



## spNNGP R Package for Nearest Neighbor Gaussian Process Models

Andrew O. Finley   
Michigan State University

Abhirup Datta   
Johns Hopkins University

Sudipto Banerjee   
University of California,  
Los Angeles

---

### Abstract

This paper describes and illustrates functionality of the **spNNGP** R package. The package provides a suite of spatial regression models for Gaussian and non-Gaussian point-referenced outcomes that are spatially indexed. The package implements several Markov chain Monte Carlo (MCMC) and MCMC-free nearest neighbor Gaussian process (NNGP) models for inference about large spatial data. Non-Gaussian outcomes are modeled using a NNGP Pólya-Gamma latent variable. **OpenMP** parallelization options are provided to take advantage of multiprocessor systems. Package features are illustrated using simulated and real data sets.

*Keywords:* MCMC, nearest neighbor Gaussian process, kriging, R.

---

## 1. Introduction

This paper introduces the **spNNGP** (Finley, Datta, and Banerjee 2022) R (R Core Team 2022) package that provides a suite of spatial regression models for Gaussian and non-Gaussian univariate outcomes observed at point-referenced two-dimensional locations. There are, by now, many R packages that provide similar basic functionality. A recent read of the “Analysis of Spatial Data” Comprehensive R Archive Network (CRAN) Task View (Bivand and Nowosad 2022) yielded  $\sim 46$  packages listed for geostatistical analysis – and this is not an exhaustive accounting of packages available for such analyses. Our software design focus and unique contribution is to provide Bayesian models and associated diagnostic and prediction functions capable of handling data sets with a large number of locations via the nearest neighbor Gaussian process (NNGP, Datta, Banerjee, Finley, and Gelfand 2016). Specifically, functions in **spNNGP** implement recent methodological and algorithmic developments presented in Finley, Datta, Cook, Morton, Andersen, and Banerjee (2019). Package **spNNGP** (Finley *et al.* 2022) is available from CRAN at <https://CRAN.R-project.org/package=spNNGP>.

There have been many recent methodological developments within the large spatial data literature that aim to deliver massively scalable spatial processes. Sun, Li, and Genton (2011) and Banerjee (2017) provide background and discussion of current work in this area. A recent contribution by Heaton *et al.* (2019) is particularly useful as it provides an overview of modeling approaches for large spatial data that are under active software development, and a comparison of these approaches based on the analysis of a common dataset in the form of a “friendly competition”. In addition to NNGP models, the comparison presented by Heaton *et al.* (2019) considered covariance tapering via the **spam** package (Furrer and Sain 2010; Furrer, Flury, and Gerber 2021), gap-filling via **gapfill** (Gerber 2021), metakriging (Guhaniyogi and Banerjee 2018), spatial partitioning (Sang, Jun, and Huang 2011; Barbian and Assunção 2017), fixed rank kriging via **FRK** (Cressie and Johannesson 2008; Zammit-Mangion and Cressie 2021), multiresolution approximation (Katzfuss 2017), stochastic partial differential equations via **INLA** (Rue *et al.* 2021), lattice kriging via **LatticeKrig** (Nychka, Bandyopadhyay, Hammerling, Lindgren, and Sain 2015; Nychka, Hammerling, Sain, Lenssen, Smirniotis, and Iverson 2019), local approximate Gaussian processes via **laGP** (Gramacy and Apley 2015; Gramacy 2016), and reduced rank predictive processes (Banerjee, Gelfand, Finley, and Sang 2008; Finley, Sang, Banerjee, and Gelfand 2009) via **spBayes** (Finley, Banerjee, and Gelfand 2015). The comparison was based on out-of-sampled predictive performance and, to a lesser extent, computing time for a moderately sized simulated and real dataset comprising 105,569 observations. Comparisons showed NNGP models yielded highly competitive predictive performance and computation time. More recently, Turek and Risser (2022) developed the **BayesNSGP** package for nonstationary Gaussian process modeling with options to use NNGPs for large data settings (Risser and Turek 2020). In a frequentist setup, fast maximum likelihood-based parameter estimation and predictions using nearest neighbor approximations to the Gaussian process likelihood are available in the **GpGp** (Guinness 2018) and **BRISC** (Saha and Datta 2022) packages on CRAN. The latter also offers inference on the spatial covariance parameters using a fast spatial bootstrap (Saha and Datta 2018). While most of the software noted above exploit sparsity in the spatial covariance or precision matrix, or pursue low-rank approximations, the **ExaGeoStat** package (Abdulah, Ltaief, Sun, Genton, and Keyes 2018) tackles decomposition of the full dense spatial covariance matrix head-on using high performance linear algebra libraries associated with various leading edge parallel architectures.

Our contribution here is many fold. We propose a novel extension for analyzing spatially correlated non-Gaussian (binary) responses using a NNGP spatial generalized linear mixed model (GLMM) by using a Pólya-Gamma prior (Polson, Scott, and Windle 2013) for the regression coefficients, which leads to an efficient data-augmented Gibbs sampler. This is, to our knowledge, the first sampler for binomial spatial-GLMM that ensures closed form Gibbs updates for most parameters using linear time and storage. We discuss model comparison and model adequacy and show how different variants of the NNGP models entail fundamentally different implementation and interpretation of model comparison metrics. This is an important pragmatic consideration for practitioners choosing among different NNGP models. Finally, we provide detailed documentation and exposition of the user-friendly **spNNGP** R package that: (1) implements NNGP model fitting algorithms for Gaussian response spatial data presented in Datta *et al.* (2016) and Finley *et al.* (2019); (2) provides NNGP models for the non-Gaussian spatial response via the Pólya-Gamma data-augmented sampler; (3) offers support functions for NNGP model fit diagnostics, summary, and prediction. We discuss code optimization

aspects that allows **spNNGP** to handle very large data sets by being judicious with memory use, taking advantage of properties of the NNGP dependence scheme to efficiently store and retrieve neighbor information, and applying parallel processing where advantageous. Core model fitting and prediction functions in **spNNGP** have achieved unprecedented scalability, delivering inference for data sets in the hundreds of millions of spatial locations.

The remainder of this article proceeds as follows. Section 2.1 provides a brief overview of Gaussian process models followed by specifics about the NNGP models in Section 2.2. In Section 2.3 we propose extensions for binomial spatial data and outline the Pólya-Gamma data-augmented Gibbs sampler. Section 2.4 discusses appropriate model comparison and adequacy measures for the different NNGP models. Section 2.5 gives a brief description of the code underlying **spNNGP** and some software development specifics for large data sets. This is followed by analysis of simulated and real data sets in Section 3 meant to provide a practical tour of some of the package’s features. Finally, Section 4 provides a brief summary with an eye toward future development.

## 2. Models and software

### 2.1. Review of Gaussian process models for spatial data

The standard geostatistical paradigm envisions each data unit as a triplet  $(\mathbf{s}_i, \mathbf{x}(\mathbf{s}_i), y(\mathbf{s}_i))$ , where  $\mathbf{s}_i$  denotes the geographical location of the measurements,  $\mathbf{x}(\mathbf{s}_i)$  is the  $p \times 1$  vector of covariates, and  $y(\mathbf{s}_i)$  is the response of interest, for  $i = 1, 2, \dots, n$ . If the covariates fail to account for all of the structured variation observed on the response, a spatial linear mixed effects model for analyzing the data is specified as

$$y(\mathbf{s}_i) = \mathbf{x}(\mathbf{s}_i)^\top \boldsymbol{\beta} + w(\mathbf{s}_i) + \epsilon(\mathbf{s}_i), \quad (1)$$

where  $\boldsymbol{\beta}$  is the  $p \times 1$  vector of regression coefficients,  $\epsilon(\mathbf{s}_i)$  is the random noise, customarily modeled as independent and identically distributed (iid) observations from  $N(0, \tau^2)$ , and  $w(\mathbf{s}_i)$  is the location-specific spatial random effect. Typically, the spatial random effects are assumed to be smooth across space, i.e.,  $w(\mathbf{s}_i)$  can be conceived of as realizations of a smooth latent surface  $\{w(\mathbf{s}) \mid \mathbf{s} \in \mathcal{D}\}$ , where  $\mathcal{D}$  denotes the geographical domain of interest. Gaussian processes (GP) are widely used in machine learning for modeling smooth functions on the many-dimensional covariate domain (Rasmussen 2003). For spatial regression, the mean function of the covariates is often specified parsimoniously via a linear model as in (1) and GPs are typically used to model the latent surface  $w(\mathbf{s})$  on the two- or three-dimensional physical domain. A GP model for the spatial surface  $\{w(\mathbf{s})\}$ , denoted by  $w(\mathbf{s}) \sim GP(0, C(\cdot, \cdot \mid \boldsymbol{\theta}))$  where  $C(\cdot, \cdot \mid \boldsymbol{\theta})$  is a covariance function, implies that for any finite set of locations  $\mathbf{s}_1, \dots, \mathbf{s}_n$ , the vector of random effects  $\mathbf{w} = (w(\mathbf{s}_1), \dots, w(\mathbf{s}_n))^\top$  follows a zero-mean multivariate Gaussian distribution with covariance matrix  $\mathbf{C}(\boldsymbol{\theta}) = \mathbf{C} = (c_{ij})$  where  $c_{ij} = C(\mathbf{s}_i, \mathbf{s}_j \mid \boldsymbol{\theta})$ . The parametric covariance function  $C(\cdot, \cdot \mid \boldsymbol{\theta})$  is often selected from the Matérn family of functions (Stein 2012), popular due to its versatility to model surfaces with varying degrees of smoothness.

The mixed effects model for the response and the GP model for the random effects can be combined into the hierarchical model

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta} + \mathbf{w}, \tau^2\mathbf{I}), \quad \mathbf{w} \sim N(\mathbf{0}, \mathbf{C}(\boldsymbol{\theta})), \quad (2)$$

where  $\mathbf{y}$  is the vector formed by stacking the  $y(\mathbf{s}_i)$ 's, and  $\mathbf{X}$  is the corresponding  $n \times p$  matrix of covariates. Equivalently, one can integrate out  $\mathbf{w}$  from (2) and write

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{C}(\boldsymbol{\theta}) + \tau^2\mathbf{I}). \quad (3)$$

Parameter estimation in a frequentist paradigm maximizes the likelihood from (3) with respect to  $\boldsymbol{\beta}$ ,  $\tau^2$  and  $\boldsymbol{\theta}$ , while in a Bayesian framework, priors are assigned to these parameters, and one can use either the hierarchical model (2) or the marginal model (3) to obtain posterior inference via Markov chain Monte Carlo (MCMC).

## 2.2. Nearest neighbor Gaussian processes for large spatial data

Gaussian processes encounter computational roadblocks when data is observed at a large number of locations. Both the joint likelihood from the hierarchical model (2) or the data likelihood from the marginalized model (3) involves a multivariate Gaussian distribution with a dense  $n \times n$  covariance matrix ( $\mathbf{C}(\boldsymbol{\theta})$  and  $\mathbf{C}(\boldsymbol{\theta}) + \tau^2\mathbf{I}$ , respectively). This task involves  $O(n^2)$  storage, and  $O(n^3)$  computations (floating point operations or FLOPs), which is infeasible when  $n$  is large.

