

# Dynamically updated spatially varying parameterisations of hierarchical Bayesian models for spatially correlated data

Mark Bass and Sujit Sahu  
University of Southampton, UK

June 14, 2016

## Abstract

Fitting hierarchical Bayesian models to large spatially correlated data sets using Markov chain Monte Carlo (MCMC) techniques is computationally expensive. Complicated covariance structures of the underlying spatial processes, together with high dimensional parameter space, mean that the number of calculations required grows cubically with the number of spatial locations at each MCMC iteration. This necessitates the need for efficient model parameterisations that hasten the convergence and improve the mixing of the associated algorithms. We consider partially centred parameterisations (PCPs) of hierarchical models which lie on a continuum between what are known as the centred parameterisation (CP) and the noncentred parameterisation (NCP). By introducing a weight matrix we remove the conditional posterior correlation between the fixed and the random effects, and hence construct a PCP which achieves immediate convergence for a three stage model, based on multiple Gaussian processes with known covariance parameters. When the covariance parameters are unknown we dynamically update the parameterisation within the sampler. We show that the PCP is robust to the data set to be modelled and that it outperforms both the CP and the NCP. The dynamically updated PCP leads to a fully automated algorithm that is illustrated with a practical data set on ozone concentration levels. The example shows the effectiveness of allowing the parameterisation to vary spatially.

**Keywords:** Gibbs sampling; Hierarchical model; Parameterisation; Rate of convergence; Spatial model.

## 1 Introduction

There is a growing interest among researchers in spatially varying coefficient (SVC) models (Hamm *et al.*, 2015; Wheeler *et al.*, 2014; Finley *et al.*, 2011; Berrocal *et al.*, 2010; Gelfand *et al.*, 2003). Conditional independencies determined by the hierarchical structure of the model facilitate the construction of Gibbs sampling type algorithms for model fitting (Gelfand and Smith, 1990). A requirement of these algorithms is the repeated inversion of dense  $n \times n$  covariance matrices, an operation of order  $O(n^3)$  in computational complexity, for  $n$  spatial locations (Cressie and Johannesson, 2008). This, coupled with high posterior correlation between model parameters and weakly identified covariance parameters, means that fitting these models is challenging and computationally expensive. To mitigate the computational expense practitioners require efficient model fitting strategies that produce Markov chains which converge quickly to the posterior distribution and exhibit low autocorrelation between successive iterates.

Parameterisation of a hierarchical model is known to affect the performance of the Markov chain Monte Carlo (MCMC) method used for inference. For normal linear hierarchical models (NLHMs) the centred parameterisation (CP) yields an efficient Gibbs sampler when the

variance of the data model is low relative to that of the random effects, and the noncentred parameterisation (NCP) yields an efficient Gibbs sampler when the variance of the data model is relatively high (Papaspiliopoulos *et al.*, 2003; Gelfand *et al.*, 1995). Where the latent variables are realisations of a spatial process with an exponential correlation function, Bass and Sahu (2016) shows that increasing the strength of correlation improves the efficiency of the CP but degrades that of the NCP. Hence the sampling efficiency of the CP and the NCP is dependent upon the typically unknown variance and spatial correlation parameters and will therefore differ across different data sets. Consequently, deciding which parameterisation to employ can be problematic.

With the aim of developing a robust parameterisation for NLHMs, Papaspiliopoulos *et al.* (2003) consider the CP and the NCP as extremes of a family of partially centred parameterisations (PCPs). They find the optimal PCP which results in a Gibbs sampler that produces independent samples from the posterior distributions of the mean parameters, but again this is conditioned on the covariance parameters. The question we look to answer in this paper is can we create a parameterisation for spatial models that is robust to the data and does not require *a priori* knowledge of the model parameters, and hence can be routinely implemented?

To address this question we write a general SVC model as a three stage NLHM. The PCP is constructed by introducing a weight matrix that allows us to eliminate the conditional posterior correlation between the global and random effects. This in turn implies immediate convergence of the associated Gibbs sampler for known variance and correlation parameters, which collectively we will call the covariance parameters. When these parameters are unknown we propose dynamically updating the weight matrix, which leads to a parameterisation that is both spatially varying and dynamically updated within the Gibbs sampler. As it is necessary to invert an  $n \times n$  matrix to compute the weight matrix, implementing the PCP results in an algorithm with longer run times than those associated with the CP and NCP. Consequently, we recommend this method for modest sized spatial datasets of hundreds, but not thousands, of observations.

For an exponential correlation function we demonstrate how the weights of partial centering depend on the covariance parameters and the sampling locations, and hence vary over the spatial domain. We show that weights are higher when the data variance is relatively low and when the spatial correlation is high. Also, higher weights are given to more densely clustered sampling locations and where the covariate values are higher. In order to judge sampling efficiency of the PCP we use well known convergence diagnostics. The performance of the PCP is shown to be robust to changes in the covariance parameters of the data generating mechanism. Moreover, the PCP converges more quickly and produces posterior samples with lower autocorrelation than either the CP or the NCP.

A related approach is the interweaving algorithm proposed by Yu and Meng (2011). The algorithm results in a Gibbs sampler that is more efficient than the worst of the CP and NCP. However, the interweaving algorithm does not guarantee immediate convergence for known covariance parameters, whereas the PCP does, see Section 2.3. Another approach is to marginalise over the random effects, thus reducing the dimension of the posterior distribution. This method 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 SVC regression 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. The scheme proposed here is fully automated.

The rest of this paper is laid out as follows: Section 2 gives details of the general spatial model and the construction and properties of its PCP. Section 3 illustrates how the weights of partial centering are influenced by the model parameters and the sampling locations. Section 4

demonstrates the sampling efficiency of the PCP by applying it to simulated spatial data sets and compares its performance to the CP and the NCP. Section 5 applies the different model parameterisations to ozone concentration data from the State of California, USA. Section 6 contains some concluding remarks.

## 2 Partial centering of spatially varying coefficient models

### 2.1 Model specification

For data observed at a set of locations  $\mathbf{s}_1, \dots, \mathbf{s}_n$  we consider the normal linear model with spatially varying regression coefficients (Gelfand *et al.*, 2003):

$$Y(\mathbf{s}_i) = \theta_0 + \beta_0(\mathbf{s}_i) + \sum_{k=1}^{p-1} \{\theta_k + \beta_k(\mathbf{s}_i)\} x_k(\mathbf{s}_i) + \epsilon(\mathbf{s}_i), \quad i = 1, \dots, n. \quad (1)$$

We model errors  $\epsilon(\mathbf{s}_i)$  as independent and normally distributed with mean zero and variance  $\sigma_\epsilon^2$ . Spatially indexed observations

$\mathbf{Y} = \{Y(\mathbf{s}_1), \dots, Y(\mathbf{s}_n)\}^T$  are conditionally independent and normally distributed as

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

where  $\mathbf{x}(\mathbf{s}_i) = \{1, x_1(\mathbf{s}_i), \dots, x_{p-1}(\mathbf{s}_i)\}^T$  is a vector containing covariate information for site  $\mathbf{s}_i$  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 Gaussian process, denoted  $\beta_k(\mathbf{s}_i)$ , which are collected into a vector  $\boldsymbol{\beta}(\mathbf{s}_i) = \{\beta_0(\mathbf{s}_i), \dots, \beta_{p-1}(\mathbf{s}_i)\}^T$ . The  $n$  realisations of the Gaussian process associated with the  $k$ th covariate are given by

$$\boldsymbol{\beta}_k = \{\beta_k(\mathbf{s}_1), \dots, \beta_k(\mathbf{s}_n)\}^T \sim N(0, \boldsymbol{\Sigma}_k), \quad k = 0, \dots, p-1,$$

where

$$\boldsymbol{\Sigma}_k = \sigma_k^2 \mathbf{R}_k, \quad \text{and} \quad (\mathbf{R}_k)_{ij} = \text{corr}\{\beta_k(\mathbf{s}_i), \beta_k(\mathbf{s}_j)\}.$$

The form of the model given in (1) is the NCP. The CP is found by introducing the variables  $\tilde{\beta}_k(\mathbf{s}_i) = \theta_k + \beta_k(\mathbf{s}_i)$ . Therefore

$$\tilde{\boldsymbol{\beta}}_k = \{\tilde{\beta}_k(\mathbf{s}_1), \dots, \tilde{\beta}_k(\mathbf{s}_n)\}^T \sim N(\theta_k \mathbf{1}, \boldsymbol{\Sigma}_k).$$

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

$$\begin{aligned} \mathbf{Y} | \tilde{\boldsymbol{\beta}} &\sim N(\mathbf{X}_1 \tilde{\boldsymbol{\beta}}, \mathbf{C}_1) \\ \tilde{\boldsymbol{\beta}} | \boldsymbol{\theta} &\sim N(\mathbf{X}_2 \boldsymbol{\theta}, \mathbf{C}_2) \\ \boldsymbol{\theta} &\sim N(\mathbf{m}, \mathbf{C}_3), \end{aligned} \quad (2)$$

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  $\mathbf{x}_k = \{x_k(\mathbf{s}_1), \dots, x_k(\mathbf{s}_n)\}^T$ . 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  $n$ ,

$$\mathbf{X}_2 = \begin{pmatrix} \mathbf{1} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{1} \end{pmatrix}.$$

The  $p$  processes are assumed independent *a priori* and so  $\mathbf{C}_2$  is block diagonal where the  $k$ th block is  $\Sigma_k$ , and so

$$\mathbf{C}_2 = \begin{pmatrix} \Sigma_0 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma_1 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma_{p-1} \end{pmatrix}.$$

The global effects  $\boldsymbol{\theta} = (\theta_0, \theta_1, \dots, \theta_{p-1})^\top$  are assumed to be independent *a priori* with the  $k$ th element assigned a Gaussian prior distribution with mean  $m_k$  and variance  $\sigma_k^2 v_k$ , hence we write  $\theta_k \sim N(m_k, \sigma_k^2 v_k)$ . Therefore  $\mathbf{m} = (m_0, \dots, m_{p-1})^\top$  and

$$\mathbf{C}_3 = \begin{pmatrix} \sigma_0^2 v_0 & 0 & \cdots & 0 \\ 0 & \sigma_1^2 v_1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_{p-1}^2 v_{p-1} \end{pmatrix}.$$

We complete the model specification by assigning prior distributions to the covariance parameters. The realisations of the  $k$ th zero mean Gaussian process,  $\beta_k$ , have a prior covariance matrix given by  $\Sigma_k = \sigma_k^2 \mathbf{R}_k$ . This prior covariance matrix is shared by the  $k$ th centred Gaussian process,  $\tilde{\beta}_k$ . The prior distributions for the variance parameters are given by

$$\sigma_k^2 \sim IG(a_k, b_k), \quad k = 0, \dots, p-1, \quad \sigma_\epsilon^2 \sim IG(a_\epsilon, b_\epsilon),$$

where we write  $X \sim IG(a, b)$  if  $X$  has a density proportional to  $x^{-(a+1)}e^{-b/x}$ .

In this paper we consider only the exponential correlation function, which is a member of the Matérn family (Handcock and Stein, 1993; Matérn, 1986) and is widely applied to spatial processes (Sahu *et al.*, 2010; Berrocal *et al.*, 2010; Sahu *et al.*, 2007; Huerta *et al.*, 2004). Therefore entries of the  $\mathbf{R}_k$  are

$$(\mathbf{R}_k)_{ij} = \text{corr}\{\beta_k(\mathbf{s}_i), \beta_k(\mathbf{s}_j)\} = \exp(-\phi_k d_{ij}).$$

where  $d_{ij} = \|\mathbf{s}_i - \mathbf{s}_j\|$  denotes the distance between  $\mathbf{s}_i$  and  $\mathbf{s}_j$  and  $\phi_k$  controls the rate of decay of the correlation. Here the strength of correlation is characterised by the effective range,  $d_k$ , which is the distance such that  $\text{corr}\{\beta_k(\mathbf{s}_i), \beta_k(\mathbf{s}_j)\} = 0.05$ . For an exponential correlation function we have that

$$d_k = -\log(0.05)/\phi_k \approx 3/\phi_k.$$

It is not possible to consistently estimate all of the variance and decay parameters under vague prior distributions (Zhang, 2004), and so we do not sample from the posterior distributions of the decay parameters. For the simulation studies in Section 4 the decay parameters are fixed, thus helping us to examine their impact upon the sampling efficiency of the PCP. For the real data example in Section 5 we perform a grid search over a range of values, selecting those that offer the best out-of-sample predictions. A grid search is equivalent to placing a discrete uniform prior distribution upon the decay parameters and is a commonly adopted approach (Sahu *et al.*, 2011; Berrocal *et al.*, 2010; Sahu *et al.*, 2007).

## 2.2 Construction of the PCP

In Section 2.1 we consider two parameterisations, non-centred and centred, that differ by the prior mean of the spatial processes. We have  $\boldsymbol{\beta} \sim N(\mathbf{0}, \mathbf{C}_2)$  for the NCP and  $\tilde{\boldsymbol{\beta}} \sim N(\mathbf{X}_2 \boldsymbol{\theta}, \mathbf{C}_2)$

for the CP where a linear shift relates the two processes such that  $\tilde{\beta} = \beta + \mathbf{X}_2\theta$ . The PCP is formed by a partial shift, which is defined by

$$\beta^w = \tilde{\beta} - (\mathbf{I} - \mathbf{W})\mathbf{X}_2\theta, \quad (3)$$

and therefore the partially centred model is written as

$$\begin{aligned} \mathbf{Y} \mid \beta^w, \theta &\sim N\{\mathbf{X}_1\beta^w + \mathbf{X}_1(\mathbf{I} - \mathbf{W})\mathbf{X}_2\theta, \mathbf{C}_1\} \\ \beta^w \mid \theta &\sim N(\mathbf{W}\mathbf{X}_2\theta, \mathbf{C}_2) \\ \theta &\sim N(\mathbf{m}, \mathbf{C}_3), \end{aligned} \quad (4)$$

where  $\beta^w = (\beta_0^{w^T}, \dots, \beta_{p-1}^{w^T})^T$ , and  $\beta_k^w = \{\beta_k^w(\mathbf{s}_1), \dots, \beta_k^w(\mathbf{s}_n)\}^T$ .

