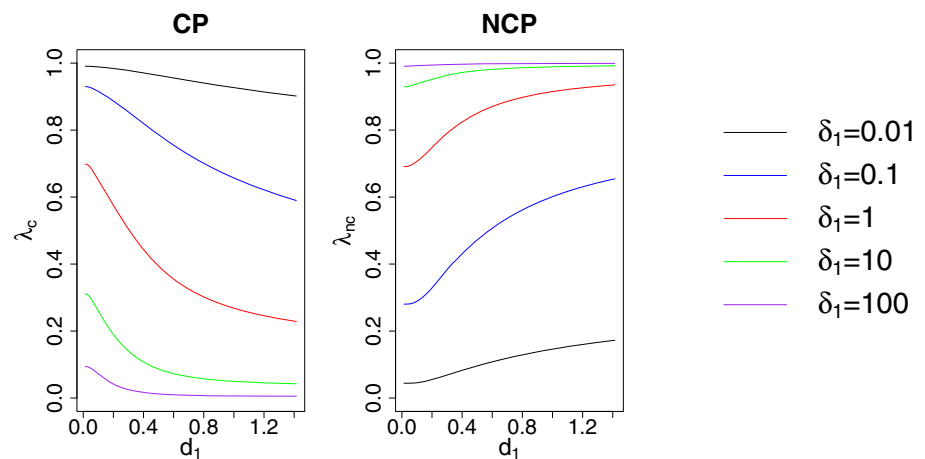
4 Practical examples with unknown covariance parameters

In this section, we focus on the practical implementation of the Gibbs sampler for the CP and the NCP for spatially varying coefficient models. The joint posterior distribution is unaffected by hierarchical centring and so inferential statements are the same under either parameterisation. However, what is affected is the efficiency of the Gibbs sampler used to make those statements.

In Sect. 3, the CP and the NCP are compared in terms of the exact convergence rates of the associated Gibbs samplers. The key assumption needed to compute these rates is that the joint posterior distribution is Gaussian with known precision matrix. Here we allow for the more common scenario that the precision matrix is known only up to a set of covariance parameters. In this case, we cannot compute the exact convergence rate. Therefore, we use the MCMC samples to assess the efficiency of the Gibbs samplers induced by the CP and the NCP. The full conditional distributions needed to construct the Gibbs samplers are given in Appendix 2.

We employ two diagnostic statistics to compare parameterisations. The first statistic we use is based on the multi-variate potential scale reduction factor (MPSRF) (Brooks and Gelman 1998). We define the $\text{MPSRF}_M(1.1)$ to be the number of iterations required for the MPSRF to fall below 1.1. To compute the $\text{MPSRF}_M(1.1)$ we run five chains of length 25,000 from widely dispersed starting values. In particular, we take values that are outside of the intervals described by pilot chains. Moreover, the same starting values are used for both the CP and the NCP. At every fifth iteration, the MPSRF is calculated and number of iterations for its value to first drop below 1.1 is the value that we record. The second statistic we use is the effective sample size (ESS) of the model parameters (Robert and Casella 2004). The ESS is computed using all 125,000 MCMC samples and gives us a measure of the Markovian dependence between successive MCMC iterates, with values of 125,000 indicating independence. There is a negligible difference in the run times for the CP and the NCP and so we do not adjust these measures by computation time.

Fig. 11 A comparison of convergence rates for the CP and the NCP at different levels of δ_1



4.1 A simulation study

We simulate data from model (6) for $n = 40$ randomly chosen locations across the unit square assuming an exponential correlation function for the spatial process. We set $\theta_0 = 0$ and generate data with five variance parameter ratios such that $\delta_0 = \sigma_0^2/\sigma_\epsilon^2 = 0.01, 0.1, 1, 10, 100$. This is done by letting $\sigma_0^2 = 1$ and varying σ_ϵ^2 accordingly. For each of the five levels of δ_0 , we have four values of the decay parameter ϕ_0 , chosen such that there is an effective range of $0, \sqrt{2}/3, 2\sqrt{2}/3$ and $\sqrt{2}$, where $\sqrt{2}$ is the maximum possible separation of two points in the unit square. Hence, there are 20 combinations of $\sigma_0^2, \sigma_\epsilon^2$ and ϕ_0 in all. Each of these combinations is used to simulate 20 datasets, and so there are 400 datasets in total.

We fix the decay parameter at its known value and sample from the marginal posterior distributions of θ_0, σ_0^2 and σ_ϵ^2 . We let hyperparameters $m_0 = 0$ and $v_0 = 10^4$. Recall that the variance parameters are given inverse gamma prior distributions with $\pi(\sigma_0^2) = IG(a_0, b_0)$ and $\pi(\sigma_\epsilon^2) = IG(a_\epsilon, b_\epsilon)$. We let $a_0 = a_\epsilon = 2$ and $b_\epsilon = b_0 = 1$, implying a prior mean of one and infinite prior variance for σ_0^2 and σ_ϵ^2 . These are common hyperparameters for inverse gamma prior distributions, see Sahu et al. (2007, 2010), Gelfand et al. (2003).

Figure 12 shows boxplots of the MPSRF_M(1.1) (top row) and the ESS of θ_0 (bottom row) for the CP. Each panel contains the results for a fixed value of δ_0 , increasing from 0.01 on the left to 100 on the right. Each panel contains four boxplots corresponding to the four effective ranges of $0, x/3, 2x/3$, and x , where $x = \sqrt{2}$. As the effective range increases, we have stronger spatial correlation between the random effects. Each boxplot is produced from the 20 values obtained for a given combination of δ_0 and ϕ_0 . We can see that the performance of the CP improves with increasing δ_0 and also with increasing strength of correlation between the random effects.

The equivalent plot for the NCP is given in Fig. 13. We can see a reverse of the pattern displayed by the CP. The performance of the NCP is worsened as δ_0 increases and the detrimental effect of increasing the strength of correlation between the random effects is also clearly evident. Therefore, δ_0 and the d_0 have the same influence on the CP and the NCP as we saw for the exact convergence rates when the variance parameters were assumed to be known in Sect. 3.2.

Figure 14 gives the ESS of σ_0^2 (top row) and the ESS of σ_ϵ^2 (bottom row) for the CP. We can see a general increasing trend in the ESS of σ_0^2 for increasing δ_0 , but a downward trend is seen for σ_ϵ^2 . However, for a fixed value of δ_0 we can see an improvement as the effective range increases, particularly in σ_ϵ^2 . This is because for the case when there is zero effective range, marginally the data variance is $(\sigma_0^2 + \sigma_\epsilon^2)\mathbf{I}$, and so increasing the effective range moves us away from the unidentifiable case which can result in poor mixing of the chains.