One of the scalable solutions is replacing the GP prior for the spatial random effects with a *nearest neighbor Gaussian process* prior (Datta *et al.* 2016). Let  $\mathcal{R} = \{\mathbf{s}_1, \dots, \mathbf{s}_n\}$  be any fixed finite set of locations in the spatial domain  $\mathcal{D}$ , and  $\mathbf{w}_{\mathcal{R}} = (w(\mathbf{s}_1), \dots, w(\mathbf{s}_n))^{\top}$ . For any location in  $\mathcal{D}$ , define neighbor sets as follows

$$\begin{aligned} N(\mathbf{s}_1) &= \{\} \text{ (empty set),} \\ N(\mathbf{s}_i) &= \min(m, i-1) \text{ nearest neighbors of } \mathbf{s}_i \text{ in } \mathbf{s}_1, \dots, \mathbf{s}_{i-1}, \text{ for } i = 2, \dots, n, \\ N(\mathbf{s}) &= m \text{ nearest neighbors of } \mathbf{s} \text{ in } \mathcal{R}. \end{aligned} \quad (4)$$

Then NNGP is specified as

$$\mathbf{w}_{\mathcal{R}} \sim \prod_{i=1}^n p(w_i | \mathbf{w}(N(\mathbf{s}_i))); \quad w(\mathbf{s}) | \mathbf{w}_{\mathcal{R}} \stackrel{\text{ind}}{\sim} p(w(\mathbf{s}) | \mathbf{w}(N(\mathbf{s}))) \text{ for all } \mathbf{s} \in \mathcal{D} \setminus \mathcal{R}. \quad (5)$$

In practice,  $\mathcal{R}$  is chosen to be the data-locations and (5) yields the NNGP approximation to the distribution of the spatial effects. If (5) is applied to the response vector, then we obtain the likelihood approximation. This is motivated from the likelihood approximation ideas in Vecchia (1988) and Stein, Chi, and Welty (2004). Datta *et al.* (2016) showed that the above construction endows a multivariate Gaussian distribution on  $\mathbf{w} = \mathbf{w}_{\mathcal{R}}$

$$\mathbf{w} \sim N(\mathbf{0}, \tilde{\mathbf{C}}(\boldsymbol{\theta})), \quad (6)$$

where  $\tilde{\mathbf{C}}(\boldsymbol{\theta})^{-1}$ , i.e., the inverse of the NNGP covariance matrix, is sparse. The extension to arbitrary locations in the bottom-row of (5) is based on kriging (conditional) distributions using only  $m$ -nearest neighbors from the data-locations. Nearest-neighbor kriging have been explored in Emery (2009) in the context of point-predictions as opposed to the full predictive (conditional) distributions specified for NNGP. The NNGP is itself a valid stochastic (Gaussian) process on the entire domain  $\mathcal{D}$  – a chief distinction from frequentist likelihood approximations or prediction equations. If probabilistic modeling of a spatial surface (observed or latent) is of primary concern, then the NNGP is a legitimate candidate and need

not be viewed as an approximation of the parent GP from which it is derived. Only when inference on the spatial covariance parameters are of interest, NNGP needs to be perceived as an approximation, as these parameters describe attributes of the parent GP. The NNGP covariance function  $\tilde{C}$  is constructed from the original covariance function  $C$  and ensures that the matrix  $\tilde{\mathbf{C}}(\boldsymbol{\theta})^{-1}$  is sparse and, consequently, the likelihood (6) can be evaluated using only  $O(n)$  storage. This makes NNGP a scalable replacement for the full GP, while delivering inference almost indistinguishable from the full GP. NNGP can, in fact, also be looked upon as a special case of a Gaussian Markov random field (GMRF, [Rue and Held 2005](#)) with a specific neighborhood structure defined through the neighbor sets of (4). In practice, however, NNGP delivers inference sufficiently close to traditional Gaussian random fields without requiring mesh-based finite element approximations of stochastic partial differential equation (SPDE) representations of Gaussian random fields (as done in [Lindgren, Rue, and Lindström 2011](#)).

### Latent NNGP

The original implementation of the NNGP model proposed in [Datta et al. \(2016\)](#) used a fully Bayesian hierarchical specification

$$N(\mathbf{y} \mid \mathbf{X}\boldsymbol{\beta} + \mathbf{w}, \tau^2\mathbf{I}) \times N(\mathbf{w} \mid \mathbf{0}, \tilde{\mathbf{C}}(\boldsymbol{\theta})) \times p(\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2) \quad (7)$$

for running an MCMC algorithm, where all the parameters  $\{\mathbf{w}, \boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$  are updated in a Gibbs sampler. Normal priors for  $\boldsymbol{\beta}$  and inverse-gamma priors for the variance components ensure that they produce conjugate full conditionals in the Gibbs sampler. The remaining covariance parameters are updated using random-walk Metropolis steps for their respective full conditional distributions. The full conditional distribution for  $\mathbf{w}$  from (7) is

$$\mathbf{w} \mid \cdot \sim N(\mathbf{B}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})/\tau^2, \mathbf{B}) \text{ where } \mathbf{B} = \tilde{\mathbf{C}}(\boldsymbol{\theta})^{-1} + \mathbf{I}/\tau^2.$$

Unfortunately, this block update of  $\mathbf{w}$  is not practical. While the full conditional precision matrix  $\tilde{\mathbf{C}}(\boldsymbol{\theta})^{-1} + \mathbf{I}/\tau^2$  has the same sparsity as  $\tilde{\mathbf{C}}(\boldsymbol{\theta})^{-1}$ , unlike  $\tilde{\mathbf{C}}(\boldsymbol{\theta})^{-1}$  its determinant cannot be calculated in  $O(n)$  FLOPs. Instead, [Datta et al. \(2016\)](#) recommended sequentially updating the full conditionals  $w_i \mid \cdot$  for  $i = 1, \dots, n$  and outlined an algorithm that accomplishes this entire sequence of updates in  $O(n)$  FLOPs. We refer to this NNGP model as the *latent NNGP*.

### Response NNGP

The latent NNGP algorithm involves running an MCMC of dimension  $O(n)$  where each of the  $n$  parameters  $w_i$  are sequentially updated. While this ensures linear scalability of the NNGP model with sample size, it can also imply very slow convergence of the high-dimensional MCMC. Instead of using NNGP for the latent Gaussian process  $w(\mathbf{s})$ , [Finley et al. \(2019\)](#) directly considered the marginal Gaussian process for the response, i.e.,  $\{y(\mathbf{s})\} \sim GP(\mathbf{x}(\mathbf{s})^\top \boldsymbol{\beta}, \Sigma(\cdot, \cdot))$  where the marginalized covariance function  $\Sigma$  is specified as  $\Sigma(\mathbf{s}_i, \mathbf{s}_j) = C(\mathbf{s}_i, \mathbf{s}_j \mid \boldsymbol{\theta}) + \tau^2 \delta(\mathbf{s}_i, \mathbf{s}_j)$ ,  $\delta$  denoting the Kronecker delta. Since the covariance function of an NNGP can be derived from any parent GP, [Finley et al. \(2019\)](#) replaced the covariance function  $\Sigma$  with its NNGP analogue  $\tilde{\Sigma}$  yielding the response model

$$\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \tilde{\Sigma}), \quad (8)$$

where  $\tilde{\Sigma}$  is the NNGP covariance matrix derived from  $\Sigma = \mathbf{C}(\boldsymbol{\theta}) + \tau^2 \mathbf{I}$ . The advantage of this *response NNGP* model over the previous latent model is that the dimension of the parameter space is drastically reduced from  $O(n)$  to  $O(1)$ . The lower dimensional NNGP tends to have improved MCMC convergence as shown in [Finley et al. \(2019\)](#). The computational scalability per MCMC iteration of the response model remains the same as in the latent model algorithm, as  $\tilde{\Sigma}^{-1}$  is sparse requiring  $O(n)$  storage and its quadratic form and determinant can be calculated using  $O(n)$  FLOPs.

### *MCMC-free inference: Conjugate NNGP model*

The conjugate NNGP algorithm is an implementation of the response NNGP model which fixes certain spatial covariance parameters leading to exact (MCMC-free) posterior Bayesian inference. Recall that under the full-GP specification, the covariance function for  $y(\mathbf{s})$  is specified as  $\Sigma(\mathbf{s}_i, \mathbf{s}_j) = C(\mathbf{s}_i, \mathbf{s}_j | \boldsymbol{\theta}) + \tau^2 \delta(\mathbf{s}_i, \mathbf{s}_j)$ . Usually, the covariance functions  $C(\cdot, \cdot | \boldsymbol{\theta})$  can be expressed as  $\sigma^2 R(\cdot, \cdot | \boldsymbol{\phi})$  where  $\sigma^2$  is the marginal variance, and  $R$  is the correlation function parameterized by  $\boldsymbol{\phi}$ , i.e.,  $\boldsymbol{\theta} = (\sigma^2, \boldsymbol{\phi})$ . Rewriting  $\tau^2 = \alpha \sigma^2$ , we have  $\Sigma(\mathbf{s}_i, \mathbf{s}_j) = \sigma^2 (R(\mathbf{s}_i, \mathbf{s}_j | \boldsymbol{\phi}) + \alpha \delta(\mathbf{s}_i, \mathbf{s}_j))$ . The conjugate NNGP model fixes  $\boldsymbol{\phi}$  and  $\alpha$ , and generates the NNGP covariance function approximation  $\tilde{M}(\cdot, \cdot | \alpha, \boldsymbol{\phi})$  of  $R(\cdot, \cdot | \boldsymbol{\phi}) + \alpha \delta(\cdot, \cdot)$ . This implies we have the following marginal model

$$\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \tilde{\mathbf{M}}) \quad (9)$$

where  $\tilde{\mathbf{M}} = \tilde{\mathbf{M}}(\boldsymbol{\phi}, \alpha)$  is a known covariance matrix once  $\boldsymbol{\phi}$  and  $\alpha$  are fixed. The model in (9) is the standard Bayesian linear model with only unknowns  $\boldsymbol{\beta}$  and  $\sigma^2$ . Using a normal-inverse-gamma prior for  $(\boldsymbol{\beta}, \sigma^2)$  leads to conjugate normal-inverse-gamma posterior distributions and hence posterior moments and other summary quantities of  $\boldsymbol{\beta}$  and  $\sigma^2$  are easily and exactly obtained. Exact posterior predictive distributions for  $y(\mathbf{s}_0)$  at a new location  $\mathbf{s}_0$  are also available. The matrix  $\tilde{\mathbf{M}}^{-1}$  (like  $\tilde{\Sigma}^{-1}$  for the response model) is sparse ( $O(n)$ ) ensuring all the posterior distributions and moments can be evaluated using  $O(n)$  memory and FLOPs. The fixed values of  $\boldsymbol{\phi}$  and  $\alpha$  are either chosen based on a variogram or can be selected in a more formal fashion using cross-validation using predictions from the model on hold-out data. We describe an implementation of the cross-validation in more detail in the *MCMC-free inference* subsection of Section 3. While the cross-validation enforces multiple runs of the conjugate model for a grid of values of  $\boldsymbol{\phi}$  and  $\alpha$ , these runs can proceed in parallel. Also,  $\boldsymbol{\phi}$  is usually one or two dimensional, hence, there are only 2 or 3 tuning parameters for the model and empirical results in [Finley et al. \(2019\)](#) suggest that a crude-resolution grid often suffices to yield accurate predictive inference. Empirical comparisons in [Finley et al. \(2019\)](#) and also in [Heaton et al. \(2019\)](#) suggest that the conjugate NNGP is orders of magnitude faster than the MCMC-based NNGP algorithms and also other competing big-spatial data methods while delivering highly competitive prediction performance.

### 2.3. Binomial response

All current implementations of the NNGP model assume that the response is Gaussian. Here we implement a Bayesian NNGP model for binomial response. Note that the response and the conjugate NNGP model explicitly rely on Gaussian distributions to derive the marginal distribution for the response  $\mathbf{y}$ . Such closed form marginal distributions are not available for non-Gaussian responses. However, conceptually we can extend the latent NNGP model

of (7) to non-Gaussian settings. Assuming  $y(\mathbf{s}_i) \stackrel{ind}{\sim} \text{Binomial}(n_i, \frac{1}{1+\exp(-\psi(\mathbf{s}_i))})$ , where  $\psi(\mathbf{s}_i) = \mathbf{x}(\mathbf{s}_i)^\top \boldsymbol{\beta} + w(\mathbf{s}_i)$  with  $w(\mathbf{s}_i)$  modeled as a NNGP, the joint likelihood for this model is given by:

$$\prod_{i=1}^n \text{Binomial}(y(\mathbf{s}_i) \mid n(\mathbf{s}_i), \frac{1}{1+\exp(-\psi(\mathbf{s}_i))}) \times N(\mathbf{w} \mid \mathbf{0}, \tilde{\mathbf{C}}(\boldsymbol{\theta})) \times p(\boldsymbol{\beta}, \boldsymbol{\theta}). \quad (10)$$

Since  $\mathbf{w}$  is given an NNGP prior, this likelihood can still be evaluated using  $O(n)$  memory and FLOPs. However, unlike the latent NNGP model for the Gaussian case, where the full conditional distributions  $w(\mathbf{s}_i) \mid \cdot$  used in the Gibbs updates were Gaussian distributions, for the non-Gaussian case these full conditionals  $w(\mathbf{s}_i) \mid \cdot$  do not belong to any standard family. Alternatively, one can resort to Metropolis-Hastings (MH) random walk updates for each  $w(\mathbf{s}_i)$ . Introducing  $n$  MH random walk updates in every iteration would exacerbate the slow convergence issues already plaguing the latent NNGP model.