Defining the PCP this way gives us tremendous flexibility in terms of parameterisation. If  $\mathbf{W}$  is the identity matrix we recover the CP and where  $\mathbf{W}$  is the zero matrix we have the NCP. If we let

$$\mathbf{W} = \begin{pmatrix} w_0\mathbf{I} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & w_1\mathbf{I} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & w_{p-1}\mathbf{I} \end{pmatrix}. \quad (5)$$

then we have the PCP investigated by Papaspiliopoulos *et al.* (2003, Section 4).

The question is how do we choose the entries of  $\mathbf{W}$  such that optimal performance of the Gibbs sampler is achieved? We answer this question by analysing the posterior correlation between global and random effects. Returning to the CP, if we apply the calculations given in Gelfand *et al.* (1995, Section 2) to our model set up it can be shown that

$$\text{cov}(\tilde{\beta}, \theta \mid \mathbf{C}_1, \mathbf{C}_2, \mathbf{C}_3, \mathbf{y}) = \mathbf{B}\mathbf{C}_2^{-1}\mathbf{X}_2\Sigma_{\theta|y} \quad (6)$$

where

$$\mathbf{B} = \text{var}(\tilde{\beta} \mid \theta, \mathbf{C}_1, \mathbf{C}_2, \mathbf{y}) = (\mathbf{X}_1^T\mathbf{C}_1^{-1}\mathbf{X}_1 + \mathbf{C}_2^{-1})^{-1},$$

and

$$\Sigma_{\theta|y} = \text{var}(\theta \mid \mathbf{C}_1, \mathbf{C}_2, \mathbf{C}_3, \mathbf{y}) = \left\{ (\mathbf{X}_1\mathbf{X}_2)^T \Sigma_Y^{-1} \mathbf{X}_1\mathbf{X}_2 + \mathbf{C}_3^{-1} \right\}^{-1}. \quad (7)$$

Therefore by substituting equation (3) into equation (6) we have that the posterior covariance of  $\beta^w$  and  $\theta$  is

$$\begin{aligned} \text{cov}(\beta^w, \theta \mid \mathbf{C}_1, \mathbf{C}_2, \mathbf{C}_3, \mathbf{y}) &= \text{cov}\{\tilde{\beta} - (\mathbf{I} - \mathbf{W})\mathbf{X}_2\theta, \theta \mid \mathbf{y}\} \\ &= \text{cov}(\tilde{\beta}, \theta \mid \mathbf{y}) - (\mathbf{I} - \mathbf{W})\mathbf{X}_2\text{var}(\theta \mid \mathbf{y}) \\ &= \mathbf{B}\mathbf{C}_2^{-1}\mathbf{X}_2\Sigma_{\theta|y} - (\mathbf{I} - \mathbf{W})\mathbf{X}_2\Sigma_{\theta|y} \\ &= \{\mathbf{B}\mathbf{C}_2^{-1} - (\mathbf{I} - \mathbf{W})\}\mathbf{X}_2\Sigma_{\theta|y}. \end{aligned} \quad (8)$$

We can see from (8) that  $\text{cov}(\beta^w, \theta \mid \mathbf{C}_1, \mathbf{C}_2, \mathbf{C}_3, \mathbf{y}) = \mathbf{0}$  when  $\mathbf{B}\mathbf{C}_2^{-1} = \mathbf{I} - \mathbf{W}$ . Therefore we define the optimal  $\mathbf{W}$  to be

$$\mathbf{W}^{opt} = \mathbf{I} - \mathbf{B}\mathbf{C}_2^{-1}. \quad (9)$$

For all that follows we drop the superscript from  $\mathbf{W}^{opt}$  and any time we refer to  $\mathbf{W}$  it will be the matrix defined in (9). By the Sherman-Woodbury-Morrison identity (Harville, 1997, Chapter 18) we can write  $\mathbf{W}$  as

$$\mathbf{W} = \mathbf{C}_2\mathbf{X}_1^T(\mathbf{C}_1 + \mathbf{X}_1\mathbf{C}_2\mathbf{X}_1^T)^{-1}\mathbf{X}_1, \quad (10)$$

which requires the inversion of matrix of order  $n$  and not of order  $np$ .

Equation (9) implies that to minimise the posterior correlation between the random effects and global effects we cannot, in general, restrict  $\mathbf{W}$  to be a diagonal matrix like that given in (5). It then follows that as  $\beta^w|\boldsymbol{\theta} \sim N(\mathbf{W}\mathbf{X}_2\boldsymbol{\theta}, \mathbf{C}_2)$  *a priori*, the prior mean of  $\beta_k^w(\mathbf{s}_i)$  will be a linear combination of all elements of  $\boldsymbol{\theta}$  and not just a proportion of  $\theta_k$ . For example, suppose we have  $p = 2$  processes,  $\beta_0^w$  and  $\beta_1^w$ , and we partition  $\mathbf{W}$  into four  $n \times n$  blocks. Then we have

$$\mathbf{W}\mathbf{X}_2\boldsymbol{\theta} = \begin{pmatrix} \mathbf{W}_{00} & \mathbf{W}_{01} \\ \mathbf{W}_{10} & \mathbf{W}_{11} \end{pmatrix} \begin{pmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix} = \begin{pmatrix} \mathbf{W}_{00}\mathbf{1}\theta_0 + \mathbf{W}_{01}\mathbf{1}\theta_1 \\ \mathbf{W}_{10}\mathbf{1}\theta_0 + \mathbf{W}_{11}\mathbf{1}\theta_1 \end{pmatrix},$$

and the  $i$ th row of  $\mathbf{W}_{00}\mathbf{1}$  is the weight assigned to  $\theta_0$  for  $\beta_0^w(\mathbf{s}_i)$ . More generally the  $i$ th row of  $\mathbf{W}_{kj}\mathbf{1}$  ( $k, j = 0, \dots, p-1$ ), is the weight assigned to  $\theta_j$  for  $\beta_k^w(\mathbf{s}_i)$ . This is illustrated in Section 5, Figure 11. Equation (9) also implies that when  $\mathbf{X}$  contains the values of spatially referenced covariates or when  $\mathbf{C}_2$  is the covariance of a spatial process, then the weights,  $\mathbf{W}\mathbf{X}_2$ , will vary over space, as demonstrated in Section 3.

### 2.3 Convergence rate of the PCP

We now look at the implication of setting  $\mathbf{W}$  according to (9) for the convergence rate of a Gibbs sampler using the PCP. For Gibbs samplers with Gaussian target distributions with known precision matrices we have analytical results for the exact convergence rate (Roberts and Sahu, 1997, Theorem 1). Convergence here is defined in terms of how rapidly the expectations of square integrable functions approach their stationary values.

Suppose that  $\boldsymbol{\xi} | \mathbf{y} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . We let  $\mathbf{Q} = \boldsymbol{\Sigma}^{-1}$  denote the posterior precision matrix. To compute the convergence rate first partition  $\mathbf{Q}$  according to a number of blocks, denoted by  $l$ , that are used for updating  $\boldsymbol{\xi}$ , i.e.,

$$\mathbf{Q} = \begin{pmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} & \cdots & \mathbf{Q}_{1l} \\ \mathbf{Q}_{21} & \mathbf{Q}_{22} & \cdots & \mathbf{Q}_{2l} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{Q}_{l1} & \mathbf{Q}_{l2} & \cdots & \mathbf{Q}_{ll} \end{pmatrix}. \quad (11)$$

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

**Theorem 2.1** (Roberts and Sahu, 1997) *A Markov chain with transition density*

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

*has a convergence rate equal to the maximum modulus eigenvalue of  $\mathbf{F}$ .*

**Corollary 2.2** *If we update  $\boldsymbol{\xi}$  in two blocks so that  $l = 2$  then*

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

*and the convergence rate is the maximum modulus eigenvalue of*

$$\mathbf{F}_{22} = \mathbf{Q}_{22}^{-1}\mathbf{Q}_{21}\mathbf{Q}_{11}^{-1}\mathbf{Q}_{12}.$$

To compute the convergence rate of the PCP we first need the posterior precision matrix of  $\beta^w$  and  $\theta$ , which we can identify by writing down  $\pi(\beta^w, \theta | y)$ , see Appendix A.1. The posterior precision matrix for the PCP is

$$Q^{pc} = \begin{pmatrix} Q_{\beta^w}^{pc} & Q_{\beta^w \theta}^{pc} \\ Q_{\theta \beta^w}^{pc} & Q_{\theta}^{pc} \end{pmatrix}, \quad (12)$$

where

$$\begin{aligned} Q_{\beta^w}^{pc} &= X_1^T C_1^{-1} X_1 + C_2^{-1}, \\ Q_{\beta^w \theta}^{pc} &= X_1^T C_1^{-1} X_1 (I - W) X_2 - C_2^{-1} W X_2, \\ Q_{\theta}^{pc} &= X_2^T (I - W)^T X_1^T C_1^{-1} X_1 (I - W) X_2 + X_2^T W^T C_2^{-1} W X_2 + C_3^{-1}. \end{aligned}$$

If we block update a Gibbs sampler according to the partitioning of the precision matrix (12), by Corollary 2.2, we have that the convergence rate of the PCP is the maximum modulus eigenvalue of the matrix

$$F_{22}^{pc} = (Q_{\theta}^{pc})^{-1} Q_{\theta \beta^w}^{pc} (Q_{\beta^w}^{pc})^{-1} Q_{\beta^w \theta}^{pc}.$$

By construction we have a  $2 \times 2$  block diagonal posterior covariance matrix for  $\beta^w$  and  $\theta$ . Therefore the precision matrix is also block diagonal and  $F_{22}^{pc}$  is null and immediate convergence is achieved. This is shown algebraically in Appendix A.2.

## 2.4 Block updating the PCP

Suppose now that we have constructed the PCP as before but we partition the partially centred random effects,  $\beta^w$ , into two disjoint sets:  $\beta_1^w$  and  $\beta_2^w$ , and update them separately in a Gibbs sampler. Partitioned accordingly, the covariance matrix of the joint posterior distribution of  $(\beta_1^w, \beta_2^w, \theta)$  is a  $3 \times 3$  block matrix given by

$$\Sigma = \begin{pmatrix} \Sigma_{\beta_1} & \Sigma_{\beta_{12}} & \mathbf{0} \\ \Sigma_{\beta_{21}} & \Sigma_{\beta_2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Sigma_{\theta} \end{pmatrix}.$$

Using results from Harville (1997, Chapter 8) we find the corresponding posterior precision matrix to be

$$Q = \begin{pmatrix} Q_{\beta_1} & Q_{\beta_{12}} & \mathbf{0} \\ Q_{\beta_{21}} & Q_{\beta_2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & Q_{\theta} \end{pmatrix}, \quad (13)$$

where

$$\begin{aligned} Q_{\beta_1} &= (\Sigma_{\beta_1} - \Sigma_{\beta_{12}} \Sigma_{\beta_2}^{-1} \Sigma_{\beta_{21}})^{-1}, \\ Q_{\beta_{12}} &= -(\Sigma_{\beta_1} - \Sigma_{\beta_{12}} \Sigma_{\beta_2}^{-1} \Sigma_{\beta_{21}})^{-1} \Sigma_{\beta_{12}} \Sigma_{\beta_2}^{-1}, \\ Q_{\beta_{21}} &= -(\Sigma_{\beta_2} - \Sigma_{\beta_{21}} \Sigma_{\beta_1}^{-1} \Sigma_{\beta_{12}})^{-1} \Sigma_{\beta_{21}} \Sigma_{\beta_1}^{-1}, \\ Q_{\beta_2} &= (\Sigma_{\beta_2} - \Sigma_{\beta_{21}} \Sigma_{\beta_1}^{-1} \Sigma_{\beta_{12}})^{-1}, \\ Q_{\theta} &= \Sigma_{\theta}^{-1}. \end{aligned}$$

Using the intermediate calculations given in Appendix A.3 it can be shown that the convergence rate corresponding to the precision matrix given by (13) is the maximum modulus eigenvalue of

$$F^{pc} = \begin{pmatrix} \mathbf{0} & -Q_{\beta_1}^{-1} Q_{\beta_{12}} & \mathbf{0} \\ \mathbf{0} & Q_{\beta_2}^{-1} Q_{\beta_{21}} Q_{\beta_1}^{-1} Q_{\beta_{12}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} = \begin{pmatrix} \mathbf{0} & \Sigma_{\beta_{12}} \Sigma_{\beta_2}^{-1} & \mathbf{0} \\ \mathbf{0} & \Sigma_{\beta_{21}} \Sigma_{\beta_1}^{-1} \Sigma_{\beta_{12}} \Sigma_{\beta_2}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix},$$

which will be zero if the posterior correlation between  $\beta_1^w$  and  $\beta_2^w$  is zero.

Alternatively, suppose that we update  $\beta^w$  as one block but partition  $\theta$  into  $\theta_1$  and  $\theta_2$ , updating them accordingly. The posterior covariance and precision matrices have the form

$$\Sigma = \begin{pmatrix} \Sigma_\beta & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \Sigma_{\theta_1} & \Sigma_{\theta_{12}} \\ \mathbf{0} & \Sigma_{\theta_{21}} & \Sigma_{\theta_2} \end{pmatrix}, \quad Q = \begin{pmatrix} Q_\beta & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & Q_{\theta_1} & Q_{\theta_{12}} \\ \mathbf{0} & Q_{\theta_{21}} & Q_{\theta_2} \end{pmatrix},$$

where

$$\begin{aligned} Q_\beta &= \Sigma_\beta^{-1}, \\ Q_{\theta_1} &= (\Sigma_{\theta_1} - \Sigma_{\theta_{12}} \Sigma_{\theta_2}^{-1} \Sigma_{\theta_{21}})^{-1}, \\ Q_{\theta_{12}} &= -(\Sigma_{\theta_1} - \Sigma_{\theta_{12}} \Sigma_{\theta_2}^{-1} \Sigma_{\theta_{21}})^{-1} \Sigma_{\theta_{12}} \Sigma_{\theta_2}^{-1}, \\ Q_{\theta_{21}} &= -(\Sigma_{\theta_2} - \Sigma_{\theta_{21}} \Sigma_{\theta_1}^{-1} \Sigma_{\theta_{12}})^{-1} \Sigma_{\theta_{21}} \Sigma_{\theta_1}^{-1}, \\ Q_{\theta_2} &= (\Sigma_{\theta_2} - \Sigma_{\theta_{21}} \Sigma_{\theta_1}^{-1} \Sigma_{\theta_{12}})^{-1}, \end{aligned}$$

and the convergence rate is the maximum modulus eigenvalue of

$$F = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & -Q_{\theta_1}^{-1} Q_{\theta_{12}} \\ \mathbf{0} & \mathbf{0} & Q_{\theta_2}^{-1} Q_{\theta_{21}} \end{pmatrix} = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Sigma_{\theta_{12}} \Sigma_{\theta_2}^{-1} \\ \mathbf{0} & \mathbf{0} & \Sigma_{\theta_{21}} \Sigma_{\theta_1}^{-1} \end{pmatrix},$$

which will be a null matrix if the two blocks of  $\theta$  are uncorrelated *a posteriori*.