Figure 15 shows the ESS of σ_0^2 and σ_ϵ^2 for the NCP. We see that the ESS of σ_0^2 is stable under changes in δ_0 and d_0 , with the exception being the case where $\delta_0 = 0.1$ and $d_0 = 0$. In this case, the results are again explained by the lack of identifiability of the variance parameters for independent random effects. The ESS of σ_ϵ^2 is reduced by increasing δ_0 . For a fixed value of δ_0 , we can see an improvement in the ESS as d_0 increases. This was also observed for σ_ϵ^2 under the CP and is similarly explained.

4.2 Real-data example

In this section, we compare the sampling efficiency of the CP and the NCP when they are fitted to a real dataset. We have annual PM10 concentrations, in micrograms per cubic metre ($\mu\text{g}/\text{m}^3$), taken from 70 monitoring sites in Greater London, UK. We use data from 50 sites for model fitting, leaving out data from 20 sites for model validation, see Fig.

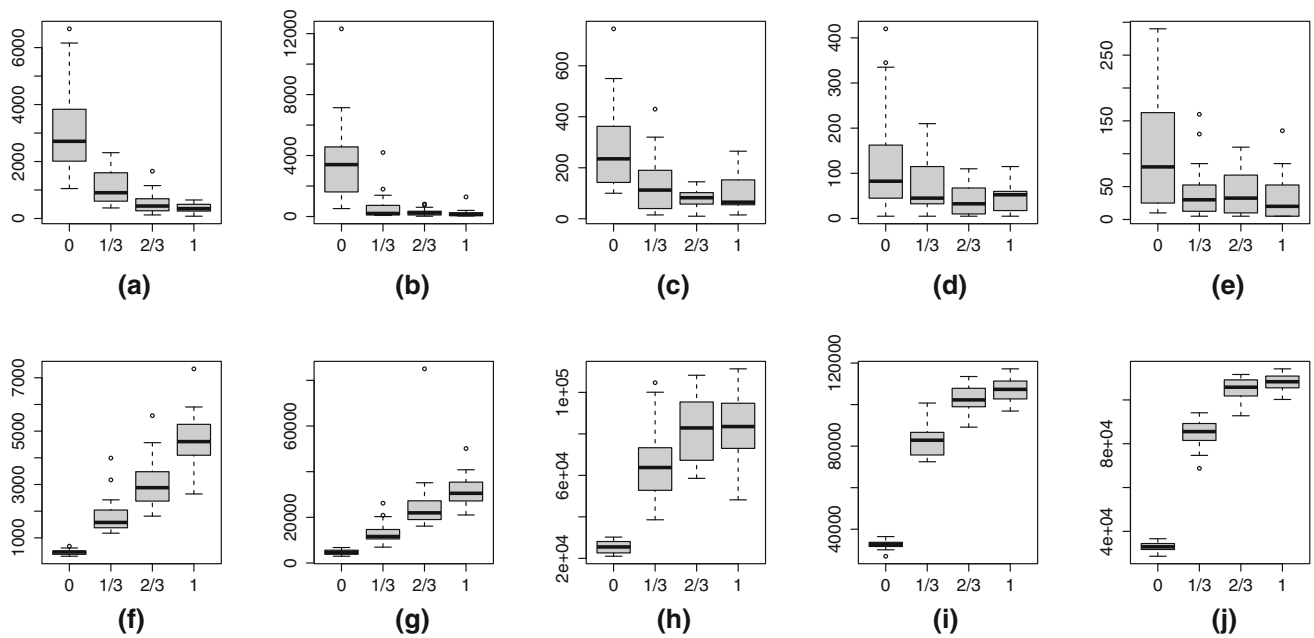


Fig. 12 MPSRF_M(1.1), panels (a–e), and the ESS of θ_0 , plots (f–j), for the CP

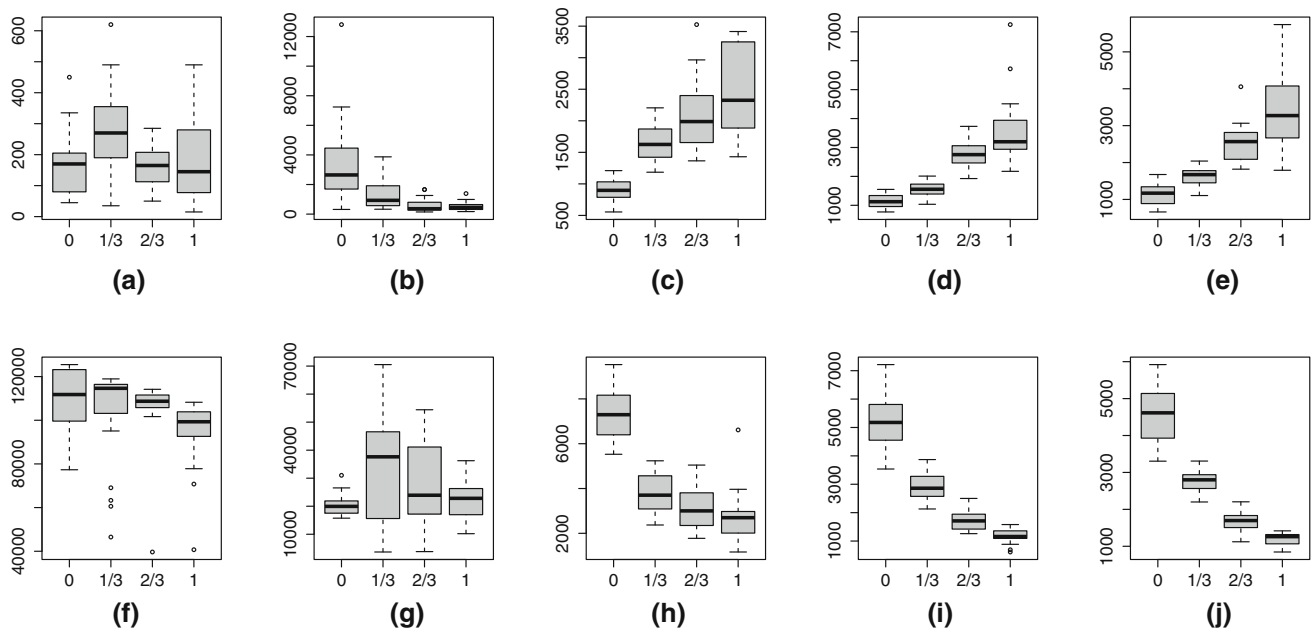


Fig. 13 MPSRF_M(1.1), panels (a–e), and the ESS of θ_0 , plots (f–j), for the NCP

16. The mean and standard deviation for the 50 data sites is 26.47 and 5.00 $\mu\text{g}/\text{m}^3$ respectively.

In addition to the observed data, we have output from the Air Quality Unified Model (AQUM), a numerical model giving air pollution predictions at 1-km grid cell resolution (Savage et al. 2013). The AQUM is used as a covariate in the model, where $x(s)$ is the AQUM output for the grid cell containing s . Therefore, we use a downscaler model as employed by Berrocal et al. (2010).