To eschew the numerous Metropolis random-walk updates, we devise a Gibbs sampler for the latent NNGP model for binomial responses exploiting the Pólya-Gamma data-augmented sampler of Polson *et al.* (2013). Again, we write  $\tilde{\mathbf{C}}(\boldsymbol{\theta}) = \sigma^2 \tilde{\mathbf{R}}(\boldsymbol{\phi})$  where  $\tilde{\mathbf{R}}$  is the NNGP approximation of the GP correlation matrix. We then assume a  $N(\boldsymbol{\mu}, \mathbf{V})$  prior for  $\boldsymbol{\beta}$ , an  $IG(a, b)$  prior for  $\sigma^2$ , and  $p(\boldsymbol{\phi})$  prior for the other covariance parameters  $\boldsymbol{\phi}$ . Letting  $\kappa(\mathbf{s}_i) = y(\mathbf{s}_i) - n(\mathbf{s}_i)/2$  we introduce the augmented data  $\boldsymbol{\omega} = (\omega(\mathbf{s}_1), \dots, \omega(\mathbf{s}_n))^\top$ . We can then write the conditional likelihood

$$\begin{aligned} p(\boldsymbol{\beta}, \boldsymbol{\theta}, \mathbf{w} \mid \mathbf{y}, \boldsymbol{\omega}) \propto & \prod_{i=1}^n \exp\left(-\frac{\omega(\mathbf{s}_i)}{2} (\mathbf{x}(\mathbf{s}_i)^\top \boldsymbol{\beta} + w(\mathbf{s}_i) - \kappa(\mathbf{s}_i)/\omega(\mathbf{s}_i))^2\right) \times \\ & \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \mathbf{w}^\top \tilde{\mathbf{R}}(\boldsymbol{\phi})^{-1} \mathbf{w}\right) \times \\ & \exp\left(-\frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\mu})^\top \mathbf{V}^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu})\right) \times \left(\frac{1}{\sigma^2}\right)^{a+1} \exp\left(\frac{b}{\sigma^2}\right) \times p(\boldsymbol{\phi}). \end{aligned} \quad (11)$$

Using the augmented data  $\boldsymbol{\omega}$ , this likelihood for the binomial mixed linear model is similar to the usual likelihood from a spatial mixed linear model with NNGP using responses  $y(\mathbf{s}_i)^* = \kappa(\mathbf{s}_i)/\omega(\mathbf{s}_i)$  and heteroskedastic variances  $\tau^2(\mathbf{s}_i) = 1/\omega(\mathbf{s}_i)$ .

An immediate consequence is that the updates for  $\boldsymbol{\beta}$  and  $\mathbf{w}$  are analogous to the updates in the NNGP with  $y(\mathbf{s}_i)$  replaced with  $y(\mathbf{s}_i)^*$  and the constant nugget  $\tau^2$  replaced by  $\tau^2(\mathbf{s}_i)$ . Note that for  $y(\mathbf{s}_i)^*$  and  $\tau^2(\mathbf{s}_i)$  are functions of the augmented data  $\boldsymbol{\omega}$  and will change in every iteration. Other than that, the updates remain exactly the same as in the latent model. The only additional update we need to do for non-Gaussian responses is that of the  $\omega(\mathbf{s}_i)$ 's. We update these using the full conditionals

$$\omega(\mathbf{s}_i) \sim PG(n(\mathbf{s}_i), \mathbf{x}(\mathbf{s}_i)^\top \boldsymbol{\beta} + w(\mathbf{s}_i)),$$

where  $PG(b, z)$  denotes the Pólya-Gamma random variable with shape parameter  $b$  and tilting parameter  $z$  (see Polson *et al.* 2013, for more details about this distribution).

This completes a Gibbs sampler for spatial GLMM with binomial responses, where, like the Gaussian case, every parameter except the covariance parameters ( $\boldsymbol{\phi}$ ) have closed form full conditionals. The memory requirements also remain  $O(n)$  thereby ensuring that the binomial NNGP model is a viable candidate for analysis of massive non-Gaussian spatial data.

## 2.4. Fitted values, replicates, and predictions

Subsequent to convergence of the MCMC for the latent or response models, we can use the posterior samples of the parameters to generate various quantities of interest like fitted values and replicates at each data location, and predictions at new locations. For the response model, samples for the fitted value of the mean at the  $i$ -th location is simply given by  $\{\mathbf{x}(\mathbf{s}_i)^\top \boldsymbol{\beta}^{(l)}\}$  where  $\boldsymbol{\beta}^{(l)}$  denotes the  $l$ -th post burn-in sample. Similarly, for the latent NNGP, the posterior distribution of the fitted values is specified by the samples  $\{\mathbf{x}(\mathbf{s}_i)^\top \boldsymbol{\beta}^{(l)} + w(\mathbf{s}_i)^{(l)}\}$ . For the MCMC-free conjugate NNGP, let  $\hat{\phi}$  and  $\hat{\alpha}$ , respectively, denote the chosen values of  $\phi$  and  $\alpha$  based on cross-validation using RMSE of predictions on hold-out (test) data. Then we can express the posterior predictive distribution as a  $t$  distribution with mean  $\mathbf{x}(\mathbf{s}_i)^\top E_{\hat{\phi}, \hat{\alpha}}(\boldsymbol{\beta} \mid \mathbf{y})$  where  $E_{\hat{\phi}, \hat{\alpha}}(\boldsymbol{\beta} \mid \mathbf{y})$  denotes the posterior mean of  $\boldsymbol{\beta}$  at  $\phi = \hat{\phi}$  and  $\alpha = \hat{\alpha}$ . The closed-form expressions for  $E_{\hat{\phi}, \hat{\alpha}}(\boldsymbol{\beta} \mid \mathbf{y})$  and the scale parameter for the  $t$  distribution for the fitted value are available in Algorithm 5 of Finley *et al.* (2019).

Following Gelman, Carlin, Stern, Dunson, Vehtari, and Rubin (2013) we consider generating replicate data at the observed data locations. This is often used for model checking and model adequacy evaluations and is based upon simulating each data point from its posterior predictive distribution. We note that while both the latent model (2) and the response model (3) lead to the same marginal distribution for  $\mathbf{y}$ , the two models differ fundamentally in their definition for replicate data. For the hierarchical model, conditional on  $\boldsymbol{\beta}$ ,  $\boldsymbol{\theta}$  and  $\mathbf{w}$ 's, the  $y(\mathbf{s}_i)$ 's are independent with iid  $N(0, \tau^2)$  distributed error terms. Hence, for each post-burn-in sample of the parameters, replicates at each data location can be generated as  $\{y(\mathbf{s}_i)^{(l)} \sim N(\mathbf{x}(\mathbf{s}_i)^\top \boldsymbol{\beta}^{(l)} + w(\mathbf{s}_i)^{(l)}, \tau^{2(l)})\}$ . For the response model, however, we can think of  $\mathbf{y}$  as a single multivariate observation from  $N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Sigma})$ . Things are further complicated when switching to the NNGP covariance matrix  $\tilde{\boldsymbol{\Sigma}}$  as, unlike  $\boldsymbol{\Sigma}$ , we cannot express  $\tilde{\boldsymbol{\Sigma}}$  as the sum of a purely spatial covariance matrix and a diagonal matrix of error variances. Hence, replicates for the marginalized model are basically new multivariate draws from  $N(\mathbf{X}\boldsymbol{\beta}^{(l)}, \tilde{\boldsymbol{\Sigma}}(\boldsymbol{\theta}^{(l)}))$ . Unlike the replicates from the hierarchical model, which are expected to exhibit strong spatial alignment with the observed data owing to the correlation induced through the use of the common  $\mathbf{w}$ 's, the replicates for the response model are only correlated with the observed data through  $\boldsymbol{\beta}$  and  $\boldsymbol{\theta}$  and hence are not expected to have similar spatial contours as the observed data. For the conjugate NNGP, which also relies on the marginalized model, the replicates are draws from  $N(\mathbf{X}\boldsymbol{\beta}_{\hat{\phi}, \hat{\alpha}}^{(l)}, \sigma_{\hat{\phi}, \hat{\alpha}}^{2(l)} \tilde{\mathbf{M}}(\hat{\phi}, \hat{\alpha}))$  with  $\{(\boldsymbol{\beta}_{\hat{\phi}, \hat{\alpha}}^{(l)}, \sigma_{\hat{\phi}, \hat{\alpha}}^{2(l)})\}$  denoting the normal-inverse-gamma posterior samples of  $\boldsymbol{\beta}$  and  $\sigma^2$  at the chosen value of  $\hat{\phi}$  and  $\hat{\alpha}$ .

Finally, we turn our attention to predictions at a new location  $\mathbf{s}$ . For the latent model, predictions are also generated hierarchically by first generating samples of  $w(\mathbf{s})^{(l)} \mid \mathbf{w}(N(\mathbf{s}))^{(l)}, \boldsymbol{\theta}^{(l)}$ . This conditional distribution is Gaussian with mean and variance being the kriging mean and variance at  $\mathbf{s}$  based on a set of nearest neighbors  $N(\mathbf{s})$ . The exact expressions for these quantities are specified in Algorithm 2 of Finley *et al.* (2019). Subsequent to generating samples of  $w(\mathbf{s})$ , we generate replicates for the response as  $y(\mathbf{s})^{(l)} \mid w(\mathbf{s})^{(l)}, \boldsymbol{\beta}^{(l)}, \tau^{2(l)} \sim N(\mathbf{x}(\mathbf{s})^\top \boldsymbol{\beta}^{(l)} + w(\mathbf{s})^{(l)}, \tau^{2(l)})$ . For the response model, the prediction algorithm is presented in Algorithm 4 of Finley *et al.* (2019) and involves directly generating samples of  $y(\mathbf{s})^{(l)} \mid y(N(\mathbf{s})), \boldsymbol{\beta}^{(l)}, \boldsymbol{\theta}^{(l)}$ . This is once again a conditional normal distribution equivalent to kriging at  $\mathbf{s}$  given its neighbors  $N(\mathbf{s})$  for the response process  $y(\cdot)$ . For the conjugate model, exact predictive distributions are available as  $y(\mathbf{s}) \mid \text{data}$  following a  $t$  distribution with location

and scale parameters provided in Algorithm 5 of [Finley \*et al.\* \(2019\)](#). Finally, for binomial responses, predictions follow the same strategy as the latent model for Gaussian responses as both models rely on the hierarchical latent variable formulation. We initially generate samples of  $w(\mathbf{s})^{(l)} \mid \mathbf{w}(N(\mathbf{s}))^{(l)}, \boldsymbol{\theta}^{(l)}$  followed by generating binomial samples for the response as  $y(\mathbf{s})^{(l)} \mid w(\mathbf{s})^{(l)}, \boldsymbol{\beta}^{(l)}, \tau^2^{(l)} \sim \text{Binomial}(n(\mathbf{s}), (1 + \exp(-(\mathbf{x}(\mathbf{s})^\top \boldsymbol{\beta}^{(l)} + w(\mathbf{s})^{(l)})))^{-1})$  where  $n(\mathbf{s})$  denotes the total count at  $\mathbf{s}$ .

## 2.5. Software features

Two model fitting functions `spNNGP` and `spConjNNGP`, with associated support functions, offer users a set of efficient regression and prediction tools that implement the methods outlined in Sections 2.2 and 2.3. Following from Section 2.2, `spNNGP` provides MCMC-based inference for the latent and response spatial regression models. The latent model can be specified for Gaussian and binomial (via the Pólya-Gamma distribution Section 2.3) outcome variables, while the response model only fits Gaussian outcome variables. The highly efficient MCMC-free conjugate NNGP algorithm used for Gaussian outcome variables described in Section 2.2 is available via the `spConjNNGP` function.

The `spNNGP` and `spConjNNGP` functions as well as their support functions (`fitted` for generating regression fitted and replicated data, `spDiag` for computing model diagnostics, and `predict` for sampling from posterior predictive distributions) are written in C/C++ using R’s foreign language interface and offer parallelization using `openMP` ([Dagum and Menon 1998](#)).

As noted previously, our aim was to provide software capable of handling data sets with 10s to 100s of millions of locations. This aim was met by minimizing memory requirements, and taking advantage of parallelization where possible. We attempted to minimize memory requirements in several ways. First, we avoid making copies of input data, e.g., the regression design matrix, spatial coordinate matrix, etc. Second, we do not hold any Euclidean distance or subsequent spatial covariance vectors and matrices in memory, but rather compute them when needed. For example, all sampling algorithms require the computation of the spatial covariance vector between each observed location and its set of nearest neighbors, and the among neighbors spatial covariance matrix. For any given observation this storage requirement is not large, i.e., at most a  $m$  length vector and upper or lower triangle of an  $m \times m$  matrix, where  $m$  is the number of neighbors (which is small, e.g.,  $m = 15$ ). However, when  $n$  is large, storing  $n$   $m$  length vectors and  $m \times m$  matrices is substantial. Therefore, we made the decision to compute these covariance vectors and matrices “on the fly” for each observation and for each iteration of the given sampler. To do this efficiently, we find each observation’s neighbor set via a fast code book search described in [Ra and Kim \(1993\)](#) which is modified to accommodate the ordering. The neighbor search is performed only once at the beginning of the program, and the result is  $n$  vectors of at most  $m$  integers that record the data row index for each observation’s neighbors (i.e., the indexes in the input outcome vector and corresponding design and coordinate matrices).