It is the relationship between convergence rate and inter-block correlation that we take advantage of when constructing the PCP. For our construction, immediate convergence is only guaranteed if the random effects and global effects are each updated as one complete block. If a greater number of blocks are used we cannot, in general, find a matrix  $W$  that will remove all cross covariances and return a convergence rate of zero. To see this first note that the posterior covariance matrix  $\Sigma_{\theta|y}$ , given in (7), is unaffected by hierarchical centering. Therefore partial centering cannot remove any posterior correlation between subsets of  $\theta$ , and so all of its elements must be updated together. Then suppose that we partition the partially centred random effects into  $l$  blocks so that  $\beta^w = (\beta_1^{wT}, \dots, \beta_l^{wT})^T$ . We find the posterior covariance between the  $ij$ th block to be

$$\begin{aligned} \text{cov}(\beta_i^w, \beta_j^w | y) &= B_{ij} + B_i C_2^{-1} X_2 \Sigma_{\theta|y} X_2^T C_2^{-1} B_{.j} - B_i C_2^{-1} X_2 \Sigma_{\theta|y} X_2^T (I - W)_{.j}^T \\ &\quad - (I - W)_{i.} X_2 \Sigma_{\theta|y} X_2^T C_2^{-1} B_{.j} + (I - W)_{i.} X_2 \Sigma_{\theta|y} X_2^T (I - W)_{.j}^T \\ &= B_{ij} + B_i C_2^{-1} X_2 \Sigma_{\theta|y} X_2^T \{C_2^{-1} B_{.j} - (I - W)_{.j}^T\} \\ &\quad + (I - W)_{i.} X_2 \Sigma_{\theta|y} X_2^T \{(I - W)_{.j}^T - C_2^{-1} B_{.j}\}, \end{aligned}$$

where  $B_{ij}$  is the  $ij$ th block of  $B = \text{var}(\tilde{\beta} | \theta, y)$ . We let  $B_i$  denote the rows of  $B$  associated with the  $i$ th block and let  $B_{.j}$  denote the columns of  $B$  associated with the  $j$ th block, with  $(I - W)_{i.}$  and  $(I - W)_{.j}$  having similar interpretations. We see that if  $(I - W)_{.j}^T = C_2^{-1} B_{.j}$  then  $\text{cov}(\beta_i^w, \beta_j^w | y) = B_{ij}$ , which is generally a non-zero matrix. Therefore we must update  $\beta^w$  as one component and  $\theta$  as another.

## 2.5 PCP for unknown variance parameters

The PCP relies on the  $W$  matrix which, by construction, removes the posterior correlation between  $\beta^w$  and  $\theta$ . However, the derivation of  $W$  is conditional on the covariance matrices,



$\mathbf{C}_1$ ,  $\mathbf{C}_2$  and  $\mathbf{C}_3$ . Therefore when the variance parameters are unknown how do we compute  $\mathbf{W}$ ? We propose a dynamically updated parameterisation that uses the most recent values to re-compute  $\mathbf{W}$  at each move of the Markov chain along a coordinate of which it is a function.

Let  $\boldsymbol{\xi}^{(t)} = \{\boldsymbol{\beta}^{w(t)}, \boldsymbol{\theta}^{(t)}, \boldsymbol{\sigma}^{2(t)}, \sigma_\epsilon^{2(t)}\}^\top$  be the current state of the Markov chain, where  $\boldsymbol{\sigma}^{2(t)} = \{\sigma_0^{2(t)}, \dots, \sigma_{p-1}^{2(t)}\}^\top$ . Also let  $\boldsymbol{\sigma}^{2(t+1)}_{-k} = \{\sigma_0^{2(t+1)}, \dots, \sigma_{k-1}^{2(t+1)}, \sigma_{k+1}^{2(t+1)}, \dots, \sigma_{p-1}^{2(t+1)}\}^\top$  be the partially updated vector without  $\sigma_k^2$ . We write  $\mathbf{W} = \mathbf{W}(\boldsymbol{\sigma}^2, \sigma_\epsilon^2)$  to highlight the dependency of  $\mathbf{W}$  on the variance parameters. Recall that in Section 2.4 we show that  $\boldsymbol{\beta}^w$  and  $\boldsymbol{\theta}$  should each be updated as one block. We obtain a new sample,  $\boldsymbol{\xi}^{(t+1)} \sim \pi(\boldsymbol{\xi} | \mathbf{y})$  as follows:

1. Sample  $\boldsymbol{\beta}^{w(t+1)} \sim \pi\{\boldsymbol{\beta}^w | \boldsymbol{\theta}^{(t)}, \boldsymbol{\sigma}^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{W}(\boldsymbol{\sigma}^{2(t)}, \sigma_\epsilon^{2(t)}), \mathbf{y}\}$ .
2. Sample  $\boldsymbol{\theta}^{(t+1)} \sim \pi\{\boldsymbol{\theta} | \boldsymbol{\beta}^{w(t+1)}, \boldsymbol{\sigma}^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{W}(\boldsymbol{\sigma}^{2(t)}, \sigma_\epsilon^{2(t)}), \mathbf{y}\}$ .
3. For  $k = 0, \dots, p-1$ ,  
Sample  $\sigma_k^{2(t+1)} \sim \pi\{\sigma_k^2 | \boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \boldsymbol{\sigma}^{2(t+1)}_{-k}, \sigma_\epsilon^{2(t)}, \mathbf{W}(\boldsymbol{\sigma}^{2(t+1)}_{-k}, \sigma_k^{2(t)}, \sigma_\epsilon^{2(t)}), \mathbf{y}\}$ .
4. Sample  $\sigma_\epsilon^{2(t+1)} \sim \pi\{\sigma_\epsilon^2 | \boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \boldsymbol{\sigma}^{2(t+1)}, \mathbf{W}(\boldsymbol{\sigma}^{2(t+1)}, \sigma_\epsilon^{2(t)}), \mathbf{y}\}$ .

The distributions of  $\sigma_k^2$  and  $\sigma_\epsilon^2$  are conditioned on their respective current values through  $\mathbf{W}$ , i.e.  $\sigma_0^{2(t+1)}$  is conditioned on  $\sigma_0^{2(t)}$ . There is no stochastic relationship between the parameters and  $\mathbf{W}$ , given the parameters  $\mathbf{W}$  is completely determined. However, we must ensure that by dynamically updating  $\mathbf{W}$  we do not disturb the stationary distribution of the Markov chains generated from the Gibbs sampler. We need to show that

$$\pi\{\boldsymbol{\xi}^{(t+1)} | \mathbf{y}\} = \int P\{\boldsymbol{\xi}^{(t+1)} | \boldsymbol{\xi}^{(t)}\} \pi\{\boldsymbol{\xi}^{(t)} | \mathbf{y}\} d\boldsymbol{\xi}^{(t)}, \quad (14)$$

where  $P\{\cdot | \cdot\}$  is the transition kernel of the chain. The definition of  $P\{\cdot | \cdot\}$  and the proof of that (14) holds is provided in Appendix A.4.

### 3 Spatially varying weights of partial centering

#### 3.1 A spatially varying intercept model

In this section we illustrate how the weights of partial centering depend upon the variance parameters but also the correlation structure of the latent processes. In particular, we will see how the weights vary across the spatial region and the impact of a spatially varying covariate. To focus on these relationships we will consider simplified versions of model (4) that have one global parameter and one latent process. We do not need to simulate data as given a set of sampling locations we can investigate the weights through the variance and decay parameters.

We begin by looking at the following model,

$$\begin{aligned} \mathbf{Y} | \boldsymbol{\beta}_0^w, \theta_0 &\sim N\{\boldsymbol{\beta}_0^w + (\mathbf{I} - \mathbf{W})\mathbf{1}\theta_0, \sigma_\epsilon^2 \mathbf{I}\} \\ \boldsymbol{\beta}_0^w | \theta_0 &\sim N(\mathbf{W}\mathbf{1}\theta_0, \sigma_0^2 \mathbf{R}_0) \\ \theta_0 &\sim N(m_0, \sigma_0^2 v_0). \end{aligned} \quad (15)$$

which has one global parameter  $\boldsymbol{\theta} = \theta_0$  and one latent spatial process  $\boldsymbol{\beta}^w = \boldsymbol{\beta}_0^w$ , and hence can be found from (4) by letting  $\mathbf{X}_1 = \mathbf{I}$ ,  $\mathbf{C}_1 = \sigma_\epsilon^2 \mathbf{I}$ ,  $\mathbf{X}_2 = \mathbf{1}$ ,  $\mathbf{C}_2 = \sigma_0^2 \mathbf{R}_0$ ,  $\mathbf{m} = m_0$  and  $\mathbf{C}_3 = \sigma_0^2 v_0$ . Therefore, using the representation of  $\mathbf{W}$  given in equation (10) we have

$$\mathbf{W} = \sigma_0^2 \mathbf{R}_0 (\sigma_\epsilon^2 \mathbf{I} + \sigma_0^2 \mathbf{R}_0)^{-1}, \quad (16)$$

and so the entries of  $\mathbf{W}$  depend on variance parameters  $\sigma_0^2$ ,  $\sigma_\epsilon^2$  and, through correlation matrix  $\mathbf{R}$ , the spatial decay parameter and the set of sampling locations,  $\mathbf{s}_1, \dots, \mathbf{s}_n$ .

Here we select sampling locations according to a pattern, such that the locations are more densely clustered in some regions of the domain than others. We consider 200 locations in the unit square, see Figure 1, which we split into nine sub-squares of equal area. We randomly select 100 points in the top left square and 25 points in the three areas to which it is adjacent. The remaining five sub-squares have five points randomly chosen within them.

We consider five variance ratios:  $\delta_0 = \sigma_0^2/\sigma_\epsilon^2 = 0.01, 0.1, 1, 10, 100$ , and three effective ranges:  $d_0 = \sqrt{2}/3, 2\sqrt{2}/3, \sqrt{2}$ , which are chosen with respect to the maximum separation of two points in the unit square,  $\sqrt{2}$ . An effective range of zero, implying independent random effects, returns weights that are the same at each location. For each of the 15 variance ratio-effective range combinations we compute  $\mathbf{W}\mathbf{X}_2$ , whose  $i$ th value is the weight assigned to  $\theta_0$  at each  $\mathbf{s}_i$   $i = 1, \dots, 200$ .

We use the `Tps` function in the R package `fields` (Furrer *et al.*, 2009) to interpolate the weights over the unit square. These interpolated plots of spatially varying weights are given in Figure 2. Each row corresponds to a value of  $\delta_0$ , from 0.01 in top row to 100 in the bottom. For each row going left to right we have increasing effective ranges,  $d_0 = \sqrt{2}/3, 2\sqrt{2}/3, \sqrt{2}$ . We can see that as the variance ratio increases the weights are higher, as they are when the effective range increases. Within each panel, the areas of higher weights are concentrated around the areas of more densely positioned sampling locations. The stronger the correlation, the farther reaching is the influence of these clusters.

### 3.2 A spatially varying slope model

We now investigate the effect of a covariate upon the spatially varying weights. To do this we look at the following model

$$\begin{aligned} \mathbf{Y} \mid \beta_1^w, \theta_1 &\sim N\{\mathbf{D}\beta_1^w + (\mathbf{I} - \mathbf{W})\mathbf{1}\theta_1, \sigma_\epsilon^2\mathbf{I}\} \\ \beta_1^w \mid \theta_1 &\sim N(\mathbf{W}\mathbf{1}\theta_1, \sigma_1^2\mathbf{R}_1) \\ \theta_1 &\sim N(m_1, \sigma_1^2v_1), \end{aligned} \tag{17}$$

where  $\mathbf{D} = \text{diag}(\mathbf{x})$  and  $\mathbf{x} = \{x(\mathbf{s}_1), \dots, x(\mathbf{s}_n)\}^T$  contains the values of a known spatially referenced covariate. We have a global slope, hence  $\boldsymbol{\theta} = \theta_1$ , and a partially centered spatial process  $\beta^w = \beta_1^w$ . Model (17) can be retrieved from model (4) by letting  $\mathbf{X}_1 = \mathbf{D}$ ,  $\mathbf{C}_1 = \sigma_\epsilon^2\mathbf{I}$ ,  $\mathbf{X}_2 = \mathbf{1}$ ,  $\mathbf{C}_2 = \sigma_1^2\mathbf{R}_1$ ,  $\mathbf{m} = m_1$  and  $\mathbf{C}_3 = \sigma_1^2v_1$ .

For model (17) the  $\mathbf{W}$  matrix is

$$\mathbf{W} = \sigma_1^2\mathbf{R}_1\mathbf{D}(\sigma_\epsilon^2\mathbf{I} + \sigma_1^2\mathbf{D}\mathbf{R}_1\mathbf{D})^{-1}\mathbf{D}.$$

To see how these weights vary across the domain we randomly select 200 points uniformly over the unit square, see Figure 3.