To stabilise the variance, we model the data on the square root scale. However, in order not to underestimate their variability, predictions are obtained on the original scale.

We fit model (1) with $p = 2$ and so we have two processes, an intercept and a slope process. Each process has a corresponding global variance and decay parameters and so $\theta = (\theta_0, \theta_1)^T$, $\sigma^2 = (\sigma_0^2, \sigma_1^2)^T$ and $\phi = (\phi_0, \phi_1)^T$. In addition, we have the data variance, σ_ϵ^2 . For the prior distribution of θ , we let $m = (0, 0)^T$ and $v_0 = v_1 = 10^4$. We

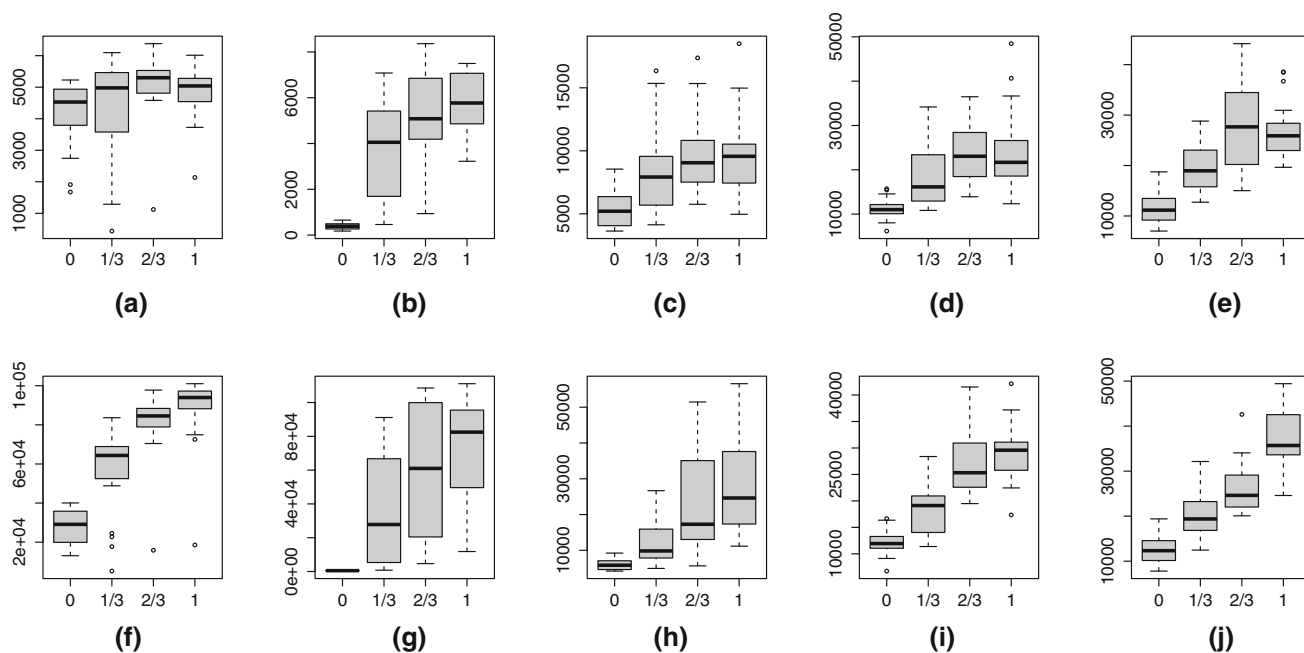


Fig. 14 ESS of σ_0^2 , panels (a–e), and σ_ϵ^2 , panels (f–j) for the CP

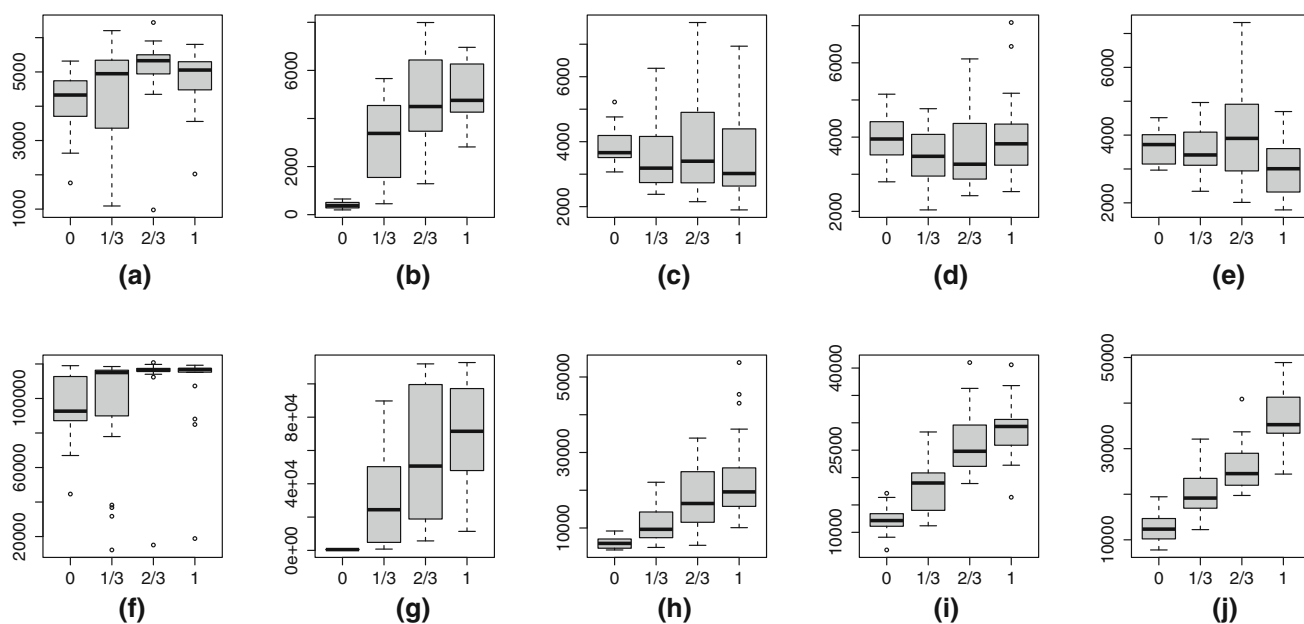


Fig. 15 ESS of σ_0^2 , panels (a–e), and σ_ϵ^2 , panels (f–j) for the NCP

let $a_0 = a_1 = a_\epsilon = 2$ and $b_0 = b_1 = b_\epsilon = 1$, so that each variance parameter is assigned an $IG(2, 1)$ prior distribution.