The ordering and subsequent selection of  $m$ -(directed) nearest neighbors can be summarized as a directed acyclic neighbor graph. This neighbor graph can be encoded as a sparse lower-triangular matrix where each row is an observation and the non-zero elements in each row hold the neighbor indexes. Within the software, this matrix is held in compressed row storage (CRS) format to minimize memory use and neighbors’ data retrieval time. The CRS format does not store zeros and organizes non-zero matrix elements ordered by row in contiguous

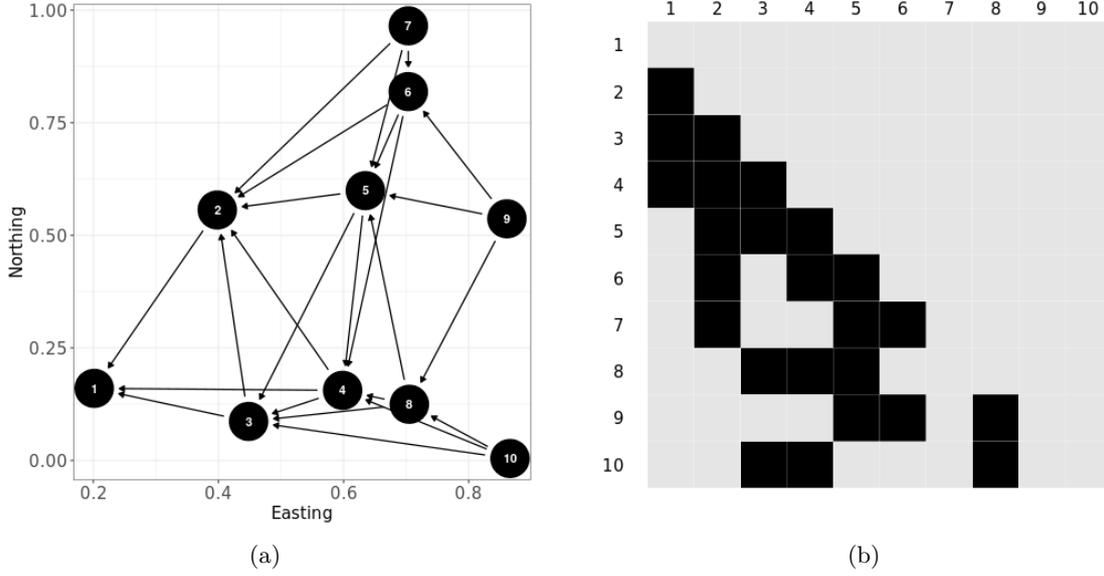


Figure 1: (a) Illustration of a neighbor graph using three neighbors for ten observations with ordering along the easting axis. (b) The sparse lower-triangular matrix (black elements are non-zero) corresponding to the graph in (a) that records the neighbor index (columns) for each observation index (rows).

memory. This is advantageous in our setting because it makes traversing all observations' neighbor indexes as simple as a single loop over the vector that holds the CRS matrix's non-zero elements. Given the index for a given neighbor, retrieving its spatial coordinates, design matrix, and outcome value is further simplified because all input data are stored in contiguous memory column-major format, which allows for fast retrieval of data using integer indexing. This straightforward looping and use of integer indexing to retrieve neighbor input data facilitates parallelization using an **openMP** `omp for pragma` at several points in the sampling algorithms.

Figure 1 provides an example of a neighbor graph and the corresponding sparse matrix that we store in CRS format. Following the sampling algorithms in [Datta \*et al.\* \(2016\)](#) and [Finley \*et al.\* \(2019\)](#), the response and conjugate models require only each observation's neighbor information; however, the latent model sampling algorithm requires information on which observations have a given observation as a neighbor. Or put another way, we need to know which neighbor sets each observation is a member of. This information is easily accessed by noting the row index of each column's non-zero elements in Figure 1(b), e.g., observation 1 is in the neighbor sets for observation 2, 3, and 4. Efficiently accessing and traversing these indexes is done in the software by converting the CRS neighbor index matrix to a compressed column storage (CCS) format, which effectively reorders the needed row indexes in contiguous memory by column.

Storing all vectors and matrices in contiguous memory column-major format also simplifies calls to Fortran basic linear algebra subprograms (**BLAS**; <https://www.netlib.org/blas>) and linear algebra package (**LAPACK**; <https://www.netlib.org/lapack>), which are used for all computationally intensive matrix operations. Hence, additional speed-up can be real-

ized if R calls a threaded implementation of **BLAS**, e.g., **openBLAS** (Zhang 2016) or Intel's math kernel library (MKL, Intel Corp. 2019), while working on a multiple processor computer. Finley *et al.* (2019) provide details about where and how each sampling algorithm uses parallelization.

### 3. Illustrations

#### 3.1. Analysis of simulated data

The basic functionality of **spNNGP** and **spConjNNGP**, along with their support functions, are illustrated using a small  $n = 5000$  simulated data set with locations distributed at random within a unit square domain. This data set was kept purposely small so the NNGP models could be compared with the full GP model. The Gaussian outcome variable was generated from (2) with latent  $\mathbf{w}$  centered on zero and covariance matrix elements  $c_{ij} = \sigma^2 \exp(-\phi \| \mathbf{s}_i - \mathbf{s}_j \|)$ , where  $\phi = 6$  and  $\sigma^2 = 1$ . The columns in the design matrix comprise an intercept and variable  $\mathbf{x}$  which was drawn from a normal distribution with mean zero and variance 1. The regression coefficients were  $\boldsymbol{\beta} = (1, -0.1)^\top$  and measurement error  $\tau^2 = 0.1$ . An additional  $n_0 = 2500$  observations on a grid were set aside to illustrate prediction.

#### *MCMC-based inference*

A call to **spNNGP** generates samples from parameters' posterior distributions and, optionally, corresponding samples of regression fitted values and replicates via composition sampling as described in Section 2.4. Similar to R's **lm** function, the desired model is passed to **spNNGP** via the **formula** argument. The model's outcome vector and design matrix components are then found in a data frame passed to the **data** argument or, if **data** is not specified, taken from the calling environment. Additionally, the locations corresponding to the observed data are passed as a matrix via the **coords** argument. One can also pass **coords** a character vector of column names in **data**'s data frame.

Specifying the latent or response model is done via the **method** argument. As detailed in Finley *et al.* (2019), for the latent model (initialized by setting **method** = "latent")  $\boldsymbol{\beta}$ ,  $\mathbf{w}$ ,  $\sigma^2$ , and  $\tau^2$  are updated from their respective full conditional distributions via a Gibbs algorithm. The spatial correlation function decay parameter  $\phi$  and, if **cov.model** = "matern" (where **cov.model** specifies the desired spatial correlation function), smoothness parameter  $\nu$  are updated via a MH algorithm. If **method** = "response", only  $\boldsymbol{\beta}$  has an efficient closed form full conditional Gibbs update and the remaining parameters  $\sigma^2$ ,  $\tau^2$ ,  $\phi$ , and perhaps  $\nu$  are updated via MH. For both models, starting values and prior distributions must be specified for all parameters (with the optional exception of  $\boldsymbol{\beta}$  and  $\mathbf{w}$ ) via the **starting** and **priors** arguments, as illustrated below. Those parameters updated via MH are transformed to have support on the real line so that a normal proposal distribution can be used. The variance of the parameter specific proposal distribution is controlled via the **tuning** argument. Tuning values should be selected to maintain the desired MH acceptance rate (in general we aim for  $\sim 25 - 50$  percent, see Roberts and Rosenthal (2009) and Gelman *et al.* (2013) for guidance). Information about acceptance rate and other model specifics is printed to the console when **verbose** = TRUE.

A NNGP model requires two user inputs: ordering of the locations, and number of nearest

neighbors  $m$  to use. The top-row of Equation 5 clearly depends on the ordering of the locations in  $\mathcal{R}$ . As defined in Equation 4, given an ordering, for a data-location corresponding to an observation, the neighbor set (conditioning set) is constructed as the set of  $m$  nearest neighbors of  $\mathbf{s}$  that conform to the user-specified ordering (i.e., directed nearest neighbors). By default, `spNNGP` and `spConjNNGP` order in increasing value of the first column of the matrix passed to `coords`, which, in most settings, produces an excellent approximation to a full GP (see supplemental experiments in [Datta et al. \(2016\)](#)); however, as demonstrated by [Guinness \(2018\)](#) improvements in parameter estimation efficiency can be achieved with different ordering designs. Functions to obtain these alternate orderings are provided in the **GpGP** R package ([Guinness 2018](#)). If the user wishes to try different orderings, they can be passed as an integer index vector via the optional `ord` argument. For illustration, we bypass the default ordering in the call to `spNNGP` below and order using the sum of locations'  $x$  and  $y$  coordinate values. We reemphasize that, subsequent to selection of any ordering, NNGP always uses nearest neighbors corresponding to that ordering. Other neighbor selection schemes such as mixture of nearest and farthest neighbors have been considered in related earlier work ([Stein et al. 2004](#)) but is not implemented in `spNNGP`. The number of neighbors to consider is controlled by the `n.neighbors` argument. As demonstrated in [Datta et al. \(2016\)](#) 15 neighbors is usually sufficient; however, depending on the data, one can achieve a good approximation to the full GP with as few as five neighbors. If  $n$  is large, selecting fewer neighbors can result in substantial decrease in runtime.

Users can choose a slow brute force or fast code book nearest neighbor search algorithm by setting the `search.type` argument to `"brute"` or `"cb"`, respectively. The fast code book search is a modified version of the algorithm detailed in [Ra and Kim \(1993\)](#). If locations do not have identical coordinate values on the axis used for the ordering then `"cb"` and `"brute"` should produce identical neighbor sets. However, if there are identical coordinate values on the axis used for ordering, then the search algorithms might produce different, but equally valid, neighbor sets. If  $n$  is large (e.g., a million or more), constructing the nearest neighbor sets can take a long time even when `search.type = "cb"`. To save time, the neighbor set can be reused in subsequent model calls if the values passed to `coords`, `ord`, and `n.neighbors` do not change. Setting `return.neighbor.info = TRUE` in `spNNGP` or `spConjNNGP` returns the necessary neighbor information in an object called `neighbor.info`. Then passing this object to the optional `neighbor.info` argument in subsequent calls to `spNNGP` or `spConjNNGP` avoids the costly nearest neighbor search. The information in `neighbor.info` can also be used if one wishes to plot the neighbor sets (see example code in the `spNNGP` manual page).

The call to `spNNGP` in the code below generates 2000 MCMC samples using the latent model (with the `n.samples` argument specifying the desired number of samples). The `n.report` argument defines the sample interval for reporting the MH acceptance rate, which in this case is once every 1000 samples. The `fit.rep = TRUE` indicates that we are requesting regression fitted and replicate data be generated via composition sampling (i.e., one-for-one with each MCMC sample). The argument `sub.sample` specifies that we only want regression fitted and replicate samples starting at MCMC sample 1000 and for each subsequent sample, i.e., `start = 1000` (the list passed to `sub.sample` can also define `end` and `thin` for additional control over MCMC chain samples used in the computations).