We generate the values of  $\mathbf{x}$  by selecting a point  $\mathbf{s}_x$ , which we may imagine to be the site of a source of pollution. We assume that the value for the observed covariate at site  $\mathbf{s}_i$  decays exponentially at rate  $\phi_x$  with increasing separation from  $\mathbf{s}_x$ , so that

$$x(\mathbf{s}_i) = \exp(-\phi_x\|\mathbf{s}_i - \mathbf{s}_x\|) \quad (i = 1, \dots, n).$$

The spatial decay parameter  $\phi_x$  is chosen such that there is an effective spatial range of  $\sqrt{2}/2$ , i.e. if  $\|\mathbf{s}_i - \mathbf{s}_x\| = \sqrt{2}/2$  then  $x(\mathbf{s}_i) = 0.05$ . The values of  $\mathbf{x}$  are standardised by subtracting their sample mean and dividing by their sample standard deviation. Figure 4 gives the interpolated

covariate surface where  $\mathbf{s}_x = (0.936, 0.117)^\top$ . We can see how the values decay with increased separation from  $\mathbf{s}_x$ .

The optimal weights are computed for 15 combinations of variance ratio  $\delta_1$  and effective range  $d_1$ , where  $\delta_1 = \sigma_1^2/\sigma_\epsilon^2 = 0.01, 0.1, 1, 10, 100$  and  $d_1 = \sqrt{2}/3, 2\sqrt{2}/3, \sqrt{2}$ . The weights are interpolated and plotted in Figure 5. The layout is the same as in Figure 2, with each row corresponding to a value of  $\delta_1$ , going from 0.01 at the top to 100 at the bottom, and increasing effective range from left to right. We see that the weights increase with increasing  $\delta_1$  or  $d_1$ . It is also clear that for locations near  $\mathbf{s}_x$ , where the values of the covariate are greatest, the weights are higher.

## 4 Simulation studies

### 4.1 Simulation example 1

In this section we use simulated data to investigate the performance of the Gibbs sampler associated with the PCP. All of the relevant posterior distributions can be found in Appendix A.5. We simulate data from model (15) for  $n = 40$  randomly chosen locations across the unit square. We let hyperparameters  $m_0 = 0$  and  $v_0 = 10^4$ .

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  $\sigma_\epsilon^2$  accordingly. For each of the five levels of  $\delta_0$  we have four values of the decay parameter  $\phi_0$ , chosen such that there is an effective range,  $d_0$ , 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  $\phi_0$  in all. Each of these combinations is used to simulate 20 datasets, and so there are 400 data sets in total.

To begin the variance parameters,  $\sigma_0^2$  and  $\sigma_\epsilon^2$ , and the decay parameter  $\phi_0 = -\log(0.05)/d_0$  are held fixed at their true values and so for each iteration of the Gibbs sampler we generate samples from the full conditional distributions of  $\beta_0^w$  and  $\theta_0$ .

The efficiency of the sampler is judged by two measures. The first statistic we use is based on the potential scale reduction factor (PSRF) (Gelman and Rubin, 1992). We define the  $\text{PSRF}_M(1.1)$  to be the number of iterations required for the PSRF to fall below 1.1. To compute the  $\text{PSRF}_M(1.1)$  we simulate multiple chains from widely dispersed starting values. In particular, we take values that are outside of the intervals described by pilot chains. At every fifth iteration the PSRF 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  $\theta_0$ , (Robert and Casella, 2004, Chapter 12).

For each data set we generate five chains of length 25,000 and compute the  $\text{PSRF}_M(1.1)$  and the ESS of  $\theta_0$  out of a total of 125,000 samples. The results are plotted in Figure 6. On the top row we have the  $\text{PSRF}_M(1.1)$  and on the bottom the ESS. The five panels in each row correspond to a value of  $\delta_0$ , with 0.01 on the left rising to 100 on the right. Within each panel we have four boxplots, one for each level of the effective range, again rising from left to right. Each boxplot consists of 20 values, for the 20 repetitions of that variance ratio-effective range combination.

Figure 6 shows that for the PCP with known variance parameters we achieve near immediate convergence and independent samples for  $\theta_0$  in all cases, and that it is robust to changes in both variance ratio and strength of correlation.

### 4.2 Simulation example 2

The second simulation study drops the assumption of known variance parameters, fixing only the decay parameter. In this case we judge performance by the  $\text{MPSRF}_M(1.1)$  which we define

to be the number of iteration needed for the multivariate PSRF (Brooks and Gelman, 1998) to fall below 1.1. In addition we record the ESS of  $\theta_0$ ,  $\sigma_0^2$  and  $\sigma_\epsilon^2$ .

Recall that 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  $\sigma_0^2$  and  $\sigma_\epsilon^2$ . These are common hyperparameters for inverse gamma prior distributions, see Sahu *et al.* (2010, 2007); Gelfand *et al.* (2003).

Figure 7 shows the MPSRF<sub>M</sub>(1.1) on the top row and the ESS of  $\theta_0$  on the bottom row for the 20 combinations of  $\delta_0$  and  $d_0$ . There is more variability in the results seen here for the MPSRF<sub>M</sub>(1.1) than we saw for the PSRF<sub>M</sub>(1.1) when the variance parameters were fixed, Figure 6. When the random effects are independent, weak identifiability of the variance parameters can effect the performance of the sampler as marginally  $\text{var}\{Y(\mathbf{s}_i)\} = \sigma_\epsilon^2 + \sigma_0^2$ . However, the robustness to changes in  $\delta_0$  remains and we still see rapid convergence in most cases. The ESS for  $\theta_0$  remains high, with a median value above 120,000 for all of the 20 combinations of  $\delta_0$  and  $d_0$ .

Boxplots of the ESS of the variance parameters are given in Figure 8, with the results for  $\sigma_0^2$  on the top row and  $\sigma_\epsilon^2$  on the bottom row. There is a suggestion that increasing  $\delta_0$  increases the ESS of  $\sigma_0^2$  and decreases the ESS of  $\sigma_\epsilon^2$ . For a fixed value of  $\delta_0$  we can see that the ESS of both variance parameters increases as the effective range increases. The stronger correlation across the random effects means that the variability seen in the data can be more easily separated between the two components.

We compare the results for the PCP with those obtained for the CP and the NCP by calculating the mean responses of each measure of performance for each of the 20 variance ratio-effective range combinations. Table 1 shows the mean MPSRF<sub>M</sub>(1.1) and mean ESS for  $\theta_0$ . The sampling efficiency of the CP and the NCP is dependent on the covariance parameters. The CP performs best for higher values of  $\delta_0$  and longer effective ranges, and the NCP performs best for lower values of  $\delta_0$  and shorter effective ranges. In contrast, the PCP is robust to changes in  $\delta_0$  and  $d_0$ . Moreover, we see that the PCP has a lower average MPSRF<sub>M</sub>(1.1) for most cases, and when it does not, the difference is less than 3%. It is clear that, in terms of the ESS of  $\theta_0$ , the PCP is superior to the CP and the NCP in all cases.

A similar comparison for the mean ESS of the variance parameters is given in Table 2. The ESS of the variance parameters is not strongly effected by the fitting methods we consider here as a reparameterisation only acts upon the mean structure of the model. However, in cases where a particular parameterisation is very inefficient, poor mixing can be observed for the variance parameters. This is seen here for the ESS of  $\sigma_0^2$  for the NCP when  $\delta_0 \geq 1$ .

## 5 Californian ozone concentration data

In this section we compare the efficiency of the CP, the NCP and the PCP when fitting model (1) to a real data set. We have ozone concentration data from the State of California. It is a spatial data set with values, in parts per billion (ppb), of the annual fourth highest daily maximum eight-hour average for the year 2008, and has been previously analysed by Gelfand *et al.* (2012). The eight-hour average for the current hour is the mean concentration of the last four hours, the current hour and the future three hours. The annual fourth highest eight-hour average is the key measure used by the U.S. Environmental Protection Agency for monitoring ozone concentrations.

Data are collected at 176 irregularly spaced locations across California. We fit the model using data from 132 sites, leaving out 44 sites for validation, see Figure 9. The mean and standard deviation for the 132 data sites is 80.35 ppb and 17.72 ppb respectively.

The spatially varying covariate we use is land use. Sites are categorised as urban or suburban

and assigned the value one, or they are categorised as rural, and assigned the value zero. Of the 132 data sites, 89 are urban or suburban with mean concentration 78.71 ppb and standard deviation 18.53 ppb. The remaining 43 rural sites have mean 83.74 ppb and standard deviation 15.59 ppb. Of the 44 validation sites 28 are urban or suburban and 16 are rural.

As we have a spatially varying coefficient we fit the different parameterisations of model (1) with  $p = 2$ . Therefore we have two spatial processes, an intercept and a slope process and so  $\tilde{\boldsymbol{\beta}} = (\tilde{\boldsymbol{\beta}}_0^T, \tilde{\boldsymbol{\beta}}_1^T)^T$  for the CP,  $\boldsymbol{\beta} = (\boldsymbol{\beta}_0^T, \boldsymbol{\beta}_1^T)^T$  for the NCP and  $\boldsymbol{\beta}^w = (\boldsymbol{\beta}_0^{wT}, \boldsymbol{\beta}_1^{wT})^T$  for the PCP. Each process has a corresponding global parameter and a variance parameter and so  $\boldsymbol{\theta} = (\theta_0, \theta_1)^T$  and  $\boldsymbol{\sigma}^2 = (\sigma_0^2, \sigma_1^2)^T$ . We use an exponential correlation function for both processes and so  $\boldsymbol{\phi} = (\phi_0, \phi_1)^T$ . In addition we have the data variance,  $\sigma_\epsilon^2$ , and so we have  $2(n + 3) + 1$  parameters to estimate for each parameterisation.

For the prior distribution of  $\boldsymbol{\theta}$  we let  $\mathbf{m} = (0, 0)^T$  and  $v_0 = v_1 = 10^4$ . We 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. To stabilise the variance and avoid negative predictions, we model the data on the square root scale, as done by Sahu *et al.* (2007) and Berrocal *et al.* (2010) when modelling ozone concentrations for the U.S.

As mentioned in Section 2.1 we estimate the spatial decay parameters by performing a grid search over a range of values for  $\phi_0$  and  $\phi_1$ . The estimates are taken to be the pair of values that minimise the prediction error with respect to the validation data. The criteria used to compute the prediction error are the mean absolute prediction error (MAPE), the root mean square prediction error (RMSPE) and the continuous ranked probability score (CRPS), see Gneiting *et al.* (2007) for example.

The greatest distance between any two of the 132 monitoring stations in California is 1190 kilometers (km) and so we select values of  $\phi_0$  and  $\phi_1$  corresponding to effective ranges of 50, 100, 250, 500 and 1000 km. 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.

We denote by  $d_0$  and  $d_1$  the effective range implied by  $\phi_0$  and  $\phi_1$  respectively. The values of the MAPE, RMSPE and CRPS for the 25 combinations of  $d_0$  and  $d_1$  are given in Table 3. We see that the prediction error is minimised for two of the three criteria when  $d_0 = 250$  and  $d_1 = 500$  and so our estimates for the spatial decay parameters are

$$\hat{\phi}_0 = -\log(0.05)/250 \approx 0.012, \quad \text{and} \quad \hat{\phi}_1 = -\log(0.05)/500 \approx 0.006.$$

For each parameterisation we generate five Markov chains of length 25,000 from the same set widely dispersed starting values. The MPSRF<sub>M</sub>(1.1) and the ESS for  $\boldsymbol{\theta} = (\theta_0, \theta_1)^T$ ,  $\boldsymbol{\sigma}^2 = (\sigma_0^2, \sigma_1^2)^T$  and  $\sigma_\epsilon^2$  are computed and given in Table 4.

We see that the CP requires far fewer iterations for the MPSRF to drop below 1.1 than the NCP, 135 versus 1405. There is a significant difference in the ESS of the mean parameters between the two parameterisations. The CP yields more than 80 times the number of effect samples for  $\theta_0$ , and more than 16 times the number of effective samples for  $\theta_1$  than the NCP. The ESS for the variance parameters is higher for the CP than the NCP in particular for  $\sigma_0^2$ . There is little difference between the CP and the PCP in terms of the MPSRF<sub>M</sub>(1.1) and the ESS of the variance parameters. The main difference lies in the ESS of the global mean parameters. The PCP returns independent samples from the marginal posterior distribution of  $\theta_0$  and near independent samples for  $\theta_1$ .

The run times for the CP and the NCP are almost the same, but updating  $\mathbf{W}$  within the sampler means that the PCP is more computationally demanding. Table 5 gives the same measures as Table 4 but adjusted for computation time. We let

$$\text{MPSRF}_t(1.1) = \text{MPSRF}_M(1.1) \times \text{time per iteration},$$

denote the computation time (in seconds) for the MPSRF to fall below 1.1, and let ESS/s denote the ESS per second.

The shorter run times of the CP give it the advantage over the PCP in terms of  $\text{MPSRF}_t(1.1)$  and ESS/s. However, in order to retain the same number of effective samples the CP will need a longer chain than the PCP. This means that more data must be stored and handled.

To obtain parameter estimates we use the PCP to run a single long chain of 50,000 iterations and discard the first 10,000. Parameter estimates and their 95% credible intervals are given in Table 6. We also include estimates and 95% credible intervals for the variance ratios  $\delta_0 = \sigma_0^2/\sigma_\epsilon^2$  and  $\delta_1 = \sigma_1^2/\sigma_\epsilon^2$ .

A negative estimate for  $\theta_1$  implies that ozone concentrations are higher in rural areas, although we see here that given the spatially correlated random effects,  $\theta_1$  is not significantly different from zero. The variances of the spatial processes are estimated to be greater than that of the pure error process as the Gaussian processes capture the spatial variation in the data. The estimates of the variance ratios  $\delta_0$  and  $\delta_1$  are approximately five and three respectively. Given these results it is not surprising that the CP outperformed the NCP here. The density plots for the model parameters are given in Figure 10.