We begin by estimating the decay parameters using an empirical Bayes method by performing a grid search over a small number of values, then choosing those values that minimise some calibration criterion. This is a common approach adopted by many authors, e.g., [Sahu et al. \(2011\)](#), [Berrocal et al. \(2010\)](#), [Sahu et al. \(2007\)](#) since [Zhang \(2004\)](#)

showed that it is not possible to consistently estimate these in a model with Matérn covariance function in the presence of other unknown parameters. In our Bayesian inference setting, this will imply weak identifiability in the posterior distribution when non-informative prior distributions are assumed.

The greatest distance between any two of the 70 monitoring stations is 96.2 kilometres (km), and so we select values

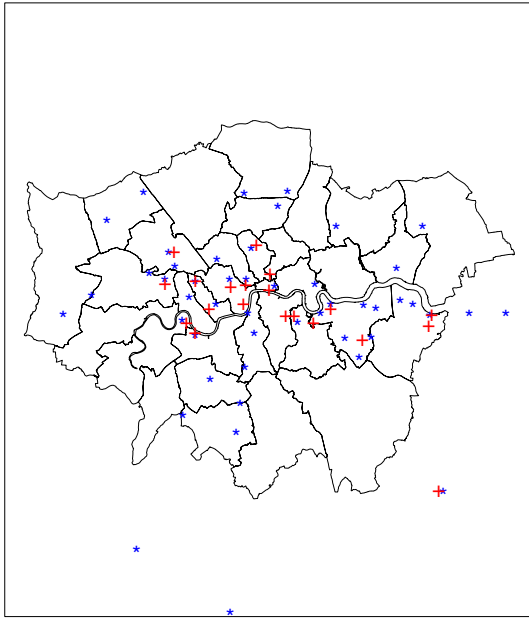


Fig. 16 Sampling locations for PM10 concentration data. *Blue stars* indicate locations used for model fitting and *red crosses* indicate locations used for model validation

of ϕ_0 and ϕ_1 corresponding to effective ranges of 5, 10, 25, 50 and 100 km. Predictions are made at the 20 validation sites and we compute the values of three measures of prediction error; the mean absolute prediction error (MAPE), the root mean square prediction error (RMSPE) and the continuous ranked probability score (CRPS), see for example [Gneiting et al. \(2007\)](#).

For each of the 25 pairs of spatial decay parameters, we generate a single chain of 25,000 iterations and discard the first 5,000. The values of the MAPE, RMSPE and CRPS are given in Table 1. Recall that d_0 and d_1 denote the effective ranges implied by ϕ_0 and ϕ_1 respectively. By all three measures, the prediction error is minimised when $d_0 = 5$ and $d_1 = 50$, and so our estimates for the spatial decay parameters are

$$\hat{\phi}_0 = -\log(0.05)/5 \approx 0.6, \quad \hat{\phi}_1 = -\log(0.05)/50 \approx 0.06.$$

We now compare the sampling efficiencies of the CP and NCP for the London PM10 data when the values of the decay parameters are fixed at the above optimal values. For each parameterisation, we generate five Markov chains of length 25,000 from the same widely dispersed starting values. The MPSRF_M(1.1) and the ESS of $\theta = (\theta_0, \theta_1)'$, $\sigma^2 = (\sigma_0^2, \sigma_1^2)'$ and σ_ϵ^2 are computed and given in Table 2. We can see that the CP requires far fewer iterations for the MPSRF to drop below 1.1 than the NCP: 275 versus 1985. The ESS of the mean parameters is greater for the CP than the NCP, especially

Table 1 Prediction error for different combinations of d_0 and d_1

d_0	d_1	MAPE	RMSPE	CRPS
5	5	5.494	6.224	3.720
	10	5.476	6.200	3.704
	25	5.416	6.164	3.701
	50	5.375	6.126	3.695
	100	5.418	6.160	3.735
10	5	5.534	6.281	3.740
	10	5.497	6.256	3.728
	25	5.480	6.230	3.733
	50	5.436	6.207	3.736
	100	5.452	6.193	3.753
25	5	5.618	6.534	3.878
	10	5.585	6.492	3.862
	25	5.492	6.351	3.801
	50	5.485	6.330	3.809
	100	5.476	6.314	3.828
50	5	5.620	6.741	4.003
	10	5.615	6.711	4.002
	25	5.549	6.572	3.944
	50	5.505	6.500	3.924
	100	5.499	6.470	3.938
100	5	5.644	6.910	4.129
	10	5.586	6.820	4.093
	25	5.541	6.658	4.035
	50	5.491	6.595	4.015
	100	5.482	6.581	4.046

for θ_1 reflecting the stronger spatial correlation for the slope process and the estimate for δ_1 . For the variance parameters σ_0^2 and σ_1^2 , the ESS for the CP is nearly double that of the NCP. The NCP achieves better mixing in the σ_ϵ^2 coordinate than does the CP.

For making inference, we generate a single long chain for 50,000 iterations and discard the first 10,000. Parameter estimates and their 95% credible intervals are given in Table 3. An estimate for the global intercept θ_0 of 5.148 reflects that the data are modelled on the square root scale. The global regression parameter θ_1 for the AQUM output is not significant given the inclusion in the model of the regression process. Of particular interest to us are the estimates of the variance parameters. Weak spatial correlation in the intercept process means that the estimate for σ_0^2 is the smallest of the three variance parameters. More of the spatial variability is explained by the intercept process, and so the estimate for σ_1^2 is the greatest of the three variance parameters. We also include estimates and 95% credible intervals for the variance ratios: $\delta_0 = \sigma_0^2/\sigma_\epsilon^2$ and $\delta_1 = \sigma_1^2/\sigma_\epsilon^2$.

Table 2 MPSRF_M(1.1) and the ESS of the model parameters

	MPSRF _M (1.1)	ESS θ_0	ESS θ_1	ESS σ_0^2	ESS σ_1^2	ESS σ_ϵ^2
CP	275	22,161	63,137	31,876	16,689	24,171
NCP	1985	15,596	2958	15,407	8275	31,806

Table 3 Parameter estimates and their 95% credible intervals (CI)

Parameter	Estimate	95% CI
θ_0	5.148	(4.942, 5.352)
θ_1	0.093	(−0.396, 0.572)
σ_0^2	0.172	(0.093, 0.299)
σ_1^2	0.224	(0.101, 0.469)
σ_ϵ^2	0.177	(0.096, 0.307)
δ_0	1.070	(0.412, 2.268)
δ_1	1.376	(0.487, 3.256)

5 Conclusion

We have compared the efficiencies of the CP and the NCP of spatial models. We find that in addition to the ratio of the variance parameters, the correlation structure between the random effects plays a key role in determining the rate of convergence.