The computer used to conduct this analysis has multicore processor and R compiled with **openMP**. Therefore, setting `n.omp.threads = 4` should decrease runtime, see Section 3.2 for more details on computing time.

```

R> n.samples <- 2000
R> starting <- list("phi" = 3 / 0.5, "sigma.sq" = 1, "tau.sq" = 1)
R> priors <- list("phi.Unif" = c(3 / 1, 3 / 0.1), "sigma.sq.IG" = c(2, 1),
+   "tau.sq.IG" = c(2, 1))
R> cov.model <- "exponential"
R> tuning <- list("phi" = 0.2)
R> ord <- order(coords[, 1] + coords[, 2])
R> sim.s <- spNNGP(formula = y ~ x, coords = coords, starting = starting,
+   tuning = tuning, priors = priors, cov.model = cov.model,
+   n.samples = n.samples, n.neighbors = 10, method = "latent", ord = ord,
+   n.omp.threads = 4, n.report = 1000, fit.rep = TRUE,
+   sub.sample = list(start = 1000), return.neighbor.info = TRUE)

## -----
##           Building the neighbor list
## -----
## -----
## Building the neighbors of neighbors list
## -----
## -----
##           Model description
## -----
## NNGP Latent model fit with 5000 observations.
##
## Number of covariates 2 (including intercept if specified).
##
## Using the exponential spatial correlation model.
##
## Using 10 nearest neighbors.
##
## Number of MCMC samples 2000.
##
## Priors and hyperpriors:
##     beta flat.
##     sigma.sq IG hyperpriors shape=2.00000 and scale=1.00000
##     tau.sq IG hyperpriors shape=2.00000 and scale=1.00000
##     phi Unif hyperpriors a=3.00000 and b=30.00000
##
## Source compiled with OpenMP support and model fit using 4 thread(s).
## -----
##           Sampling
## -----
## Sampled: 1000 of 2000, 50.00%
## Report interval Metrop. Acceptance rate: 46.50%
## Overall Metrop. Acceptance rate: 46.50%
## -----
## Sampled: 2000 of 2000, 100.00%

```

```
## Report interval Metrop. Acceptance rate: 33.70%
## Overall Metrop. Acceptance rate: 40.10%
## -----
```

As seen in the output above the call to `spNNGP` prints some basic model and sampler information. The output also notes Source compiled with `OpenMP` support and model fit using 4 thread(s). `spNNGP` functions will throw a warning if R was not compiled with `openMP` support and `n.omp.threads` is set to a value greater than 1.

Objects returned by `spNNGP` and other functions in the package use S3 methods `summary` and `print`. As illustrated later, S3 methods for `fitted` and `residuals` are also implemented. The `print` method prints the initial call to the function as well as some model and sampler specifics. The default behavior of `summary` is to print posterior summaries for model parameters using samples from the second half of the MCMC chains. The values for arguments `sub.sample` and `quantiles` passed to `summary` allow for finer control on posterior summaries. Alternatively the user can access all posterior samples as `coda` (Plummer, Best, Cowles, and Vines 2006) `mcmc` class objects in the `spNNGP` return objects' `p.beta.samples`, `p.theta.samples`, and `p.w.samples`, which as the names suggest hold  $\beta$ ,  $\theta$ , and  $\mathbf{w}$  samples, respectively.

```
R> summary(sim.s)

##
## Call:
## spNNGP(formula = y ~ x, coords = coords, method = "latent",
##   n.neighbors = 10, starting = starting, tuning = tuning,
##   priors = priors, cov.model = cov.model, n.samples = n.samples,
##   n.omp.threads = 4, ord = ord, return.neighbor.info = TRUE,
##   fit.rep = TRUE, sub.sample = list(start = 1000), n.report = 1000)
##
## Model class is NNGP, method latent, family gaussian.
##
## Model object contains 2000 MCMC samples.
##
## Chain sub.sample:
## start = 1000
## end = 2000
## thin = 1
## samples size = 1001
##           2.5%   25%   50%   75%   97.5%
## (Intercept) 0.5648 0.5948 0.6182 0.6757 0.7781
## x           -0.1066 -0.0966 -0.0912 -0.0852 -0.0756
## sigma.sq    0.8535 0.9446 1.0013 1.1486 1.6633
## tau.sq      0.2292 0.2378 0.2428 0.2486 0.2581
## phi         3.5656 4.7293 5.8478 6.2808 6.9279
```

As requested by setting `fit.rep = TRUE` the `spNNGP` return object also holds samples for the regression fitted values, labeled `y.hat.samples`, and replicated data, labeled `y.rep.samples`.

For convenience, the median and lower and upper 95% credible intervals for MCMC samples at each location are provided in `y.hat.quant`s and `y.rep.quant`s. These samples and corresponding summaries can be accessed directly in the `spNNGP` return object or extracted using the S3 `fitted` method. Beyond simply extracting the regression fitted values and replicated data from `spNNGP` and other model objects in the package, the `fitted` function performs additional composition sampling if the requested MCMC sample subset differs from the one initially specified in the model call. For example, the initial call to `spNNGP` specified `sub.sample = list(start = 1000)`, but if later we decide we want regression fitted values and replicated data for every 10th MCMC sample starting at sample 100, a call to `fitted(sim.s, sub.sample = list(start = 100, thin = 10))` would generate the desired subset (the `residuals` function provides the same behavior).

### *MCMC-free inference*

The conjugate model is called using the `spConjNNGP` function. Fixed  $\alpha$ ,  $\phi$ , and perhaps  $\nu$  are specified using a named vector passed to the `theta.alpha` argument. Alternatively, a  $K$ -fold cross-validation (where  $K$  is set via the `k-fold` argument) is used to discover the “optimal” set of parameters if `theta.alpha` is passed a matrix with columns named `alpha`, `phi`, and perhaps `nu`. The “optimal” set of parameter values (i.e., a row in `theta.alpha`) is the one that minimizes the average value of the specified scoring rule over the  $K$  folds. This scoring rule is set via the `score.rule` argument with options `"rmspe"` and `"crps"` for root mean squared prediction error (RMSPE) and continuous ranked probability score (CRPS, [Gneiting and Raftery 2007](#)). The  $K$ -fold cross-validation progress is printed to the screen as illustrated below. Once the optimal parameter set is identified, a final model is fit using all the available data. The description of this final model is given in the `Model description` section followed by the optimal set of  $\alpha$ ,  $\phi$ , and, if the Matérn correlation model is used,  $\nu$ .

The printout following the call to `spConjNNGP` below also includes a section called `Computing replicates` that reports on the exact sampling from the model parameters and regression fitted values posterior distributions, and generation of replicated data. Posterior sampling and generation of replicated data is optional and controlled by the `fit.rep` and `n.samples` arguments, with `n.samples` being set to the number of desired samples to collect. When `fit.rep = TRUE`, `spConjNNGP` effectively calls the S3 method `fitted` function for the conjugate model class. Hence, if posterior samples are not collected in the initial call to `spConjNNGP` or a different number of samples is needed, then a call to `fitted` using the `spConjNNGP` object will generate the required samples.

```
R> theta.alpha <- as.matrix(expand.grid(seq(0.01, 1, length.out = 15),
+   seq(3, 30, length.out = 15)))
R> colnames(theta.alpha) <- c("alpha", "phi")
R> sim.c <- spConjNNGP(y ~ x, coords = as.matrix(coords),
+   cov.model = "exponential", sigma.sq.IG = c(2, 0.5 * var(y)),
+   n.neighbors = 15, ord = ord, theta.alpha = theta.alpha,
+   k.fold = 2, score.rule = "rmspe", fit.rep = TRUE, n.samples = 200,
+   n.omp.threads = 4)

## -----
##           Starting k-fold
```

```

## -----
##
|
|                                     | 0%
|
| *****                            | 50%
|
| *****                            | 100%
## -----
##           Model description
## -----
## NNGP Conjugate model fit with 5000 observations.
##
## Number of covariates 2 (including intercept if specified).
##
## Using the exponential spatial correlation model.
##
## Using 15 nearest neighbors.
##
## Source compiled with OpenMP support and model fit using 4 thread(s).
## -----
## Priors and hyperpriors:
##           beta flat.
##           sigma.sq IG hyperpriors shape=2.00000 and scale=0.56696
## -----
##           Estimation for parameter set(s)
## Set phi=4.92857 and alpha=0.22214
## -----
##           Computing replicates
## -----
## NNGP Response model fit with 5000 observations.
##
## Number of covariates 2 (including intercept if specified).
##
## Using the exponential spatial correlation model.
##
## Using 15 nearest neighbors.
##
## Number of MCMC samples 200.
##
## Source compiled with OpenMP support and model fit using 4 thread(s).
## -----
##           Sampling
## Sampled: 100 of 200, 50.00%
## Sampled: 200 of 200, 100.00%

```

Results from the  $K$ -fold cross-validation are returned by `spConjNNGP` and held in the `k.fold.scores`

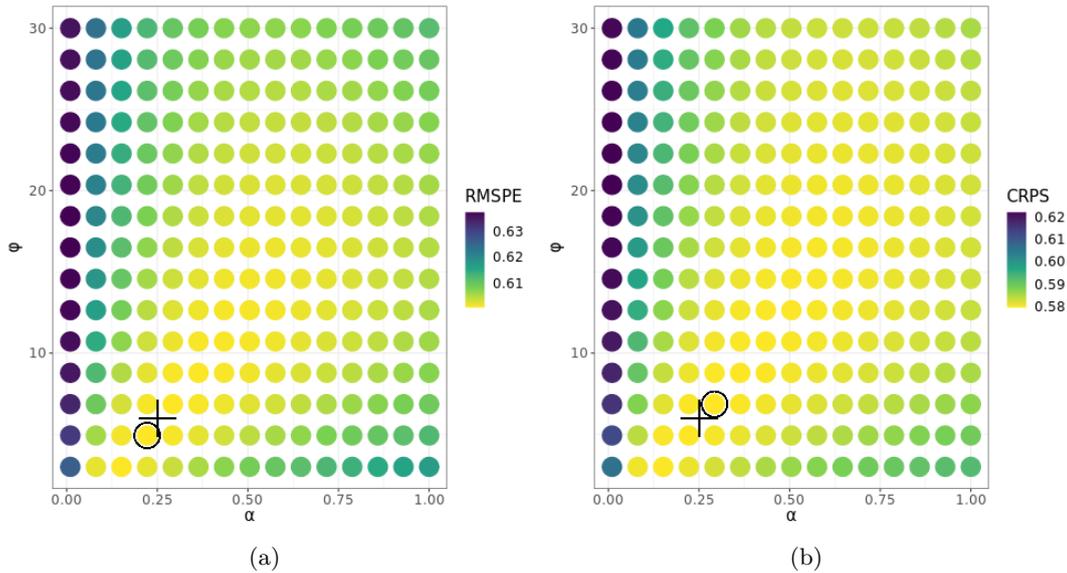


Figure 2: Simulated data analysis `spConjNNGP` parameter search grid and  $K$ -fold cross-validation results for RMSPE (a) and CRPS (b) scoring rule. The “optimal” parameter combination is circled. A plus symbol identifies the “true”  $\alpha$  and  $\phi$  used to generate the data.

	True	Model type			
		Full GP	Latent	Response	Conjugate
$\beta_0$	1	0.69 (0.24, 1.09)	0.62 (0.56, 0.78)	0.7 (0.27, 1.1)	0.62 (0.12, 1.19)
$\beta_x$	-0.1	-0.09 (-0.11, -0.07)	-0.09 (-0.11, -0.08)	-0.09 (-0.11, -0.08)	-0.09 (-0.11, -0.08)
$\sigma^2$	1	0.78 (0.67, 0.98)	1 (0.85, 1.66)	0.98 (0.72, 1.32)	1.11 (1.06, 1.15)
$\phi$	6	7.58 (5.93, 9.19)	5.85 (3.57, 6.93)	5.95 (4.53, 8.16)	4.93
$\tau^2$	0.25	0.24 (0.23, 0.26)	0.24 (0.23, 0.26)	0.24 (0.23, 0.26)	0.25 (0.24, 0.26)

Table 1: Simulated data analysis posterior summaries of median and 95% credible interval from the three model types.

matrix which is a copy of `theta.alpha` with additional columns for  $K$ -fold RMSPE and CRPS as illustrated below.

```
R> head(sim.c$k.fold.scores, n = 3)
```

```
##      phi      alpha      rmspe      crps
## [1,]   3 0.01000000 0.6262589 0.3567326
## [2,]   3 0.08071429 0.6019226 0.3399307
## [3,]   3 0.15142857 0.6005768 0.3387949
```

`k.fold.scores` is useful for assessing predictive performance sensitivity to choice of covariance parameters. Figure 2(a) was created by plotting the search grid parameter values

and their resulting minimum RMSPE. Note, the call to `spConjNNGP` specifies `score.rule = "rmspe"` which means the “optimal”  $\alpha$ ,  $\phi$ , and, if `cov.model = "matern"`,  $\nu$  will be the set that minimizes RMSPE. If one sets `score.rule = "crps"` the `k.fold.scores` can be used to identify the set of covariance parameters that minimize CRPS (the result of which is illustrated in Figure 2(b)).

The S3 `summary` method provides parameter point estimates using the optimal set and, if `n.samples` is specified, posterior summaries. The `sub.sample` argument can be used in `summary` if the `spConjNNGP fit.rep = FALSE` or if you wish for a posterior summary for a different number of `n.samples`. If `sub.sample` is specified in `summary`, the function returns a matrix of the requested number of samples. The output from `summary(sim.c)` which uses the initial number of 200 samples is given later in Table 1.