Figure 11 shows the spatially varying parameterisation that results from the PCP. Contained are interpolated maps of the average weight of partial centering given by each process to each global parameter for all 176 sampling locations. The top row displays the weights for  $\theta_0$  on the left and  $\theta_1$  on the right, for the intercept process,  $\beta_0^w$ . The bottom row has the same but for the slope process,  $\beta_1^w$ . We can see that for all plots higher weights are given to areas where there are clusters of sampling locations with lower weights to the north and east of the State. Looking first at the off diagonal panes, we can see that both processes give a low weighting their opposing global parameter. For the diagonal plots we see that the intercept process gives a high weight to  $\theta_0$  everywhere, reflecting the relatively large value of  $\sigma_0^2$  and explaining why the ESS of  $\theta_0$  for the CP is nearly as high as that obtained by the PCP. The weights given by the slope process to  $\theta_1$  show the most spatial variation and demonstrate clearly why the neither the CP or NCP could match the sampling efficiency of the PCP  $\theta_1$  and highlight the effectiveness of allowing the parameters to vary spatially.

## 6 Discussion

We have investigated the performance of a PCP for the spatially varying coefficients model. We are able to parameterise the model in such a way that we remove the posterior covariance between the random and global effects and produce a Gibbs sampler which converges immediately. The construction is conditioned on the covariance matrices in the model. We have shown that the parameterisation can be updated dynamically within the Gibbs sampler for the case when these matrices are known only up to a set of covariance parameters which must be estimated.

The optimal weights of partial centering are shown to vary over the spatial domain, with higher weights given to locations where the data is more informative about the latent surface. Therefore, higher weights are found when the data precision is relatively high, or there is some clustering of locations. We also saw higher weights for locations where the value of the covariate was higher.

We make it clear that although the interpolated plots given in Sections 3 and 5 are informative they do not represent a true surface in the sense that the interpolated values are not estimates of a true value of the weights at an unsampled location. Indeed, if we were to obtain a further measurement at a new location then the values of the weights at the existing locations would change.

Our investigations show that the unlike the CP and the NCP, the performance of the PCP is robust to changes in the variance and correlation parameters. Swift convergence and independent, or near independent samples from the posterior distributions of the mean parameters are achieved for all of the data sets we considered, whether it was simulated or real data.

The PCP requires us to update all of the random effects in one block and all of the global effects in another. It is a computationally intensive strategy which we recommend for modest sized spatial datasets. However, for larger datasets many practitioners turn to Gaussian predictive process (GPP) models (Banerjee *et al.*, 2008), as implemented in R software packages spBayes (Finley *et al.*, 2015) and spTimer (Bakar and Sahu, 2015). As the method proposed in this paper can be applied to any model that can be written as a three stage NLHM, we believe that the PCP could be used in conjunction with GPP models thus broadening its applicability.

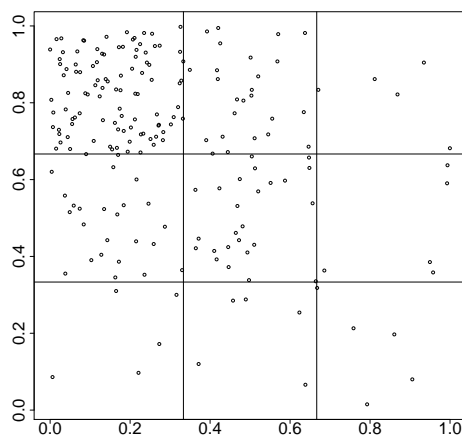


Figure 1: Patterned sampling locations. 100 top left; 25 in top middle, middle left and middle middle; five top right, middle right and bottom third.

Table 1: Means of the  $\text{MPSRF}_M(1.1)$  and the ESS of  $\theta_0$  for 20 variance ratio-effective range combinations for the CP, the NCP and the PCP

$\delta_0$	$d_0/\sqrt{2}$	$\text{MPSRF}_M(1.1)$			ESS of $\theta_0$		
		CP	NCP	PCP	CP	NCP	PCP
0.01	0	3064.50	172.00	<b>163.25</b>	463	108819	<b>124988</b>
	1/3	1115.75	278.25	<b>251.25</b>	1821	103342	<b>116659</b>
	2/3	544.50	166.25	<b>154.00</b>	3055	105397	<b>125108</b>
	1	366.50	184.75	<b>175.00</b>	4652	94657	<b>123730</b>
0.1	0	3528.25	3455.50	<b>2464.50</b>	4707	20420	<b>125272</b>
	1/3	<b>607.00</b>	1305.50	624.25	13336	33347	<b>108922</b>
	2/3	271.25	592.75	<b>251.50</b>	25884	28959	<b>109987</b>
	1	211.00	518.25	<b>203.75</b>	31938	22361	<b>113191</b>
1	0	274.50	910.50	<b>187.50</b>	25523	7353	<b>124945</b>
	1/3	134.50	1639.00	<b>118.75</b>	65927	3785	<b>120325</b>
	2/3	78.25	2092.00	<b>76.00</b>	82100	3148	<b>121013</b>
	1	<b>101.00</b>	2473.75	103.00	84742	2722	<b>115700</b>
10	0	123.50	1140.25	<b>105.75</b>	32578	5226	<b>125091</b>
	1/3	<b>79.75</b>	1556.25	<b>79.75</b>	83586	2918	<b>123055</b>
	2/3	45.25	2824.25	<b>44.50</b>	102306	1734	<b>124341</b>
	1	<b>49.50</b>	3542.25	50.50	107261	1177	<b>122774</b>
100	0	104.75	1124.75	<b>102.75</b>	32891	4596	<b>125388</b>
	1/3	42.75	1607.50	<b>42.50</b>	84755	2772	<b>124120</b>
	2/3	41.75	2544.75	<b>41.50</b>	104941	1671	<b>124214</b>
	1	<b>32.75</b>	3427.75	33.00	108050	1186	<b>124670</b>

Table 2: Means of the ESS of  $\sigma_0^2$  and ESS of  $\sigma_\epsilon^2$  for 20 variance ratio-effective range combinations for the CP, the NCP and the PCP

$\delta_0$	$d_0/\sqrt{2}$	ESS of $\sigma_0^2$			ESS of $\sigma_\epsilon^2$		
		CP	NCP	PCP	CP	NCP	PCP
0.01	0	4138	4078	<b>4244</b>	27753	95416	<b>98456</b>
	1/3	<b>4312</b>	4268	4295	56647	<b>95916</b>	94741
	2/3	5061	<b>5071</b>	5044	81338	<b>111550</b>	110636
	1	<b>4804</b>	4737	4758	88184	<b>108623</b>	107762
0.1	0	383	391	<b>397</b>	532	505	<b>545</b>
	1/3	3696	3027	<b>3719</b>	36659	30190	<b>39819</b>
	2/3	5368	4873	<b>5406</b>	60338	58930	<b>62573</b>
	1	5784	5095	<b>5790</b>	72669	70080	<b>73138</b>
1	0	5527	3868	<b>5557</b>	6188	6208	<b>6321</b>
	1/3	<b>8333</b>	3557	8297	12400	11488	<b>12412</b>
	2/3	<b>9430</b>	3945	<b>9430</b>	<b>22846</b>	18310	22828
	1	<b>9714</b>	3712	9711	<b>28853</b>	23222	28770
10	0	11270	3945	<b>11352</b>	12216	12308	<b>12495</b>
	1/3	<b>18435</b>	3543	18413	<b>18226</b>	18009	18167
	2/3	<b>23692</b>	3699	23632	<b>26945</b>	26373	26825
	1	<b>24595</b>	4015	24572	<b>29228</b>	28790	29177
100	0	11423	3652	<b>11457</b>	12427	12495	<b>12690</b>
	1/3	<b>19708</b>	3559	19668	<b>20118</b>	19930	20085
	2/3	<b>27706</b>	3929	27649	<b>26197</b>	25924	26120
	1	<b>26937</b>	3012	26936	<b>37121</b>	36499	37087



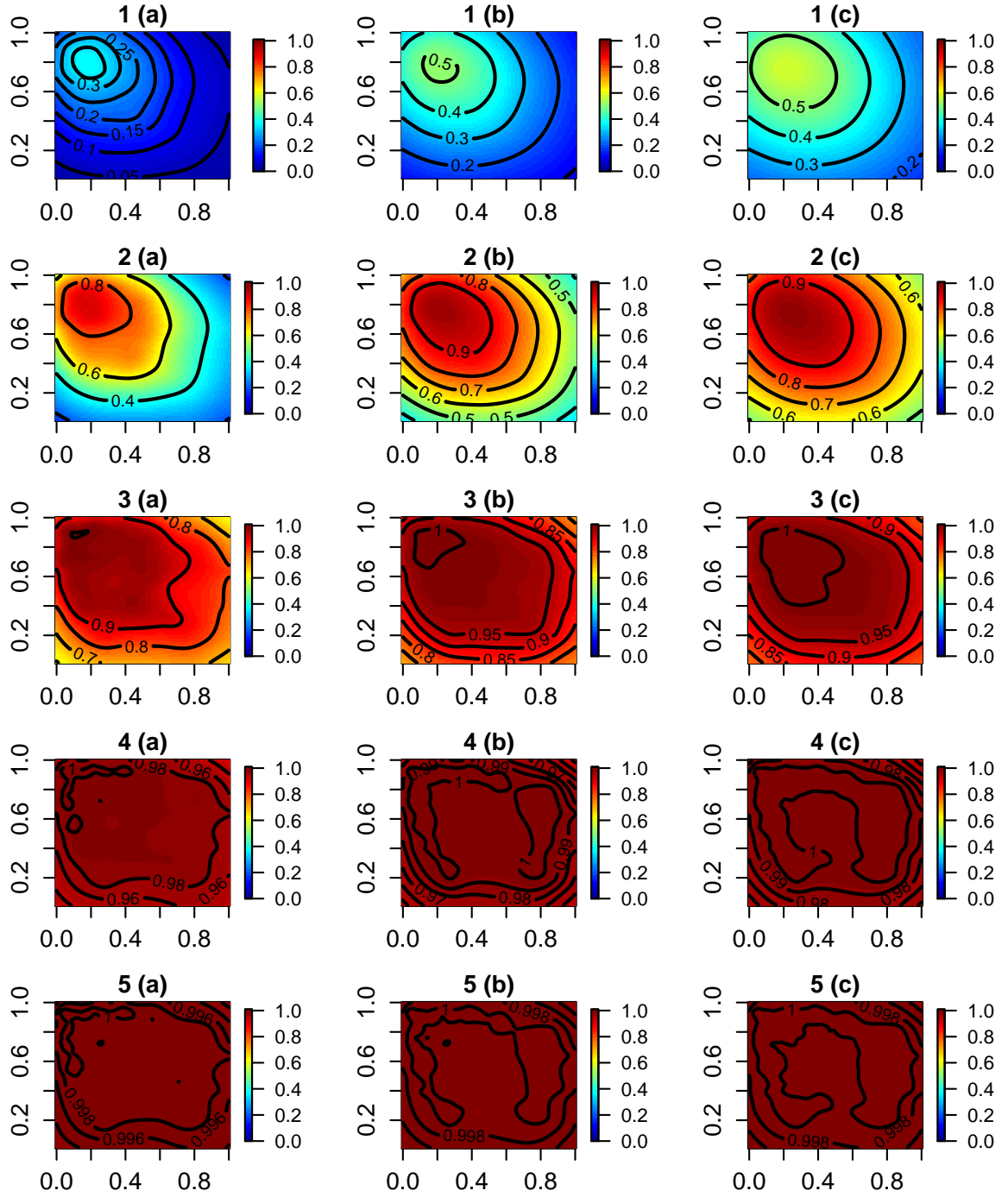


Figure 2: Interpolated surfaces of weights for the PCP for 15 combinations of variance ratio  $\delta_0$  and effective range  $d_0$ . Panels are given an alpha-numeric label. Numbers refer to the five values of  $\delta_0 = 0.01, 0.1, 1, 10, 100$ . Letters (a), (b) and (c) refer to three values of  $d_0 = \sqrt{2}/3, 2\sqrt{2}/3, \sqrt{2}$ .

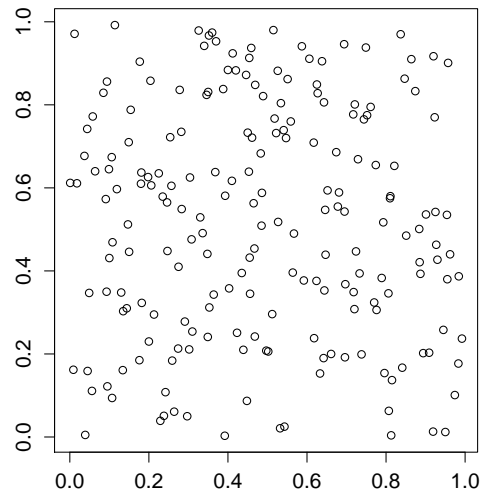


Figure 3: 200 randomly selected locations within the unit square.

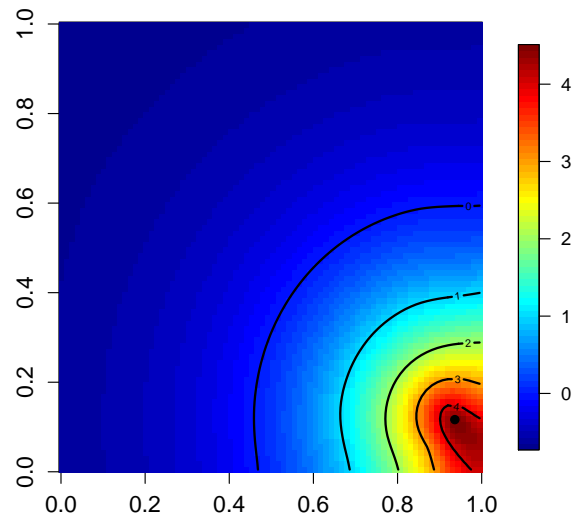


Figure 4: Interpolated surface of  $\mathbf{x}$  for the uniformly sampled data locations given in Figure 3.