For known variance and correlation parameters, the exact rate of convergence has been examined. We have shown that for spatial models with an exponential correlation function, increasing the variance of the random effects relative to that of the data, as well as increasing the strength of correlation, works to hasten the convergence of the CP but slows the convergence of the NCP. However, when the covariance matrix is tapered to remove long-range correlation, convergence for the CP is hindered, but convergence for the NCP is favoured. Introducing geometric anisotropy to strengthen the correlation in one direction has, for randomly selected locations, a similar effect to strengthening it in all directions; the CP is favoured and the NCP hindered. Both these results are consistent with the notion that the performance of CP is improved in the presence of greater spatial correlation but the performance of the NCP is worsened. We have seen that, as the smoothness parameter in the Matérn correlation function is increased, both the CP and the NCP are slower to converge, and in the presence of moderate spatial correlation, both parameterisations will fail to converge when the sample size is large enough.

When the variance parameters are unknown, the sampling efficiencies of the parameterisations are compared via the MPSRF_M(1.1) and the ESS of the unknown model parameters. We have seen that the relationships between the sampling efficiencies of the respective parameterisations, and the ratio of the variance parameters and the strength of spa-

tial correlation, still hold for unknown variance parameters. The CP performs better when the data precision is relatively high and when the correlation is strong. In contrast to this, the NCP performs best when the data are less informative and the correlation is weak.

Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

Appendix 1: Computing the exact convergences for the CP and NCP

For Gibbs samplers with Gaussian target distributions with known precision matrices, we have analytic results for the exact convergence rate (Roberts and Sahu 1997). We let ξ denote the set of all mean parameters in the model: $\xi = (\beta^T, \theta^T)^T$. Suppose that $\xi | y \sim N(\mu, \Sigma)$, and let $Q = \Sigma^{-1}$ be the posterior precision matrix (PPM). Further suppose that ξ is partitioned into l blocks for updating within the Gibbs sampler. To compute the convergence rate first, partition Q according to the l blocks where the ij th block is denoted by Q_{ij} , $i, j = 1, \dots, l$.

Let $A = I - \text{diag}(Q_{11}^{-1}, \dots, Q_{ll}^{-1})Q$ and $F = (I - L)^{-1}U$, where L is the block lower triangular matrix of A , and $U = A - L$. Roberts and Sahu (1997) show that the Markov chain induced by the Gibbs sampler with components block updated according to the above blocking scheme, has a Gaussian transition density with mean $E\{\xi^{(t+1)}|\xi^{(t)}\} = F\xi^{(t)} + f$, where $f = (I - F)\mu$ and covariance matrix $\Sigma - F\Sigma F^T$. Their observation leads to the following theorem:

Theorem 5.1 (Roberts and Sahu 1997, Theorem 1) A Markov chain with transition density

$$N\left\{F\xi^{(t)} + f, \Sigma - F\Sigma F^T\right\},$$

has a convergence rate equal to the maximum modulus eigenvalue of F .

Corollary 5.2 If we update ξ in two blocks so that $l = 2$, then

$$Q = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix}, \quad F = \begin{pmatrix} 0 & -Q_{11}^{-1}Q_{12} \\ 0 & Q_{22}^{-1}Q_{21}Q_{11}^{-1}Q_{12} \end{pmatrix},$$

and the convergence rate is the maximum modulus eigenvalue of

$$F_{22} = Q_{22}^{-1}Q_{21}Q_{11}^{-1}Q_{12}.$$

We use Theorem 5.1 to compare the convergence rates of Gibbs samplers associated with the CP and the NCP. First, we must compute the PPM for each parameterisation. To identify the PPM for the CP write

$$\begin{aligned} \pi(\tilde{\beta}, \theta | y) &\propto \pi(Y | \tilde{\beta})\pi(\tilde{\beta} | \theta)\pi(\theta) \\ &\propto \exp \left\{ -\frac{1}{2} \left[(Y - X_1\tilde{\beta})^T C_1^{-1} (Y - X_1\tilde{\beta}) \right. \right. \\ &\quad \left. \left. + (\tilde{\beta} - X_2\theta)^T C_2^{-1} (\tilde{\beta} - X_2\theta) \right. \right. \\ &\quad \left. \left. + (\theta - m)^T C_3^{-1} (\theta - m) \right] \right\} \\ &= \exp \left\{ -\frac{1}{2} \left[\dots + \tilde{\beta}^T (X_1^T C_1^{-1} X_1 + C_2^{-1}) \tilde{\beta} \right. \right. \\ &\quad \left. \left. - 2\tilde{\beta}^T C_2^{-1} X_2 \theta \right. \right. \\ &\quad \left. \left. + \theta^T (X_2^T C_2^{-1} X_2 + C_3^{-1}) \theta + \dots \right] \right\}, \end{aligned}$$

where the last equation only includes the terms containing both $\tilde{\beta}$ and θ . Therefore, the posterior precision matrix for the CP is given by

$$Q^c = \begin{pmatrix} X_1^T C_1^{-1} X_1 + C_2^{-1} & -C_2^{-1} X_2 \\ -X_2^T C_2^{-1} & X_2^T C_2^{-1} X_2 + C_3^{-1} \end{pmatrix}.$$

By corollary 5.2, if we update all random effects as one block and all global effects as another, then the convergence rate for the CP is the maximum modulus eigenvalue of

$$F_{22}^c = \left(X_2^T C_2^{-1} X_2 + C_3^{-1} \right)^{-1} X_2^T C_2^{-1} \left(X_1^T C_1^{-1} X_1 + C_2^{-1} \right)^{-1} C_2^{-1} X_2.$$