### *Model diagnostics and prediction*

When passed a `spNNGP` or `spCongNNGP` object the `spDiag` function returns a list that, depending on the model `method`, includes some or all of the following elements:

- DIC** a data frame holding the Deviance information criterion (DIC) and associated values defined by Spiegelhalter, Best, Carlin, and Van Der Linde (2002). The DIC data frame includes rows labeled DIC the criterion (lower is better), D a goodness of fit, and pD the effective number of parameters.
- WAIC** a data frame holding Watanabe-Akaike information criteria (WAIC) and associated values. The WAIC data frame includes rows labeled LPPD log pointwise predictive density, P.1 penalty term defined in unnumbered equation above Equation 11 in Gelman *et al.* (2013), P.2 an alternative penalty term defined in Equation 11, and the criteria WAIC.1 and WAIC.2 (lower is better) computed using P.1 and P.2, respectively.
- GPD** a data frame holding the values needed to compute the predictive criterion  $D = G + P$  defined by Gelfand and Ghosh (1998). The GPD data frame includes rows labeled G a goodness of fit, P a penalty term, and D the criterion (lower is better).
- GRS** a scoring rule, see Equation 27 in Gneiting and Raftery (2007) for details, with larger values of GRS indicating better model fit.

Among the four model comparison metrics, DIC and WAIC rely on the assumption that given all the parameters (including latent ones), the data points are conditionally independent. The response and conjugate NNGP models do not preserve this conditional independence structure and do not provide samples from the latent effect. Hence, it is not appropriate to compute WAIC and DIC for them. Comparisons across the models using GPD and GRS scores require generating replicate data. We have discussed in Section 2.4 how replicates have fundamentally different interpretation for the latent and response models. For the former, the replicates are generated conditional on the latent random effects and hence are spatially correlated with the original data, whereas for the marginalized response model, the replicate is simply a new realization of a multivariate Gaussian distribution with the same mean and covariance structure as the original data. Hence, generally the GPD and GRS scores will be better for the latent models (as is evident in Table 2). It is not advisable to compare the latent and response NNGP models using GPD and GRS as they represent different principles of replication. Finally, we can compare the response model with the conjugate model using

	Model			
	Full GP	Latent	Response	Conjugate
WAIC.1	8036.12	8034.85	–	–
WAIC.2	8501.09	8499.76	–	–
P.1	915.32	917.91	–	–
P.2	1147.8	1150.36	–	–
LPPD.2	–3102.75	–3099.52	–	–
DIC	8439.63	8438.64	–	–
pD	1318.82	1321.7	–	–
L	–2900.99	–2897.62	–	–
G	897.32	894.85	5645.59	5912.89
P	1539.07	1539.525	5234.422	6005.025
D	2436.39	2434.377	10880.01	11917.92
GRS	2975.25	2985.82	–5662.55	–5918.45

Table 2: Simulated data analysis output from calls to `spDiag` for the three model types.

GPD and GRS as they both use the same form of replicates. However, the conjugate model uses cross-validation to tune hyper-parameters, violating the principles of all these model comparison metrics tailored for classical Bayesian procedures, and it is difficult to interpret the model comparison values for it.

The code below calls `spDiag` for the `spNNGP` latent and `spConjNNGP` model output `sim.s` and `sim.c`, respectively. Additionally, for comparison, we ran `spNNGP` for the response model (i.e., by setting `method = "response"`) and called the resulting object `sim.r`. Fit diagnostics for all three models are given in Table 2. However, as shown in Table 1, all models recover the “true” parameter values well, and Figure 4 shows all models produce comparable predictive surfaces.

```
R> s.diag <- spDiag(sim.s)
R> r.diag <- spDiag(sim.r)
R> c.diag <- spDiag(sim.c)
```

Given the discordance in the interpretation of the traditional model comparison metrics for the NNGP models, a pragmatic way to compare them is based on their predictive performance on a hold out set. The `predict` function is used to generate posterior predictive samples for new locations with associated covariates, given a `spNNGP` or `spConjNNGP` object. The code below generates posterior predictive samples for the  $n_0 = 2500$  holdout locations using the latent model. The latent model is the only `method` that provides posterior predictive samples for the latent effect `w`. These samples are held in the `p.w.0` matrix in the `s.pred` object generated below, and a summary of these samples along with the “true” `w` are given in Figure 3. The `predict` function was also called for the response `sim.r` and conjugate `sim.c` model objects to generate posterior predictive samples for the holdout locations (held in output object `p.y.0`) along with subsequent surface summaries in Figure 4.

```
R> s.pred <- predict(sim.s, X.0 = cbind(1, x.ho), coords.0 = coords.ho,
+   sub.sample = list(start = 1000, thin = 10),
+   n.omp.threads = 4, n.report = 1000)
```

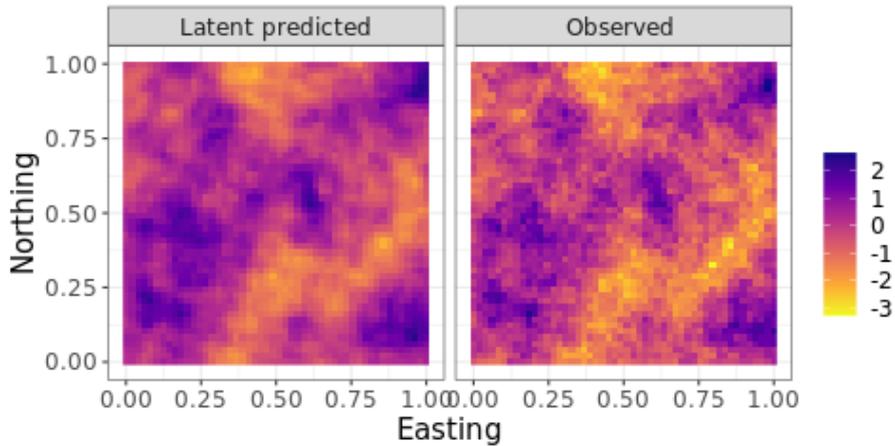


Figure 3: Observed and predicted simulated data  $w$  over the grid of holdout locations.

```
## -----
##           Prediction description
## -----
## NNGP Latent model fit with 5000 observations.
##
## Number of covariates 2 (including intercept if specified).
##
## Using the exponential spatial correlation model.
##
## Using 10 nearest neighbors.
##
## Number of MCMC samples 101.
##
## Predicting at 2500 locations.
##
##
## Source compiled with OpenMP support and model fit using 4 threads.
## -----
##           Predicting
## -----
## Location: 1000 of 2500, 40.00%
## Location: 2000 of 2500, 80.00%
## Location: 2500 of 2500, 100.00%
```

### 3.2. Timing given $n$ and number of cores

Here we provide a brief overview of the relationship between  $n$ , model type, and number of cores. The computer used for these runtime experiments (and analysis in subsequent sections) is running a linux operating system with a AMD Ryzen Threadripper 3990X 64-Core Processor (128 threads) and R compiled with **openMP** with thread-enabled MKL 2019.5.281 build with `MKL_DEBUG_CPU_TYPE=5` (Intel Corp. 2019). Timings generated by `proc.time()` is returned

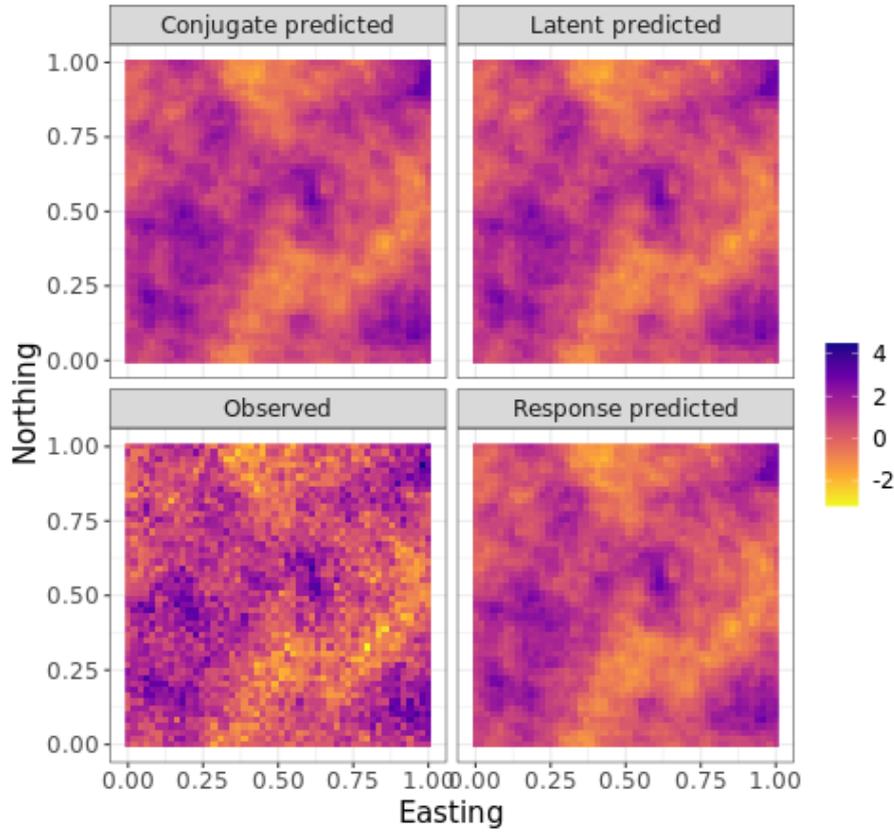


Figure 4: Observed and predicted simulated data  $\mathbf{y}$  over the grid of holdout locations.

by core `spNNGP` functions in the `run.time` vector. Figure 5(a) provides a sense of runtime (i.e., “wall time”) needed to collect 1000 MCMC samples for  $n = 100000$  using the response and latent algorithms for a range of cores. Execution time for the conjugate model is about equal to one MCMC iteration of the response model. As this figure shows, for this computer and  $n$ , there is little speed-up beyond  $\sim 40$  cores mostly due to communication overhead. Then fixing the number of cores at 40, Figure 5(b) gives a sense of computing time required for different size  $n$ .

### 3.3. Analysis of forest canopy height

In this section we analyze a forest canopy height dataset at  $n = 188,717$  locations using `spNNGP`. Digital maps of forest structure are key inputs to many ecosystem and Earth system modeling efforts (Finney 2004; Hurtt *et al.* 2004; Stratton 2006; Lefsky 2010; Klein, Randin, and Korner 2015). These and similar applications seek inference about forest canopy height variables and predictions that can be propagated through computer models of ecosystem function to yield more robust error quantification. Given the scientific and applied interest in forest structure, there is increasing demand for wall-to-wall (i.e., complete domain coverage) forest canopy height data at national and biome scales. Next generation LiDAR systems capable of large-scale mapping of forest canopy characteristics, such as ICESat-2 (Abdalati *et al.* 2010; ICESat-2 2015), Global Ecosystem Dynamics Investigation LiDAR (GEDI 2014), and NASA Goddard’s LiDAR, Hyperspectral, and Thermal (G-LiHT) Airborne Imager (Cook

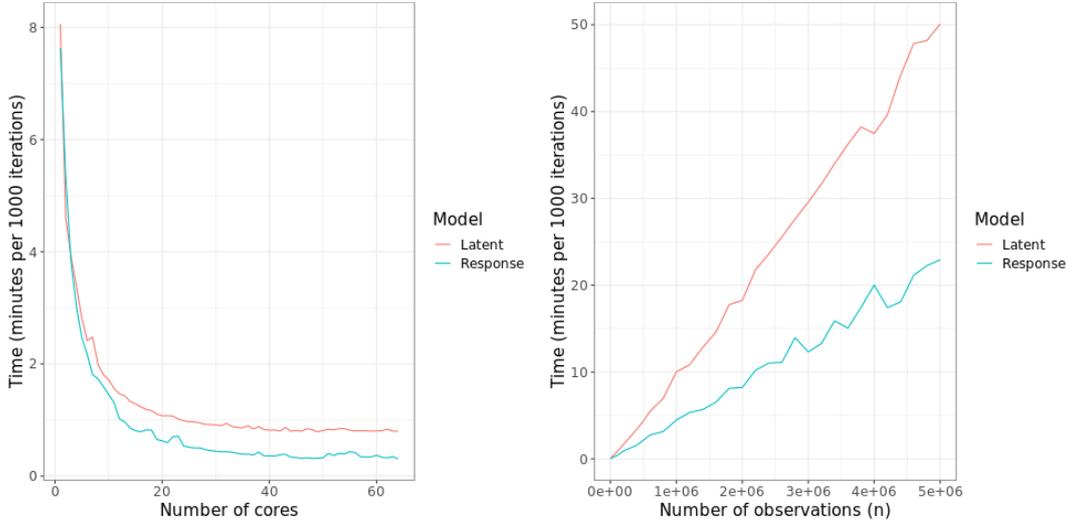


Figure 5: (a) Runtime for 1000 MCMC iterations for  $n = 100000$  and different number of cores. (b) Runtime for 1000 MCMC iterations using 40 cores and  $n$  from 1000 to 5 million.