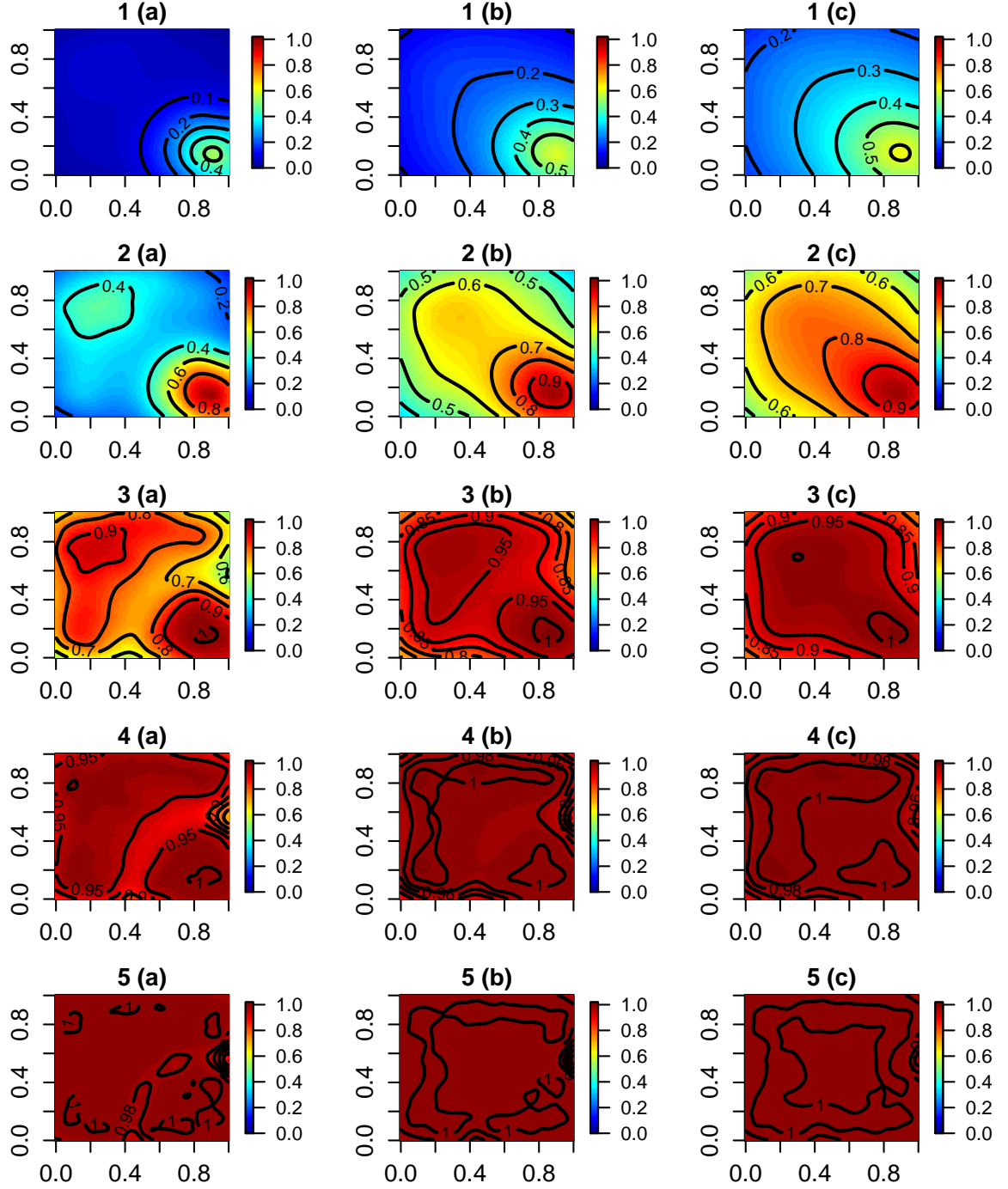


Figure 5: Interpolated surfaces of weights for the PCP for 15 combinations of variance ratio  $\delta_0$  and effective range  $d_1$ . Panels are given an alpha-numeric label. Numbers refer to the five values of  $\delta_1 = 0.01, 0.1, 1, 10, 100$ . Letters (a), (b) and (c) refer to three values of  $d_1 = \sqrt{2}/3, 2\sqrt{2}/3, \sqrt{2}$ .

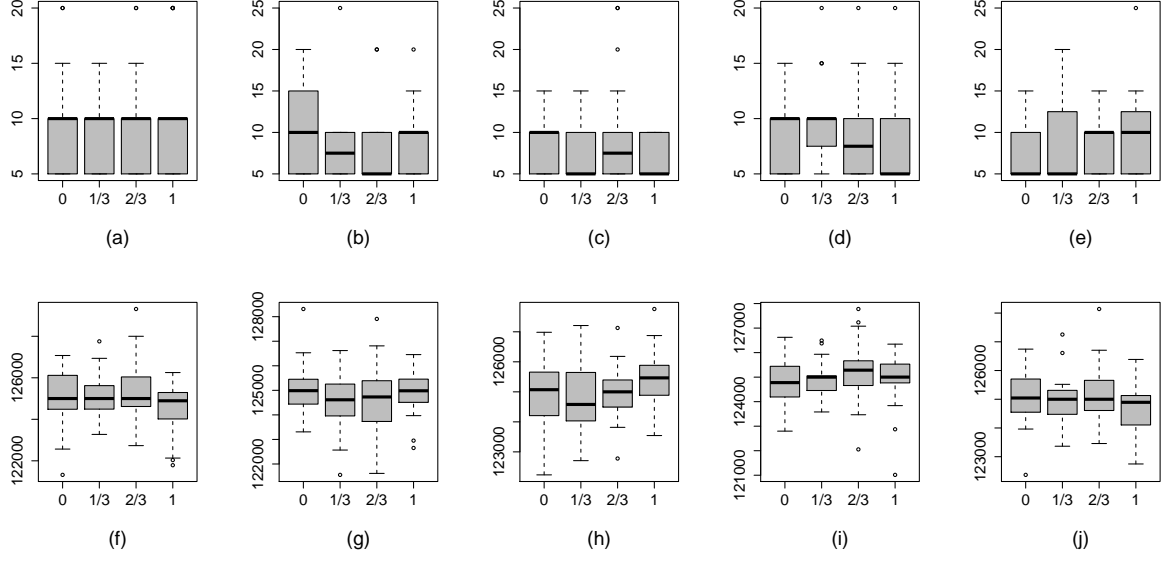


Figure 6:  $\text{PSRF}_M(1.1)$  and the ESS of  $\theta_0$  based on for the PCP with known variance parameters. Results are based on 5 chains each with 25,000 iterations. Plots (a)–(e) give the  $\text{PSRF}_M(1.1)$ , plots (f)–(j) the ESS of  $\theta_0$ . L–R  $\delta_0 = 0.01, 0.1, 1, 10, 100$ . Within each plot effective ranges of  $0, \sqrt{2}/3, 2\sqrt{2}/3$  and  $\sqrt{2}$  are used.

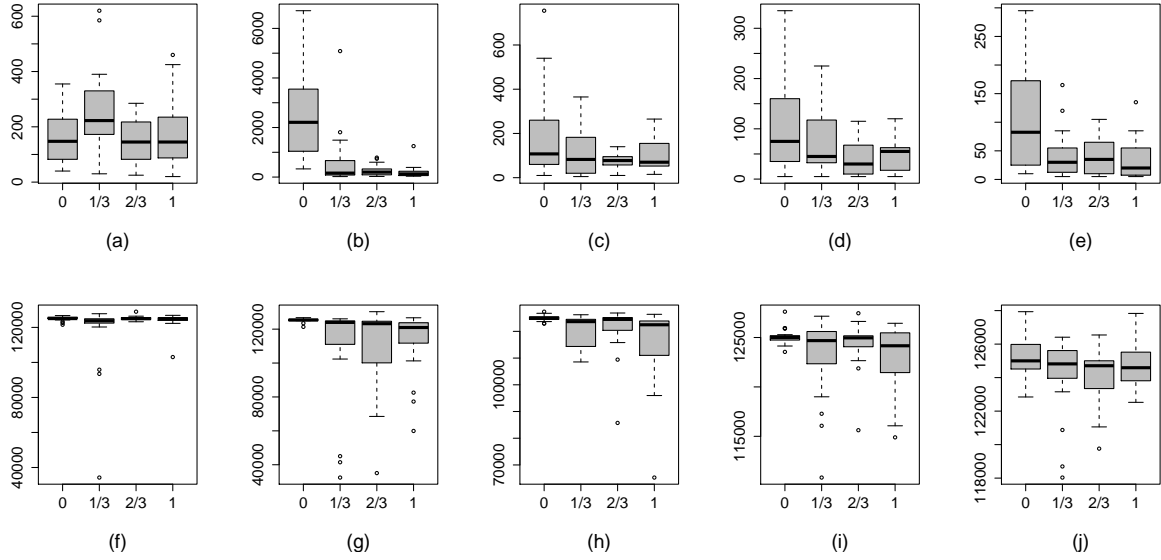


Figure 7:  $\text{MPSRF}_M(1.1)$  and the ESS of  $\theta_0$  for the PCP with unknown variance parameters. Results are based on 5 chains each with 25,000 iterations. Plots (a)–(e) give the  $\text{MPSRF}_M(1.1)$ , plots (f)–(j) the ESS of  $\theta_0$ . L–R  $\delta_0 = 0.01, 0.1, 1, 10, 100$ . Within each plot effective ranges of  $0, \sqrt{2}/3, 2\sqrt{2}/3$  and  $\sqrt{2}$  are used.

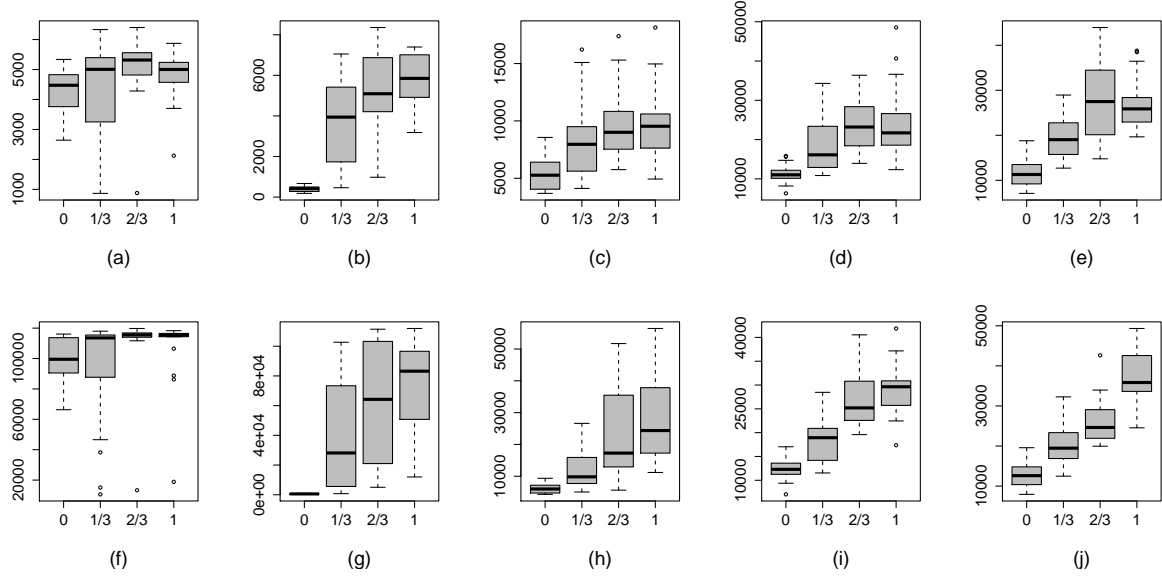


Figure 8: ESS of  $\sigma_0^2$  and  $\sigma_\epsilon^2$  for the PCP with unknown variance parameters. Results are based on 5 chains each with 25,000 iterations. Plots (a)–(e) give the ESS of  $\sigma_0^2$ , plots (f)–(j) the ESS of  $\sigma_\epsilon^2$ . L–R  $\delta_0 = 0.01, 0.1, 1, 10, 100$ . Within each plot effective ranges of 0,  $\sqrt{2}/3$ ,  $2\sqrt{2}/3$  and  $\sqrt{2}$  are used.

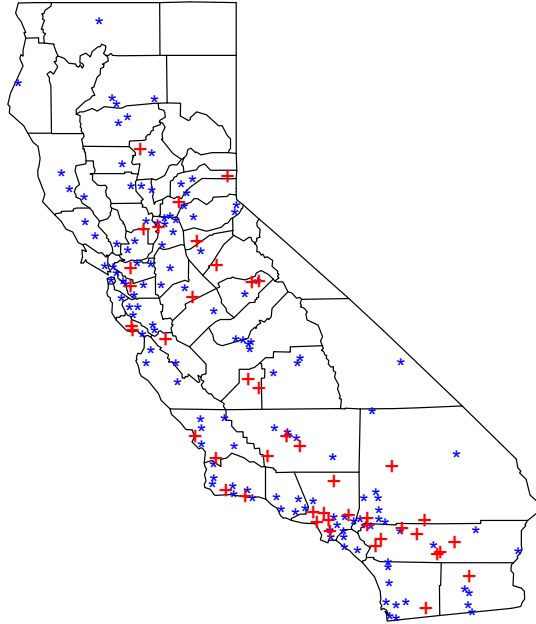


Figure 9: Sampling locations for Californian ozone concentration data. Blue stars indicate the locations of the 132 data sites, red crosses indicate the locations of the 44 validation sites

Table 3: Prediction error for different combinations of  $d_0$  and  $d_1$ 

$d_0$	$d_1$	MAPE	RMSPE	CRPS
50	50	15.84	19.45	11.15
	100	15.83	19.41	11.15
	250	15.87	19.40	11.14
	500	15.91	19.41	11.16
	1000	15.89	19.41	11.15
100	50	14.92	18.58	10.52
	100	14.97	18.58	10.55
	250	14.99	18.53	10.54
	500	14.98	18.50	10.53
	1000	15.01	18.53	10.54
250	50	<b>14.63</b>	18.39	10.42
	100	14.69	18.39	10.44
	250	14.70	18.33	10.43
	500	14.65	<b>18.27</b>	<b>10.39</b>
	1000	14.66	18.28	10.40
500	50	15.37	19.10	11.00
	100	15.36	19.06	10.99
	250	15.29	18.96	10.93
	500	15.30	18.94	10.93
	1000	15.28	18.93	10.93
1000	50	16.17	20.20	11.98
	100	16.24	20.22	11.99
	250	16.17	20.04	11.90
	500	16.19	20.05	11.91
	1000	16.24	20.05	11.94