Similarly for the NCP, we find the PPM by writing

$$\begin{aligned} \pi(\beta, \theta | y) &\propto \pi(Y | \beta, \theta)\pi(\beta | \theta)\pi(\theta) \\ &\propto \exp \left\{ -\frac{1}{2} \left[(Y - X_1\beta - X_1X_2\theta)^T C_1^{-1} \right. \right. \\ &\quad \left. \left. (Y - X_1\beta - X_1X_2\theta) + \right. \right. \end{aligned}$$

$$\begin{aligned} &\left. + \beta^T C_2^{-1} \beta + (\theta - m)^T C_3^{-1} (\theta - m) \right] \Big\} \\ &= \exp \left\{ -\frac{1}{2} \left[\dots + \beta^T (X_1^T C_1^{-1} X_1 + C_2^{-1}) \beta \right. \right. \\ &\quad \left. \left. + 2\beta^T X_1^T C_2^{-1} X_1 X_2 \theta \right. \right. \\ &\quad \left. \left. + \theta^T (X_2^T X_1^T C_1^{-1} X_1 X_2 + C_3^{-1}) \theta + \dots \right] \right\}, \end{aligned}$$

and hence, we have

$$Q^{nc} = \begin{pmatrix} X_1^T C_1^{-1} X_1 + C_2^{-1} & X_1^T C_1^{-1} X_1 X_2 \\ X_2^T X_1^T C_1^{-1} X_1 & X_2^T X_1^T C_1^{-1} X_1 X_2 + C_3^{-1} \end{pmatrix}.$$

By Corollary 5.2, the convergence rate of the Gibbs sampler for the NCP is the maximum modulus eigenvalue of

$$F_{22}^{nc} = \left(X_2^T X_1^T C_1^{-1} X_1 X_2 + C_3^{-1} \right)^{-1} X_2^T X_1^T C_1^{-1} X_1 \left(X_1^T C_1^{-1} X_1 + C_2^{-1} \right)^{-1} X_1^T C_1^{-1} X_1 X_2.$$

Appendix 2: Full conditional distributions

(a) Posterior distributions for the CP

In this section, we give the joint posterior and full conditional distributions for the CP of model (1). We denote by $\sigma^2 = (\sigma_0^2, \dots, \sigma_{p-1}^2)^T$ the vector containing the variance parameters of the random effects and let $\phi = (\phi_0, \dots, \phi_{p-1})^T$ contain the decay parameters. We let $\xi = (\tilde{\beta}^T, \theta^T, \sigma^2, \sigma_\epsilon^2, \phi^T)^T$ contain all np random effects, p global effects, $p+1$ variance parameters and p decay parameters. The joint posterior distribution of ξ is given by

$$\begin{aligned} \pi(\xi | y) &\propto \pi(Y | \tilde{\beta}, \sigma_\epsilon^2)\pi(\tilde{\beta} | \theta, \sigma^2, \phi)\pi(\theta | \sigma^2)\pi(\sigma^2)\pi(\sigma_\epsilon^2)\pi(\phi) \\ &\propto \prod_{k=0}^{p-1} (\sigma_k^2)^{-(n/2+1/2+a_k+1)} |R_k|^{-1/2} (\sigma_\epsilon^2)^{-(n/2+a_\epsilon+1)} \\ &\quad \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left[\left(Y - \sum_{k=0}^{p-1} D_k \tilde{\beta}_k \right)^T \left(Y - \sum_{k=0}^{p-1} D_k \tilde{\beta}_k \right) + 2b_\epsilon \right] \right\} \\ &\quad \exp \left\{ -\frac{1}{2} \sum_{k=0}^{p-1} (\tilde{\beta}_k - \theta_k \mathbf{1})^T \Sigma_k^{-1} (\tilde{\beta}_k - \theta_k \mathbf{1}) \right\} \\ &\quad \exp \left\{ -\frac{1}{2} \sum_{k=0}^{p-1} \frac{1}{\sigma_k^2} \left(\frac{(\theta_k - m_k)^2}{v_k} + 2b_k \right) \right\} \prod_{k=0}^{p-1} \pi(\phi_k), \quad (29) \end{aligned}$$

where D_0 is defined to be the identity matrix I .

We use Gibbs sampling to sample from $\pi(\xi | y)$ for the CP, given in (29). We assume that the random effects will be block updated according to their process, i.e. we jointly update the n -dimensional vector β_k , for $k = 0, \dots, p-1$.

All other parameters in ξ are updated as single univariate components. The full conditional distributions that we need for the CP are given below:

- The full conditional distribution for the centred spatially correlated random effects $\tilde{\beta}_k, k = 0, \dots, p-1$, is

$$\tilde{\beta}_k | \tilde{\beta}_{-k}, \theta, \sigma^2, \sigma_\epsilon^2, \phi, \mathbf{y} \sim N(\mathbf{m}_k^*, \Sigma_k^*),$$

where we denote by $\tilde{\beta}_{-k}$ the vector of all random effects $\tilde{\beta}$ without the realisations of the k th process $\tilde{\beta}_k$ and

$$\Sigma_k^* = \left(\frac{1}{\sigma_\epsilon^2} \mathbf{D}_k^T \mathbf{D}_k + \Sigma_k^{-1} \right)^{-1},$$

$$\mathbf{m}_k^* = \Sigma_k^* \left[\frac{1}{\sigma_\epsilon^2} \mathbf{D}_k \left(\mathbf{y} - \sum_{\substack{j=0 \\ j \neq k}}^{p-1} \mathbf{D}_j \tilde{\beta}_j \right) + \Sigma_k^{-1} \theta_k \mathbf{1} \right].$$

- The full conditional distribution for the global effects $\theta_k, k = 0, \dots, p-1$, for the CP is

$$\theta_k | \tilde{\beta}, \theta_{-k}, \sigma^2, \sigma_\epsilon^2, \phi, \mathbf{y} \sim N(m_k^*, v_k^*),$$

where

$$v_k^* = \left(\mathbf{1}^T \Sigma_k^{-1} \mathbf{1} + \frac{1}{\sigma_k^2 v_k} \right)^{-1},$$

$$m_k^* = v_k^* \left(\mathbf{1}^T \Sigma_k^{-1} \tilde{\beta}_k + \frac{m_k}{\sigma_k^2 v_k} \right).$$

- The full conditional distribution for the random effects variance $\sigma_k^2, k = 0, \dots, p-1$, for the CP is

$$\sigma_k^2 | \tilde{\beta}, \theta, \sigma_{-k}^2, \sigma_\epsilon^2, \phi, \mathbf{y} \sim IG \left\{ \frac{n+1}{2} + a_k, \frac{1}{2} \left[\left(\tilde{\beta}_k - \theta_k \mathbf{1} \right)^T \mathbf{R}_k^{-1} \left(\tilde{\beta}_k - \theta_k \mathbf{1} \right) + \frac{(\theta_k - m_k)^2}{v_k} + 2b_k \right] \right\}.$$

- The full conditional distribution for data variance σ_ϵ^2 for the CP is

$$\sigma_\epsilon^2 | \tilde{\beta}, \theta, \sigma^2, \phi, \mathbf{y} \sim IG \left\{ \frac{n}{2} + a_\epsilon, \frac{1}{2} \left[\left(\mathbf{Y} - \sum_{k=0}^{p-1} \mathbf{D}_k \tilde{\beta}_k \right)^T \left(\mathbf{Y} - \sum_{k=0}^{p-1} \mathbf{D}_k \tilde{\beta}_k \right) + 2b_\epsilon \right] \right\}.$$