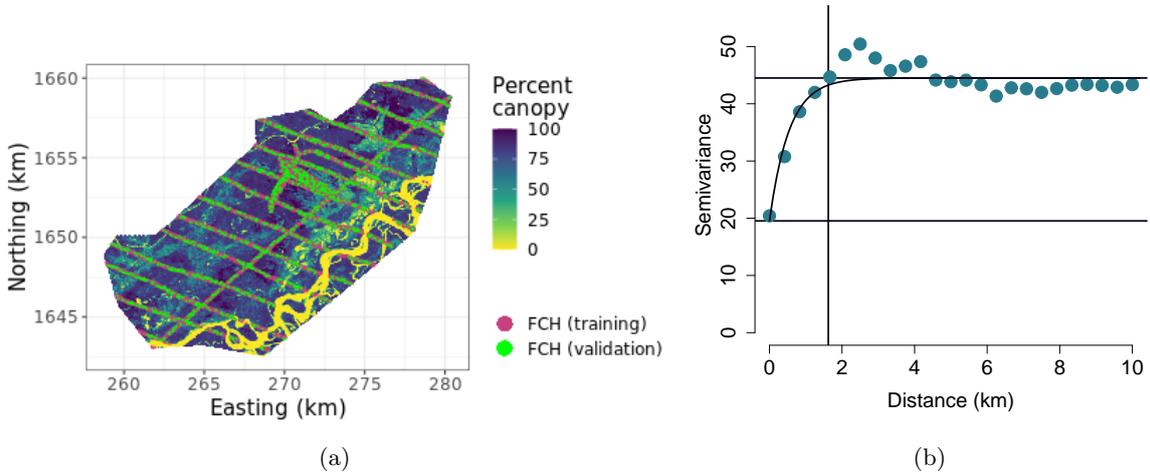


Figure 6: (a) Bonanza Creek experimental forest (BCEF) with G-LiHT LiDAR forest canopy height (FCH) estimates for model training and validation. The underlying map is the percent tree cover (PTC). (b) Semivariogram of BCEF non-spatial regression model residuals. Exponential covariance function estimate denoted by the curved line with associated estimates for  $\tau^2$ ,  $\sigma^2$ , and the effective spatial range are given by the lower horizontal, upper horizontal, and vertical lines, respectively.

*et al.* 2013), sample forest features using LiDAR instruments in long transects or cluster designs (see, e.g., the strips of LiDAR in Figure 6(a)). These next generation systems yield LiDAR data over the desired large spatial extents; however, the sparseness of the LiDAR sampling designs means prediction is required to deliver the desired wall-to-wall data products. Our goal is to create high spatial resolution forest canopy height predictions, with accom-

panying uncertainty estimates for the Bonanza Creek experimental forest (BCEF; <https://www.lter.uaf.edu>) located in interior Alaska, USA. The BCEF domain delineated for this study, Figure 6(a), is  $\sim 21,000$  ha and includes a section of the Tanana River floodplain along the southeastern border. The BCEF is a mixture of non-forest and forest vegetation featuring white spruce, black spruce, tamarack, quaking aspen, and balsam poplar trees mixed with willow and alder shrubland species Bonanza Creek LTER (2019). Figure 6(a) also shows location of the  $n = 188,717$  G-LiHT LiDAR forest canopy height (FCH) estimates. These data are included in the **spNNGP** package.

The 188,717 FCH estimates come from the G-LiHT LiDAR point cloud summarized to a  $13 \times 13$  m grid cell size (G-LiHT 2019). Over each grid cell, the maximum canopy height (i.e., FCH) was estimated using the 100th percentile height of the point cloud. A Landsat derived percent tree cover (PTC) data product developed by Hansen, Potapov, Moore, Hancher, Turubanova, Tyukavina, Thau, Stehman, Goetz, Loveland, Kommareddy, Egorov, Chini, Justice, and Townshend (2013), shown as the underlying surface in Figure 6(a) is used as a predictor variable for FCH. PTC is the percent tree cover estimates for peak growing season in 2010 and was created using a regression tree model applied to Landsat 7 ETM+ annual composites.

A semivariogram of the non-spatial regression model residuals can inform how the residual spatial/non-spatial variance (i.e., outcome variance not explained by the regression mean) is partitioned and the spatial range, see, e.g., Chapter 5 in Banerjee, Carlin, and Gelfand (2014) for details. Here we consider the residuals from

$$\mathbf{y} = \beta_0 + \beta_{PTC}\mathbf{x} + \boldsymbol{\varepsilon}, \quad (12)$$

where  $\mathbf{y}$  is the vector of observed FCH estimates,  $\beta_0$  is an intercept,  $\beta_{PTC}$  is the slope coefficient associated with the PTC predictor variable denoted as  $\mathbf{x}$ , and  $\boldsymbol{\varepsilon}$  is the  $n \times 1$  vector following  $N(\mathbf{0}, \tau^2 \mathbf{I}_n)$ . In the subsequent analyses we use an exponential spatial correlation function that approaches zero as the distance between locations increases. Therefore we define the distance,  $d_0$ , at which this correlation drops to 0.05 as the “effective spatial range”, which allows us to solve  $\phi = -\log(0.05)/d_0 \approx 3/d_0$ . Using the **variog** and **variofit** functions in the **geoR** package (Ribeiro Jr and Diggle 2020), the semivarogram and empirical parameter estimates for the BCEF are given in Figure 6(b). Due to computational constraints we used a random subset of 25,000 residuals from (12) to generate the variogram.

### *Estimation and prediction*

Here we consider the non-spatial, latent, response, and conjugate models for BCEF data. Posterior samples for the non-spatial regression model were generated using the **bayesLMRef** function in **spBayes**. Models are assessed using output from **spDiag**. Further, out-of-sample predictive performance was assessed by fitting the models to 100,000 observations (selected at random from the 188,717) and then predicting for the remaining holdout 88,717 observations. Finally, the models are used to predict FCH for a grid of 237,617 locations over the BCEF where PTC was recorded and resulting maps are compared.

```
R> n.samples <- 5000
R> starting <- list("phi" = 3 / 2, "sigma.sq" = 40, "tau.sq" = 1)
R> priors <- list("phi.Unif" = c(3 / 10, 3 / 0.1), "sigma.sq.IG" = c(2, 40),
+   "tau.sq.IG" = c(2, 10))
R> cov.model <- "exponential"
```

```

R> tuning <- list("phi" = 0.02)
R> bcef.s <- spNNGP(FCH ~ PTC, coords = c("x", "y"), data = BCEF.mod,
+   starting = starting, method = "latent", n.neighbors = 10,
+   tuning = tuning, priors = priors, cov.model = cov.model,
+   n.samples = n.samples, n.omp.threads = 40, n.report = 2500,
+   fit.rep = TRUE, sub.sample = list(start = 4000, thin = 10))

## -----
##           Building the neighbor list
## -----
## -----
## Building the neighbors of neighbors list
## -----
## -----
##           Model description
## -----
## NNGP Latent model fit with 100000 observations.
##
## Number of covariates 2 (including intercept if specified).
##
## Using the exponential spatial correlation model.
##
## Using 10 nearest neighbors.
##
## Number of MCMC samples 5000.
##
## Priors and hyperpriors:
##     beta flat.
##     sigma.sq IG hyperpriors shape=2.00000 and scale=40.00000
##     tau.sq IG hyperpriors shape=2.00000 and scale=10.00000
##     phi Unif hyperpriors a=0.30000 and b=30.00000
##
## Source compiled with OpenMP support and model fit using 40 thread(s).
## -----
##           Sampling
## -----
## Sampled: 2500 of 5000, 50.00%
## Report interval Metrop. Acceptance rate: 35.12%
## Overall Metrop. Acceptance rate: 35.12%
## -----
## Sampled: 5000 of 5000, 100.00%
## Report interval Metrop. Acceptance rate: 33.36%
## Overall Metrop. Acceptance rate: 34.24%
## -----

```

For brevity, the output from subsequent calls to `spNNGP` for the response model and `spConjNNGP` are suppressed by setting `verbose = FALSE`.

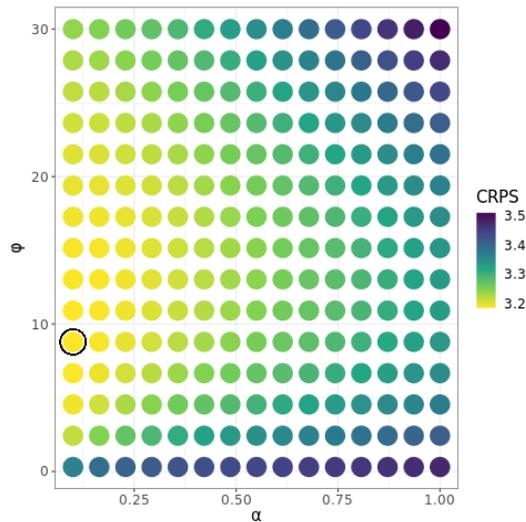


Figure 7: BCEF data analysis `spConjNNGP` parameter search grid and  $K$ -fold cross-validation results using CRPS scoring rules. The “optimal” parameter combination is circled.

	Model			
	Non-Spatial	Latent	Response	Conjugate
$\beta_0$	1.36 (1.18, 1.49)	11.09 (10.64, 11.92)	9.94 (9.41, 10.47)	10.45 (10.07, 10.85)
$\beta_{PTC}$	0.2 (0.2, 0.2)	0.03 (0.03, 0.03)	0.03 (0.03, 0.03)	0.03 (0.03, 0.04)
$\sigma^2$	–	43.48 (41.71, 45.68)	42.55 (40.6, 45.02)	34.66 (34.33, 34.89)
$\phi$	–	6.91 (6.53, 7.22)	7.19 (6.75, 7.52)	4.93
$\tau^2$	43.39 (43.12, 43.72)	3.46 (3.35, 3.58)	3.37 (3.27, 3.47)	3.47 (3.43, 3.49)

Table 3: BCEF data analysis posterior summaries of median and 95% credible interval from the three model types.

```
R> tuning <- list("phi" = 0.01, "sigma.sq" = 0.01, "tau.sq" = 0.005)
R> bcef.r <- spNNGP(FCH ~ PTC, coords = c("x", "y"), data = BCEF.mod,
+   starting = starting, method = "response", n.neighbors = 10,
+   tuning = tuning, priors = priors, cov.model = cov.model,
+   n.samples = n.samples, n.omp.threads = 40, n.report = 2500,
+   fit.rep = TRUE, sub.sample = list(start = 4000, thin = 10),
+   verbose = FALSE)
R> theta.alpha <- as.matrix(expand.grid(seq(0.1, 1, length.out = 15),
+   seq(3 / 10, 3 / 0.1, length.out = 15)))
R> colnames(theta.alpha) <- c("alpha", "phi")
R> bcef.c <- spConjNNGP(FCH ~ PTC, coords = c("x", "y"), data = BCEF.mod,
+   cov.model = "exponential", sigma.sq.IG = c(2, 40), n.neighbors = 10,
+   theta.alpha = theta.alpha, k.fold = 2, score.rule = "crps",
+   fit.rep = TRUE, n.samples = 200, n.omp.threads = 40, verbose = FALSE)
```

Posterior summaries are given in Table 3. As suggested by the exploratory variogram analysis, and now confirmed with formal model estimates, the spatial range is quite short, e.g., the

	Model type			
	Non-Spatial	Latent	Response	Conjugate
WAIC.1	660845	435074.6	–	–
WAIC.2	660845	462473.1	–	–
P.1	2.35	27501.35	–	–
P.2	2.38	41200.63	–	–
LPPD.2	–330420.1	–190035.9	–	–
DIC	660845.3	464395.5	–	–
pD	2.69	56822.25	–	–
L	–330420	–175375.5	–	–
G	4387671	152329.7	7086265	6403563
P	4336921	542156.8	4119386	3502314
D	8724591	694486.5	11205651	9905877
GRS	–479230.7	–196616.3	–545860.2	–539539.6
CRPS	3.79	1.57	1.9	1.53
RMSPE	6.62	3.01	3.39	2.94
CI Cover	94.42	85.8	72.56	–
CI Width	24.73	7.82	5.76	–

Table 4: BCEF data analysis model fit via `spDiag` and out-of-sampled prediction diagnostics for the non-spatial and three spatial model types fit to the BCEF data. The last four rows were calculated using prediction for the holdout set. The row labeled CI Cover is the percent of 95% posterior predictive distribution credible intervals that cover the observed holdout value. The row labeled CI Width is the average width of the 95% posterior predictive credible interval.

latent model estimate of the median effective spatial range is  $\sim 0.43$  km (i.e.,  $-\log(0.05)/6.91$ ). Despite this short range, the spatial variance is large relative to the non-spatial variance. Such results are not surprising given the BCEF’s composition and structure are the result of myriad large and small spatial scale biotic (e.g., insect disturbance) and abiotic (e.g., soil, topography, climate, wind, fire) factors that cause spatially complex mortality and regrowth patterns. Formal model fit diagnostics provided in Table 4 suggests the addition of the latent spatial effect does improve fit compared with the non-spatial model.