Table 4: MPSRF $_M(1.1)$  and the ESS of the model parameters

	MPSRF $_M(1.1)$	ESS $\theta_0$	ESS $\theta_1$	ESS $\sigma_0^2$	ESS $\sigma_1^2$	ESS $\sigma_\epsilon^2$
CP	135	103129	56718	29539	5271	16836
NCP	1405	1252	3242	23371	4333	15797
PCP	160	125000	121995	28348	5082	15028

Table 5: MPSRF $_t(1.1)$  and ESS/s of the model parameters

	MPSRF $_t(1.1)$	ESS/s $\theta_0$	ESS/s $\theta_1$	ESS/s $\sigma_0^2$	ESS/s $\sigma_1^2$	ESS/s $\sigma_\epsilon^2$
CP	2.6	42.7	23.5	12.2	2.2	7.0
NCP	27.2	0.5	1.3	9.7	1.8	6.5
PCP	5.4	29.6	28.9	6.7	1.2	3.6

Table 6: Parameter estimates and their 95% credible intervals (CI)

Parameter	Estimate	95% CI
$\theta_0$	8.654	(8.197, 9.010)
$\theta_1$	-0.176	(-0.784, 0.397)
$\sigma_0^2$	0.677	(0.456, 0.958)
$\sigma_1^2$	0.360	(0.143, 0.768)
$\sigma_\epsilon^2$	0.137	(0.081, 0.218)
$\delta_0$	5.329	(2.428, 9.970)
$\delta_1$	2.808	(0.914, 6.716)

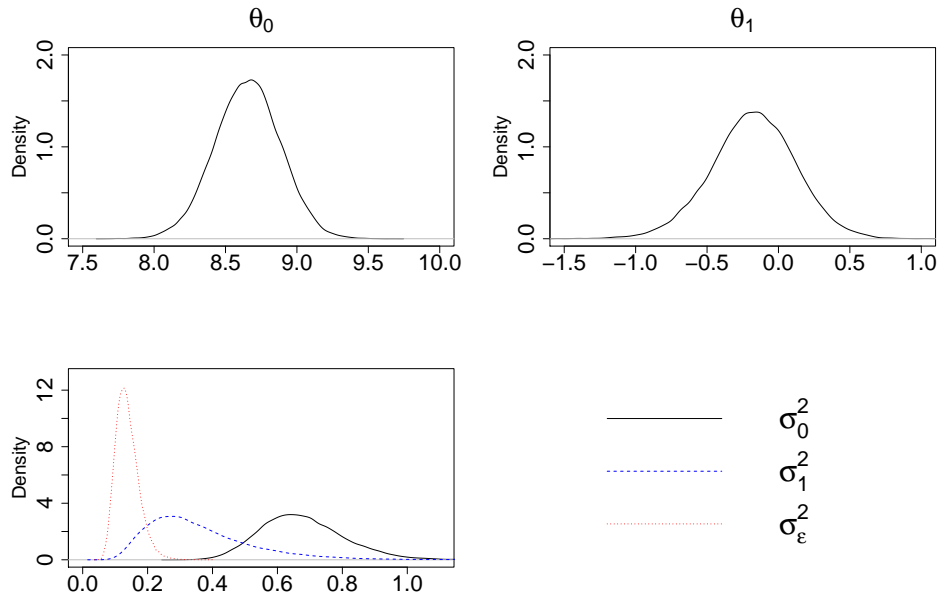


Figure 10: Density plots of model parameters for Californian ozone concentration data.

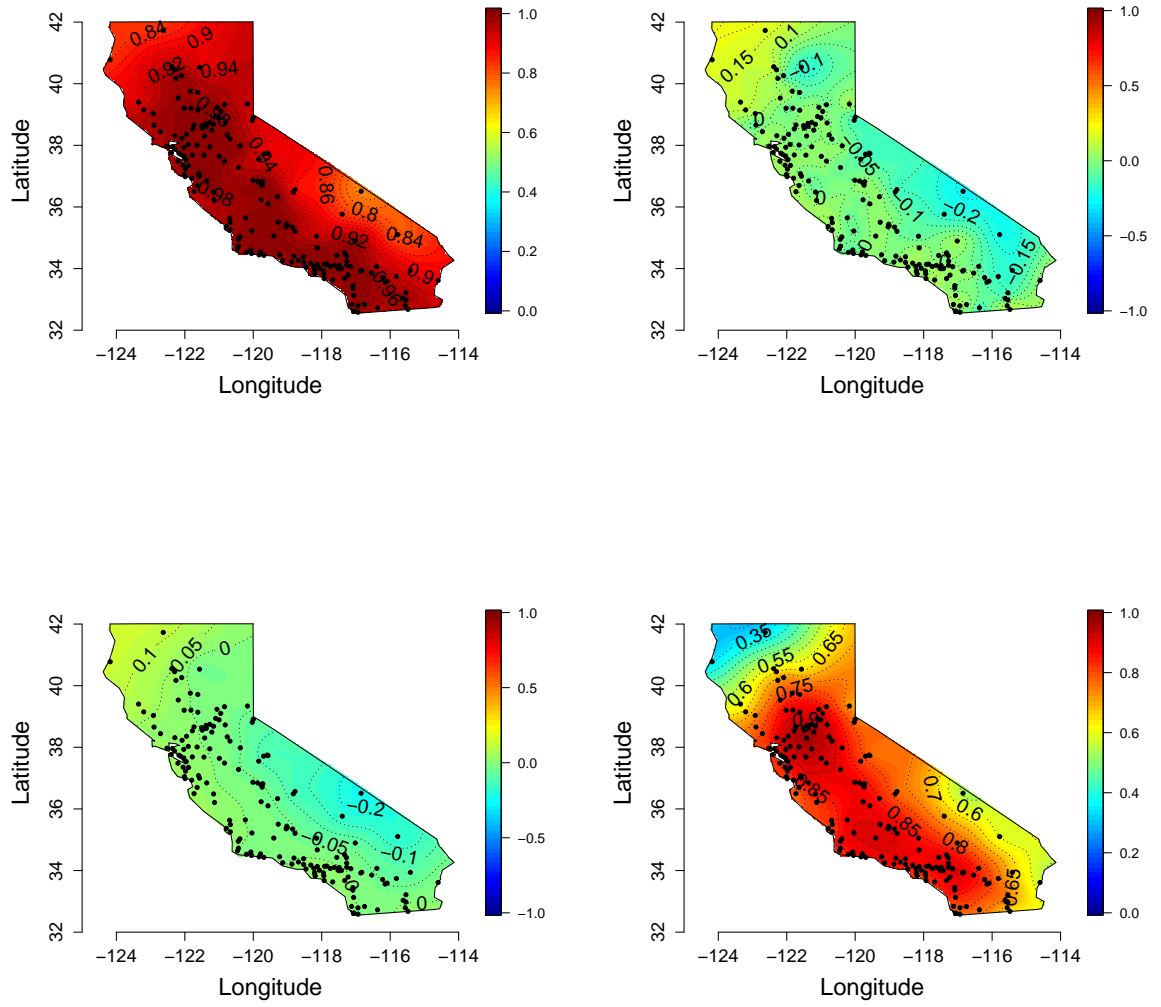


Figure 11: Spatially varying weights for Californian ozone concentration data. On the top row, left to right, are the weights assigned to  $\theta_0$  and  $\theta_1$  respectively for the intercept process  $\beta_0^w$  and on the bottom row are the equivalent plots for the slope process  $\beta_1^w$ .



## A Appendix

### A.1 Identifying the posterior precision matrix of the PCP

To identify the precision matrix of the joint posterior distribution of  $\beta^w$  and  $\theta$  for the PCP we write:

$$\begin{aligned}
\pi(\beta^w, \theta | \mathbf{y}) &\propto \pi(\mathbf{Y} | \beta^w, \theta) \pi(\beta^w | \theta) \pi(\theta) \\
&\propto \exp \left( -\frac{1}{2} \left[ \{ \mathbf{Y} - \mathbf{X}_1 \beta^w - \mathbf{X}_1 (\mathbf{I} - \mathbf{W}) \mathbf{X}_2 \}^T \mathbf{C}_1^{-1} \{ \mathbf{Y} - \mathbf{X}_1 \beta^w \right. \right. \\
&\quad \left. \left. - \mathbf{X}_1 (\mathbf{I} - \mathbf{W}) \mathbf{X}_2 \} + (\beta^w - \mathbf{W} \mathbf{X}_2 \theta)^T \mathbf{C}_2^{-1} (\beta^w - \mathbf{W} \mathbf{X}_2 \theta) \right. \right. \\
&\quad \left. \left. + (\theta - \mathbf{m})^T \mathbf{C}_3^{-1} (\theta - \mathbf{m}) \right] \right) \\
&= \exp \left( -\frac{1}{2} \left[ \dots + \beta^{wT} (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 + \mathbf{C}_2^{-1}) \beta^w \right. \right. \\
&\quad \left. \left. + 2 \beta^{wT} \{ \mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 (\mathbf{I} - \mathbf{W}) \mathbf{X}_2 - \mathbf{C}_2^{-1} \mathbf{W} \mathbf{X}_2 \} \theta + \theta^T \{ \mathbf{X}_2^T (\mathbf{I} - \mathbf{W})^T \right. \right. \\
&\quad \left. \left. \mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 (\mathbf{I} - \mathbf{W}) \mathbf{X}_2 + \mathbf{X}_2^T \mathbf{W}^T \mathbf{C}_2^{-1} \mathbf{W} \mathbf{X}_2 + \mathbf{C}_3^{-1} \} \theta + \dots \right] \right).
\end{aligned}$$

The entries of the precision matrix can then be read off of the final expression.

### A.2 Convergence rate of PCP for known variance parameters

Consider  $\mathbf{Q}_{\beta^w \theta}^{pc}$  and substitute  $\mathbf{W}$  from equation (9), then we have

$$\begin{aligned}
\mathbf{Q}_{\beta^w \theta}^{pc} &= \mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 (\mathbf{I} - \mathbf{W}) \mathbf{X}_2 - \mathbf{C}_2^{-1} \mathbf{W} \mathbf{X}_2 \\
&= (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1) \left\{ (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 + \mathbf{C}_2^{-1})^{-1} \mathbf{C}_2^{-1} \right\} \mathbf{X}_2 - \mathbf{C}_2^{-1} \left\{ \mathbf{I} - (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 \right. \\
&\quad \left. + \mathbf{C}_2^{-1})^{-1} \mathbf{C}_2^{-1} \right\} \mathbf{X}_2 \\
&= \left\{ (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1) (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 + \mathbf{C}_2^{-1})^{-1} \mathbf{C}_2^{-1} + \mathbf{C}_2^{-1} (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 + \mathbf{C}_2^{-1})^{-1} \mathbf{C}_2^{-1} \right. \\
&\quad \left. - \mathbf{C}_2^{-1} \right\} \mathbf{X}_2 \\
&= \left\{ (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 + \mathbf{C}_2^{-1}) (\mathbf{X}_1^T \mathbf{C}_1^{-1} \mathbf{X}_1 + \mathbf{C}_2^{-1})^{-1} \mathbf{C}_2^{-1} - \mathbf{C}_2^{-1} \right\} \mathbf{X}_2 \\
&= \left\{ \mathbf{C}_2^{-1} - \mathbf{C}_2^{-1} \right\} \mathbf{X}_2 \\
&= \mathbf{0}.
\end{aligned}$$

Therefore by setting  $\mathbf{W} = \mathbf{I} - \mathbf{B} \mathbf{C}_2^{-1}$ ,  $\mathbf{F}_{22}^{pc}$  becomes the null matrix and immediate convergence follows.

### A.3 Convergence rate of a three component Gibbs sampler

It can be shown that a Gibbs sampler with Gaussian target distribution with precision matrix given by

$$\mathbf{Q} = \begin{pmatrix} Q_{11} & Q_{12} & Q_{13} \\ Q_{21} & Q_{22} & Q_{23} \\ Q_{31} & Q_{32} & Q_{33} \end{pmatrix},$$

has a convergence rate which is equal to the maximum modulus eigenvalue of

$$F = \begin{pmatrix} 0 & -Q_{11}^{-1}Q_{12} & -Q_{11}^{-1}Q_{13} \\ 0 & Q_{22}^{-1}Q_{21}Q_{11}^{-1}Q_{12} & Q_{22}^{-1}Q_{21}Q_{11}^{-1}Q_{13} - Q_{22}^{-1}Q_{23} \\ 0 & F_{32} & F_{33} \end{pmatrix},$$

where

$$\begin{aligned} F_{32} &= (Q_{33}^{-1}Q_{31} - Q_{33}^{-1}Q_{32}Q_{22}^{-1}Q_{21})Q_{11}^{-1}Q_{12}, \\ F_{33} &= (Q_{33}^{-1}Q_{31} - Q_{33}^{-1}Q_{32}Q_{22}^{-1}Q_{21})Q_{11}^{-1}Q_{13} + Q_{33}^{-1}Q_{32}Q_{22}^{-1}Q_{23}. \end{aligned}$$

#### A.4 Stationarity of the PCP

To demonstrate that stationarity is preserved we let  $p = 1$  in model (4). The transition kernel of the Markov chain is

$$\begin{aligned} P\{\boldsymbol{\xi}^{(t+1)}|\boldsymbol{\xi}^{(t)}\} &= \pi\{\boldsymbol{\beta}^{w(t+1)}|\theta_0^{(t)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \pi\{\theta_0^{(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \\ &\quad \pi\{\sigma_0^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \\ &\quad \pi\{\sigma_\epsilon^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\}. \end{aligned}$$

We have dropped the  $\mathbf{W}$ 's to save space, conditioning the variance parameters on their current values where necessary. It follows that

$$\begin{aligned} &\int P\{\boldsymbol{\xi}^{(t+1)}|\boldsymbol{\xi}^{(t)}\} \pi(\boldsymbol{\xi}^{(t)}|\mathbf{y}) d\boldsymbol{\xi}^{(t)} \\ &= \int \pi\{\sigma_\epsilon^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \\ &\quad \pi\{\sigma_0^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \pi\{\theta_0^{(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \\ &\quad \underbrace{\pi\{\boldsymbol{\beta}^{w(t+1)}|\theta_0^{(t)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \left[ \int \pi\{\boldsymbol{\beta}_0^{w(t)}, \theta_0^{(t)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}|\mathbf{y}\} d\boldsymbol{\beta}_0^{w(t)} \right]}_{= \pi\{\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}|\mathbf{y}\}} d\theta_0^{(t)} d\sigma_0^{2(t)} d\sigma_\epsilon^{2(t)} \\ &= \int \pi\{\sigma_\epsilon^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \pi\{\sigma_0^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \\ &\quad \underbrace{\pi\{\theta_0^{(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \left[ \int \pi\{\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}|\mathbf{y}\} d\theta_0^{(t)} \right]}_{= \pi\{\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}|\mathbf{y}\}} d\sigma_0^{2(t)} d\sigma_\epsilon^{2(t)} \\ &= \int \pi\{\sigma_\epsilon^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \\ &\quad \underbrace{\left[ \int \pi\{\sigma_0^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \pi\{\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}|\mathbf{y}\} d\sigma_0^{2(t)} \right]}_{= \pi\{\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t)}|\mathbf{y}\}} d\sigma_\epsilon^{2(t)} \\ &= \int \pi\{\sigma_\epsilon^{2(t+1)}|\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \pi\{\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t)}|\mathbf{y}\} d\sigma_\epsilon^{2(t)} \\ &= \pi\{\boldsymbol{\beta}_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t+1)}|\mathbf{y}\} \\ &= \pi(\boldsymbol{\xi}^{(t+1)}|\mathbf{y}), \end{aligned} \tag{18}$$

and hence stationarity is preserved. The above argument can easily be extended for  $p > 1$  or to include other correlation parameters if they are being modelled.