Posterior distributions for the NCP

We now look at the joint posterior and full conditional distributions of the model parameters for the NCP. For the NCP, we have $\xi = (\beta^T, \theta^T, \sigma^2, \sigma_\epsilon^2, \phi^T)^T$, and

$$\begin{aligned} \pi(\xi | \mathbf{y}) &\propto \pi(\mathbf{Y} | \beta, \theta, \sigma_\epsilon^2) \\ &\quad \pi(\beta | \sigma^2, \phi) \pi(\theta | \sigma^2) \pi(\sigma^2) \pi(\sigma_\epsilon^2) \pi(\phi) \\ &\propto \prod_{k=0}^{p-1} (\sigma_k^2)^{-(n/2+1/2+a_k+1)} |\mathbf{R}_k|^{-1/2} (\sigma_\epsilon^2)^{-(n/2+a_\epsilon+1)} \\ &\quad \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left[\left(\mathbf{Y} - \sum_{k=0}^{p-1} (\mathbf{D}_k \beta_k + \mathbf{x}_k \theta_k) \right)^T \left(\mathbf{Y} - \sum_{k=0}^{p-1} (\mathbf{D}_k \beta_k + \mathbf{x}_k \theta_k) \right) + 2b_\epsilon \right] \right\} \\ &\quad \exp \left\{ -\frac{1}{2} \sum_{k=0}^{p-1} \beta_k^T \Sigma_k^{-1} \beta_k \right\} \\ &\quad \exp \left\{ -\frac{1}{2} \sum_{k=0}^{p-1} \frac{1}{\sigma_k^2} \left(\frac{(\theta_k - m_k)^2}{v_k} + 2b_k \right) \right\} \\ &\quad \prod_{k=0}^{p-1} \pi(\phi_k), \end{aligned} \quad (30)$$

where we define \mathbf{x}_0 to be the vector of ones.

The full conditional distributions for the NCP are given below:

- The full conditional distribution for the non-centred spatially correlated random effects $\beta_k, k = 0, \dots, p-1$, is

$$\beta_k | \beta_{-k}, \theta, \sigma^2, \sigma_\epsilon^2, \phi, \mathbf{y} \sim N(\mathbf{m}_k^*, \Sigma_k^*),$$

where

$$\Sigma_k^* = \left(\frac{1}{\sigma_\epsilon^2} \mathbf{D}_k^T \mathbf{D}_k + \Sigma_k^{-1} \right)^{-1},$$

$$\mathbf{m}_k^* = \Sigma_k^* \left[\frac{1}{\sigma_\epsilon^2} \mathbf{x}_k^T \left(\mathbf{y} - \sum_{\substack{j=0 \\ j \neq k}}^{p-1} \mathbf{D}_j \beta_j - \sum_{j=0}^{p-1} \mathbf{x}_j \theta_j \right) \right].$$

- The full conditional distribution for the global effects $\theta_k, k = 0, \dots, p-1$, for the NCP is

$$\theta_k | \beta, \theta_{-k}, \sigma^2, \sigma_\epsilon^2, \phi, \mathbf{y} \sim N(m_k^*, v_k^*),$$

where

$$v_k^* = \left(\frac{1}{\sigma_\epsilon^2} \mathbf{x}_k^T \mathbf{x}_k + \frac{1}{\sigma_k^2 v_k} \right)^{-1},$$

$$m_k^* = v_k^* \left[\frac{1}{\sigma_\epsilon^2} \mathbf{x}_k^T \left(\mathbf{y} - \sum_{j=0}^{p-1} \mathbf{D}_j \boldsymbol{\beta}_j - \sum_{\substack{j=0 \\ j \neq k}}^{p-1} \mathbf{x}_j \theta_j \right) + \frac{m_k}{\sigma_k^2 v_k} \right].$$

- The full conditional distribution for the random effects variance σ_k^2 , $k = 0, \dots, p-1$, for the NCP is

$$\sigma_k^2 | \boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\sigma}_{-k}^2, \boldsymbol{\phi}, \mathbf{y} \sim IG \left\{ \frac{n+1}{2} + a_k, \frac{1}{2} \left(\boldsymbol{\beta}_k^T \mathbf{R}_k^{-1} \boldsymbol{\beta}_k + \frac{(\theta_k - m_k)^2}{v_k} + 2b_k \right) \right\}.$$

- The full conditional distribution for the data variance σ_ϵ^2 for the NCP is

$$\sigma_\epsilon^2 | \boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\sigma}^2, \boldsymbol{\phi}, \mathbf{y} \sim IG \left\{ \frac{n}{2} + a_\epsilon, \frac{1}{2} \left[\left(\mathbf{Y} - \sum_{k=0}^{p-1} (\mathbf{D}_k \boldsymbol{\beta}_k + \mathbf{x}_k \theta_k) \right)^T \left(\mathbf{Y} - \sum_{k=0}^{p-1} (\mathbf{D}_k \boldsymbol{\beta}_k + \mathbf{x}_k \theta_k) \right) + 2b_\epsilon \right] \right\}.$$