If we update  $\mathbf{W}$  and the end of each complete pass of the sampler then the stationarity condition (14) does not hold. For instance, consider  $\sigma_\epsilon^2$ , which is conditioned on  $\sigma_0^2$  through  $\mathbf{W}$ . If  $\mathbf{W}$  is not recalculated using  $\sigma_0^{2(t+1)}$  then  $\sigma_\epsilon^{2(t+1)}$  is conditioned and  $\sigma_0^{2(t)}$ , and consequently

$$\int \pi\{\sigma_\epsilon^{2(t+1)}|\beta_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t)}, \sigma_\epsilon^{2(t)}, \mathbf{y}\} \pi\{\beta_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t)}|\mathbf{y}\} d\sigma_\epsilon^{2(t)} \\ \neq \pi\{\beta_0^{w(t+1)}, \theta_0^{(t+1)}, \sigma_0^{2(t+1)}, \sigma_\epsilon^{2(t+1)}|\mathbf{y}\},$$

but equality is required to complete step (18) in the string of equalities proving stationarity.

## A.5 Joint posterior and full conditional distributions

We begin here by writing down the joint posterior distribution of the parameters in model (4). We let  $\boldsymbol{\xi} = (\beta^{wT}, \boldsymbol{\theta}^T, \boldsymbol{\sigma}^{2T}, \sigma_\epsilon^2)^T$  be the vector containing all  $np$  partially centred random effects,  $p$  global effects,  $p$  random effect variances, the data variance and  $p$  decay parameters for the correlation functions. The joint posterior for  $\boldsymbol{\xi}$  is

$$\begin{aligned} \pi(\boldsymbol{\xi}|\mathbf{y}) &\propto \pi(\mathbf{Y}|\beta^w, \boldsymbol{\theta}, \sigma_\epsilon^2) \pi(\beta^w|\boldsymbol{\theta}, \boldsymbol{\sigma}^2) \pi(\boldsymbol{\theta}|\boldsymbol{\sigma}^2) \pi(\boldsymbol{\sigma}^2) \pi(\sigma_\epsilon^2) \\ &\propto \prod_{k=0}^{p-1} (\sigma_k^2)^{-(n/2+1/2+a_k+1)} |\mathbf{R}_k|^{-1/2} (\sigma_\epsilon^2)^{-(n/2+a_\epsilon+1)} \\ &\quad \exp \left\{ -\frac{1}{2\sigma_\epsilon^2} \left( [\mathbf{Y} - \mathbf{X}_1\{\beta^w + (\mathbf{I} - \mathbf{W})\mathbf{X}_2\boldsymbol{\theta}\}]^T [\mathbf{Y} - \mathbf{X}_1\{\beta^w + \right. \\ &\quad \left. (\mathbf{I} - \mathbf{W})\mathbf{X}_2\boldsymbol{\theta}\}] + 2b_\epsilon \right) \right\} \exp \left\{ -\frac{1}{2} (\beta^w - \mathbf{W}\mathbf{X}_2\boldsymbol{\theta})^T \mathbf{C}_2^{-1} (\beta^w - \mathbf{W}\mathbf{X}_2\boldsymbol{\theta}) \right\} \\ &\quad \exp \left[ -\frac{1}{2} \sum_{k=0}^{p-1} \frac{1}{\sigma_k^2} \left\{ \frac{(\theta_k - m_k)^2}{v_k} + 2b_k \right\} \right], \end{aligned}$$

where a description of the prior distributions  $\pi(\boldsymbol{\sigma}^2)$  and  $\pi(\sigma_\epsilon^2)$  can be found in Section 2.1.

It is argued in Section 2.3 that we must jointly update the  $\beta^w$ 's and jointly update  $\boldsymbol{\theta}$  and this is reflected in the conditional distributions given below.

- The full conditional distribution of  $\beta^w$  is

$$\beta^w|\boldsymbol{\theta}, \boldsymbol{\sigma}^2, \sigma_\epsilon^2, \mathbf{y} \sim N(\mathbf{m}_\beta^*, \mathbf{C}_2^*),$$

where

$$\begin{aligned} \mathbf{C}_2^* &= (\sigma_\epsilon^{-2} \mathbf{X}_1^T \mathbf{X}_1 + \mathbf{C}_2^{-1})^{-1}, \\ \mathbf{m}_\beta^* &= \mathbf{C}_2^* [\sigma_\epsilon^{-2} \{\mathbf{y} - \mathbf{X}_1(\mathbf{I} - \mathbf{W})\mathbf{X}_2\boldsymbol{\theta}\} + \mathbf{C}_2^{-1} \mathbf{W}\mathbf{X}_2\boldsymbol{\theta}]. \end{aligned}$$

- The full conditional distribution of  $\boldsymbol{\theta}$  is

$$\boldsymbol{\theta}|\beta^w, \boldsymbol{\sigma}^2, \sigma_\epsilon^2, \mathbf{y} \sim N(\mathbf{m}_\theta^*, \mathbf{C}_3^*),$$

where

$$\begin{aligned} \mathbf{C}_3^* &= [\sigma_\epsilon^{-2} \{\mathbf{X}_1(\mathbf{I} - \mathbf{W})\mathbf{X}_2\}^T \{\mathbf{X}_1(\mathbf{I} - \mathbf{W})\mathbf{X}_2\} + (\mathbf{W}\mathbf{X}_2)^T \mathbf{C}_2^{-1} \mathbf{W}\mathbf{X}_2 + \mathbf{C}_3^{-1}]^{-1}, \\ \mathbf{m}_\theta^* &= \mathbf{C}_3^* [\sigma_\epsilon^{-2} \{\mathbf{X}_1(\mathbf{I} - \mathbf{W})\mathbf{X}_2\}^T (\mathbf{y} - \mathbf{X}_1\beta^w) + (\mathbf{W}\mathbf{X}_2)^T \mathbf{C}_2^{-1} \beta^w + \mathbf{C}_3^{-1} \mathbf{m}]. \end{aligned}$$

- The full conditional distribution of  $\sigma_k^2$ ,  $k = 0, \dots, p-1$ , is

$$\sigma_k^2 | \boldsymbol{\beta}^w, \boldsymbol{\theta}, \boldsymbol{\sigma}^2_{-k}, \sigma_\epsilon^2, \mathbf{y} \sim IG \left[ \frac{n+1}{2} + a_k, \frac{1}{2} \left\{ \left( \boldsymbol{\beta}_k^w - \sum_{m=0}^{p-1} \mathbf{W}_{km} \theta_k \mathbf{1} \right)^T \mathbf{R}_k^{-1} \left( \boldsymbol{\beta}_k^w - \sum_{m=0}^{p-1} \mathbf{W}_{km} \theta_k \mathbf{1} \right) + \frac{(\theta_k - m_k)^2}{v_k} + 2b_k \right\} \right],$$

where  $\mathbf{W}_{km}$  denotes the  $km$ th,  $n \times n$  block of  $\mathbf{W}$ .

- The full conditional distribution of  $\sigma_\epsilon^2$  is

$$\sigma_\epsilon^2 | \boldsymbol{\beta}^w, \boldsymbol{\theta}, \boldsymbol{\sigma}^2, \mathbf{y} \sim IG \left\{ \frac{n}{2} + a_\epsilon, \frac{1}{2} \left( [\mathbf{Y} - \mathbf{X}_1 \{ \boldsymbol{\beta}^w + (\mathbf{I} - \mathbf{W}) \mathbf{X}_2 \boldsymbol{\theta} \}]^T [\mathbf{Y} - \mathbf{X}_1 \{ \boldsymbol{\beta}^w + (\mathbf{I} - \mathbf{W}) \mathbf{X}_2 \boldsymbol{\theta} \}] + 2b_\epsilon \right) \right\}.$$

## References

- Bakar, K. S. and Sahu, S. K. (2015). spTimer: Spatio-temporal Bayesian modeling using R. *Journal of Statistical Software*, **63**(15), 1–32.
- Banerjee, S., Gelfand, A. E., Finley, A. O., and Sang, H. (2008). Gaussian predictive process models for large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **70**(4), 825–848.
- Bass, M. R. and Sahu, S. K. (2016). Efficient parameterisations of gaussian process based models for bayesian computation using MCMC. Technical report, University of Southampton.
- Berrocal, V. J., Gelfand, A. E., and Holland, D. M. (2010). A spatio-temporal downscaler for output from numerical models. *Journal of Agricultural, Biological, and Environmental statistics*, **15**(2), 176–197.
- Brooks, S. P. and Gelman, A. (1998). General methods for monitoring convergence of iterative simulations. *Journal of Computational and Graphical Statistics*, **7**(4), 434–455.
- Cressie, N. and Johannesson, G. (2008). Fixed rank kriging for very large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **70**(1), 209–226.
- Finley, A. O., Banerjee, S., and MacFarlane, D. W. (2011). A hierarchical model for quantifying forest variables over large heterogeneous landscapes with uncertain forest areas. *Journal of the American Statistical Association*, **106**(493), 31–48.
- Finley, A. O., Banerjee, S., and Gelfand, A. E. (2015). spBayes for large univariate and multivariate point-referenced spatio-temporal data models. *Journal of Statistical Software*, **63**(13), 1–28.
- Furrer, R., Nychka, D., and Sain, S. (2009). fields: Tools for spatial data. *R package version*, **6**(11).
- Gelfand, A. E. and Smith, A. F. (1990). Sampling-based approaches to calculating marginal densities. *Journal of the American statistical association*, **85**(410), 398–409.
- Gelfand, A. E., Sahu, S. K., and Carlin, B. P. (1995). Efficient parameterisations for normal linear mixed models. *Biometrika*, **82**(3), 479–488.
- Gelfand, A. E., Kim, H.-J., Sirmans, C., and Banerjee, S. (2003). Spatial modeling with spatially varying coefficient processes. *Journal of the American Statistical Association*, **98**(462), 387–396.
- Gelfand, A. E., Sahu, S. K., and Holland, D. M. (2012). On the effect of preferential sampling in spatial prediction. *Environmetrics*, **23**(7), 565–578.
- Gelman, A. and Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statistical Science*, **7**(4), 457–472.
- Gneiting, T., Balabdaoui, F., and Raftery, A. E. (2007). Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **69**(2), 243–268.
- Hamm, N., Finley, A., Schaap, M., and Stein, A. (2015). A spatially varying coefficient model for mapping PM10 air quality at the European scale. *Atmospheric Environment*, **102**, 393–405.

- Handcock, M. S. and Stein, M. L. (1993). A Bayesian analysis of kriging. *Technometrics*, **35**(4), 403–410.
- Harville, D. A. (1997). *Matrix Algebra from a Statistician's Perspective*. Springer-Verlag New York.
- Huerta, G., Sansó, B., and Stroud, J. R. (2004). A spatiotemporal model for Mexico City ozone levels. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, **53**(2), 231–248.
- Matérn, B. (1986). *Spatial Variation*. Springer Verlag, Berlin, 2nd. edition.
- Papaspiliopoulos, O., Roberts, G. O., and Sköld, M. (2003). Non-centered parameterisations for hierarchical models and data augmentation (with discussion). In *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*, pages 307–326. Oxford University Press, USA.
- Robert, C. O. and Casella, G. (2004). *Monte Carlo Statistical Methods*. Springer-Verlag New York, 2nd. edition.
- Roberts, G. O. and Sahu, S. K. (1997). Updating schemes, correlation structure, blocking and parameterization for the Gibbs sampler. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **59**(2), 291–317.
- Sahu, S. K., Gelfand, A. E., and Holland, D. M. (2007). High resolution space–time ozone modeling for assessing trends. *Journal of the American Statistical Association*, **102**(480), 1221–1234.
- Sahu, S. K., Gelfand, A. E., and Holland, D. M. (2010). Fusing point and areal level space–time data with application to wet deposition. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, **59**(1), 77–103.
- Sahu, S. K., Yip, S., and Holland, D. M. (2011). A fast Bayesian method for updating and forecasting hourly ozone levels. *Environmental and Ecological Statistics*, **18**(1), 185–207.
- Wheeler, D. C., Páez, A., Spinney, J., and Waller, L. A. (2014). A Bayesian approach to hedonic price analysis. *Papers in Regional Science*, **93**(3), 663–683.
- Yu, Y. and Meng, X.-L. (2011). To center or not to center: That is not the question: an Ancillarity–Sufficiency Interweaving Strategy (ASIS) for boosting MCMC efficiency. *Journal of Computational and Graphical Statistics*, **20**(3), 531–570.
- Zhang, H. (2004). Inconsistent estimation and asymptotically equal interpolations in model-based geostatistics. *Journal of the American Statistical Association*, **99**(465), 250–261.