References

- Abramowitz, M., Stegun, I.A.: Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables. Courier Dover Publications, New York (1972)
- Apputhurai, P., Stephenson, A.G.: Spatiotemporal hierarchical modelling of extreme precipitation in Western Australia using anisotropic Gaussian random fields. *Environ. Ecol. Stat.* **20**(4), 667–677 (2013)
- Banerjee, S., Gelfand, A.E., Finley, A.O., Sang, H.: Gaussian predictive process models for large spatial data sets. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **70**(4), 825–848 (2008)
- Banerjee, S., Finley, A.O., Waldmann, P., Ericsson, T.: Hierarchical spatial process models for multiple traits in large genetic trials. *J. Am. Stat. Assoc.* **105**(490), 506–521 (2010)
- Banerjee, S., Carlin, B.P., Gelfand, A.E.: Hierarchical Modeling and Analysis for Spatial Data, 2nd edn. CRC Press, Boca Raton (2015)
- Bass, M.R., Sahu, S.K.: Dynamically updated spatially varying parameterizations of hierarchical Bayesian models for spatially correlated data. Submitted (2016)
- Berrocal, V.J., Gelfand, A.E., Holland, D.M.: A spatio-temporal down-scaler for output from numerical models. *J. Agric. Biol. Environ. Stat.* **15**(2), 176–197 (2010)
- Brooks, S.P., Gelman, A.: General methods for monitoring convergence of iterative simulations. *J. Comput. Graph. Stat.* **7**(4), 434–455 (1998)
- Cressie, N., Johannesson, G.: Fixed rank kriging for very large spatial data sets. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **70**(1), 209–226 (2008)
- Ecker, M.D., Gelfand, A.E.: Bayesian modeling and inference for geometrically anisotropic spatial data. *Math. Geol.* **31**(1), 67–83 (1999)
- Filippone, M., Zhong, M., Girolami, M.: A comparative evaluation of stochastic-based inference methods for gaussian process models. *Mach. Learn.* **93**(1), 93–114 (2013)
- Furrer, R., Genton, M.G., Nychka, D.: Covariance tapering for interpolation of large spatial datasets. *J. Comput. Graph. Stat.* **15**(3), 502–523 (2006)
- Gelfand, A.E., Smith, A.F.: Sampling-based approaches to calculating marginal densities. *J. Am. Stat. Assoc.* **85**(410), 398–409 (1990)
- Gelfand, A.E., Sahu, S.K., Carlin, B.P.: Efficient parameterisations for normal linear mixed models. *Biometrika* **82**(3), 479–488 (1995)
- Gelfand, A.E., Sahu, S.K., Carlin, B.P.: Efficient parameterizations for generalized linear mixed models, (with discussion). In: Bernardo, J.M., Berger, J.O., Dawid, A.P., Smith, A.F.M. (eds.) *Bayesian Statistics 5*, pp. 165–180. Oxford University Press, Oxford (1996)
- Gelfand, A.E., Kim, H.-J., Sirmans, C., Banerjee, S.: Spatial modeling with spatially varying coefficient processes. *J. Am. Stat. Assoc.* **98**(462), 387–396 (2003)
- Gelman, A., Rubin, D.B.: Inference from iterative simulation using multiple sequences. *Stat. Sci.* **7**(4), 457–472 (1992)
- Gneiting, T., Balabdaoui, F., Raftery, A.E.: Probabilistic forecasts, calibration and sharpness. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **69**(2), 243–268 (2007)
- Handcock, M.S., Stein, M.L.: A Bayesian analysis of kriging. *Technometrics* **35**(4), 403–410 (1993)
- Harville, D.A.: Matrix Algebra from a Statistician's Perspective. Springer, New York (1997)
- Horn, R.A., Johnson, C.R.: Matrix Analysis. Cambridge University Press, Cambridge (2012)
- Huerta, G., Sansó, B., Stroud, J.R.: A spatiotemporal model for Mexico City ozone levels. *J. R. Stat. Soc. Ser. C (Appl. Stat.)* **53**(2), 231–248 (2004)
- Imai, K., van Dyk, D.: A bayesian analysis of the multinomial probit model using marginal data augmentation. *J. Econom.* **124**, 311–334 (2005)
- Kaufman, C.G., Schervish, M.J., Nychka, D.W.: Covariance tapering for likelihood-based estimation in large spatial data sets. *J. Am. Stat. Assoc.* **103**(484), 1545–1555 (2008)
- Liu, J.S., Wu, Y.N.: Parameter expansion for data augmentation. *J. Am. Stat. Assoc.* **94**, 1264–1274 (1999)
- Matérn, B.: Spatial Variation, 2nd edn. Springer, Berlin (1986)
- Meyn, S.P., Tweedie, R.L.: Markov Chains and Stochastic Stability. Springer, London (1993)
- Papaspiliopoulos, O., Roberts, G.: Stability of the Gibbs sampler for Bayesian hierarchical models. *Ann. Stat.* **36**(1), 95–117 (2008)
- Papaspiliopoulos, O., Roberts, G.O., Sköld, M.: Non-centered parameterisations for hierarchical models and data augmentation (with discussion). *Bayesian Statistics 7* (Bernardo, JM and Bayarri, MJ and Berger, JO and Dawid, AP and Heckerman, D and Smith, AFM and West, M): Proceedings of the Seventh Valencia International Meeting, pp. 307–326. Oxford University Press, USA (2003)
- Rasmussen, C.E., Williams, C.K.I.: Gaussian Processes for Machine Learning. MIT press, Cambridge, MA (2006)
- Robert, C.O., Casella, G.: Monte Carlo Statistical Methods, 2nd edn. Springer, New York (2004)
- Roberts, G.O., Sahu, S.K.: Updating schemes, correlation structure, blocking and parameterization for the Gibbs sampler. *J. R. Stat. Soc. Ser. B (Stat. Methodol.)* **59**(2), 291–317 (1997)
- Sahu, S.K., Roberts, G.O.: On convergence of the em algorithm and the gibbs sampler. *Stat. Comput.* **9**, 55–64 (1999)

- Sahu, S.K., Gelfand, A.E., Holland, D.M.: High resolution space-time ozone modeling for assessing trends. *J. Am. Stat. Assoc.* **102**(480), 1221–1234 (2007)
- Sahu, S.K., Gelfand, A.E., Holland, D.M.: Fusing point and areal level space-time data with application to wet deposition. *J. R. Stat. Soc. Ser. C (Appl. Stat.)* **59**(1), 77–103 (2010)
- Sahu, S.K., Yip, S., Holland, D.M.: A fast Bayesian method for updating and forecasting hourly ozone levels. *Environ. Ecol. Stat.* **18**(1), 185–207 (2011)
- Savage, N., Agnew, P., Davis, L., Ordóñez, C., Thorpe, R., Johnson, C., O'Connor, F., Dalvi, M.: Air quality modelling using the met office unified model (aqum os24-26): model description and initial evaluation. *Geosci. Model Dev.* **6**(2), 353–372 (2013)
- Schabenberger, O., Gotway, C.A.: *Statistical Methods for Spatial Data Analysis*. CRC Press, Boca Raton (2004)
- Stein, M.L.: *Interpolation of Spatial Data: Some Theory for Kriging*. Springer, New York (1999)
- van Dyk, D., Meng, X.-L.: The art of data augmentation (with discussion). *J. Comput. Graph. Stat.* **10**, 1–50 (2001)
- Yu, Y., Meng, X.-L.: To center or not to center: That is not the question: an Ancillarity-Sufficiency Interweaving Strategy (ASIS) for boosting MCMC efficiency. *J. Comput. Graph. Stat.* **20**(3), 531–570 (2011)
- Zhang, H.: Inconsistent estimation and asymptotically equal interpolations in model-based geostatistics. *J. Am. Stat. Assoc.* **99**(465), 250–261 (2004)
- Zimmerman, D.L.: Another look at anisotropy in geostatistics. *Math. Geol.* **25**(4), 453–470 (1993)