Lastly, and as illustrated in the synthetic data analysis, a call to `predict` yields posterior predictive samples for the entire domain of interest which in this case is sampling over the grid of  $n_0 = 237617$  locations where only PTC was recorded. A surface of the posterior distributions’ mean and variance are given in Figures 8(a) and 8(b), respectively. Locations where observations are available to inform prediction (i.e., along flight lines) are clearly delineated by high prediction precision in the prediction variance surfaces Figure 8(b). Given the relatively short spatial range this information borrowing to inform prediction does not extend too far off of the flight lines. These surfaces along with the out-of-sampled prediction performance metrics given in the last four rows in Table 4 suggest there is not too much difference among the spatial models.

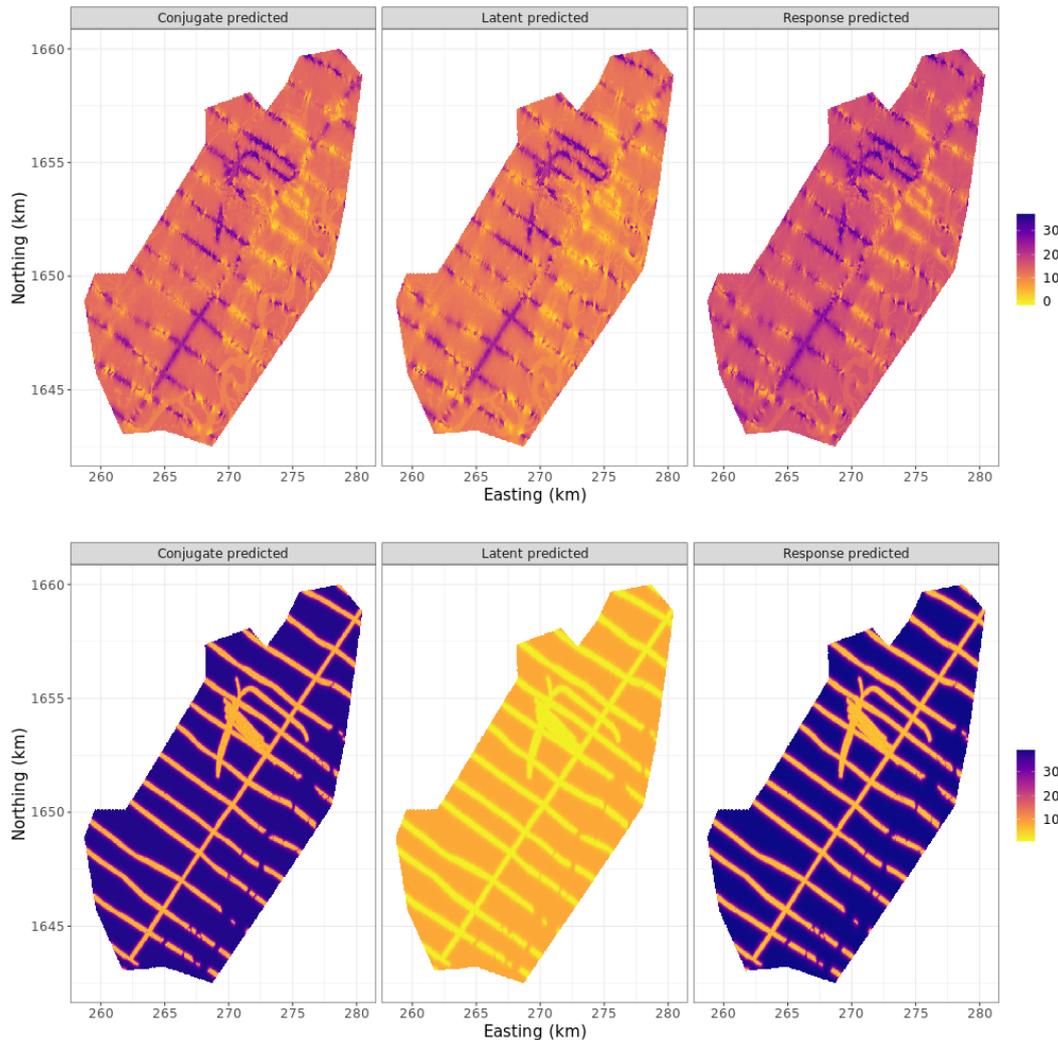


Figure 8: BCEF data analysis posterior predictive distribution mean (a) and variance (b).

### 3.4. Analysis of species distributions

In this section, we present analysis of species distribution using the binomial NNGP model. Species distribution models (SDMs) project the outcome of community assembly processes dispersal, the abiotic environment, and biotic factors onto geographic space (Guisan and Zimmermann 2000; Pulliam 2000). Here, we reanalyze data recently presented in Lany, Zarnetske, Finley, and McCullough (2020) to develop a SDM for eastern hemlock (*Tsuga canadensis* L.) coded as TSCA in subsequent analysis. The data comprise hemlock occurrence (binomial outcome) on 17,743 forest stands across Michigan, USA. A set of predictors were also observed at each stand and subsequently used to explain the probability of hemlock occurrence. Predictor variables included minimum winter temperature (MIN), maximum summer temperature (MAX), total precipitation in the coldest quarter of the year (WIP), total precipitation in the warmest quarter of the year (SUP), annual actual evapotranspiration (AET) and annual climatic water deficit (DEF).

We consider three candidate models: (1) non-spatial logistic regression using the Pólya-Gamma data-augmented sampler of Polson *et al.* (2013) as implemented in the our `spNNGP` `PGLogit` function; (2) logistic regression with a space-varying Gaussian Predictive Process (GPP) random effect Banerjee *et al.* (2008) using the `spBayes`' `spGLM` function; (3) logistic regression with space-varying NNGP again via the Pólya-Gamma sampler invoked using the `family = "binomial"` argument in the `spNNGP` function as illustrated in the code below. Computational benefits of GPP models arise from the estimation of the latent Gaussian process at set of locations referred to as knots, where the number of knots is typically much smaller than  $n$ . When applying a GPP the analysis must strike an acceptable balance between computing time, which is governed by the number of knots, and Gaussian process being approximated. More details about implementing a GPP can be found in Banerjee *et al.* (2008), Finley *et al.* (2009), Guhaniyogi, Finley, Banerjee, and Gelfand (2011), and Stein (2014). Here, for simplicity, we selected 25 knots on a fixed grid over the domain as it yielded a reasonable approximation to the desired GP with a computing time close that that of the NNGP. All three models were fit using  $n = 15,000$  observations and assessed using goodness of fit metrics and out-of-sampled predictive performance based on a holdout set of  $n_0 = 2,743$ . Like a `spNNGP` class object, the return object from the `PGLogit` function can be passed to `spDiag` to yield model diagnostics as illustrated in the code below. Both `PGLogit` and `spNNGP` with `faimly = "binomial"` accept different number of trials for each location, i.e.,  $n_i$  in Section 2.3, which is passed via the `weights` argument; however, for the current setting we accept the default of all weights equal to 1, i.e., each stand yields either presence or absence of hemlock.

```
R> n.samples <- 10000
R> starting <- list("phi" = 3 / 50, "sigma.sq" = 10)
R> tuning <- list("phi" = 0.05)
R> priors <- list("phi.Unif" = c(3 / 300, 3 / 10), "sigma.sq.IG" = c(2, 10))
R> cov.model <- "exponential"
R> m.s <- spNNGP(TSCA ~ MIN + MAX + SUP + WIP + AET + DEF,
+   coords = c("long", "lat"), data = mi.mod, family = "binomial",
+   method = "latent", n.neighbors = 10, starting = starting,
+   tuning = tuning, priors = priors, cov.model = cov.model,
+   n.samples = n.samples, n.report = 5000, fit.rep = TRUE,
+   sub.sample = list(start = 9000), n.omp.threads = 40)

## -----
##           Building the neighbor list
## -----
## -----
## Building the neighbors of neighbors list
## -----
## -----
##           Model description
## -----
## NNGP Latent model fit with 15000 observations.
##
## Number of covariates 7 (including intercept if specified).
```

	Model		
	Non-Spatial	GPP 25 knots	NNGP Latent
$\beta_0$	-2.65 (-2.72, -2.59)	-3.6 (-3.71, -3.45)	-4.77 (-5.08, -4.47)
$\beta_{MIN}$	0.51 (0.39, 0.64)	0.74 (0.56, 0.9)	0.13 (-0.22, 0.42)
$\beta_{MAX}$	-0.21 (-0.3, -0.12)	-0.12 (-0.26, 0.01)	-0.02 (-0.36, 0.21)
$\beta_{SUP}$	0.19 (0.11, 0.27)	-0.12 (-0.23, -0.03)	-0.16 (-0.44, 0.09)
$\beta_{WIP}$	-0.01 (-0.1, 0.07)	0.12 (-0.01, 0.25)	0.01 (-0.27, 0.32)
$\beta_{AET}$	-0.41 (-0.54, -0.27)	-0.32 (-0.49, -0.16)	-0.31 (-0.52, -0.1)
$\beta_{DEF}$	-0.44 (-0.53, -0.34)	-0.35 (-0.48, -0.23)	-0.28 (-0.42, -0.13)
$\sigma^2$	-	4.12 (2.33, 8.46)	8.45 (6.95, 9.79)
$\phi$	-	0.01 (0.01, 0.01)	0.08 (0.06, 0.1)

Table 5: Michigan Eastern Hemlock SDM analysis posterior summaries of median and 95% credible interval from the three candidate models.

```
## Using the exponential spatial correlation model.
##
## Using 10 nearest neighbors.
##
## Number of MCMC samples 10000.
##
## Priors and hyperpriors:
##     beta flat.
##     sigma.sq IG hyperpriors shape=2.00000 and scale=10.00000
##     phi Unif hyperpriors a=0.01000 and b=0.30000
##
## Source compiled with OpenMP support and model fit using 40 thread(s).
## -----
##                               Sampling
## -----
## Sampled: 5000 of 10000, 50.00%
## Report interval Metrop. Acceptance rate: 39.96%
## Overall Metrop. Acceptance rate: 39.96%
## -----
## Sampled: 10000 of 10000, 100.00%
## Report interval Metrop. Acceptance rate: 39.78%
## Overall Metrop. Acceptance rate: 39.87%
## -----
```

Parameter estimates for the three models are given in Table 5, and shows that a number of the predictor variables help explain the probability of hemlock occurrence across the study area. Due to possible spatial confounding (Hanks, Schliep, Hooten, and Hoeting 2015), we should interpret the sign and significance of these regression parameters in the spatial models with caution. Within sample fit diagnostics given in Table 6 and out-of-sample receiver operating characteristic (ROC) prediction curves (based on  $n_0 = 2,743$  holdout locations) given in Figure 9, suggest the latent spatial variables improve the species distribution models over that of the non-spatial model. Among the spatial models, the NNGP model outperforms

	Model		
	Non-Spatial	GPP 25 knots	NNGP Latent
WAIC.1	7494.56	6889.75	5252.05
WAIC.2	7494.6	6890.07	5511.65
P.1	6.8	27.92	623.02
P.2	6.82	28.07	752.82
LPPD.2	-3740.48	-3416.96	-2003
DIC	7494.58	6888.86	5410.72
pD	6.82	27.02	781.7

Table 6: Michigan Eastern Hemlock SDM analysis model fit via `spDiag` and out-of-sample prediction diagnostics for the non-spatial and two spatial models.

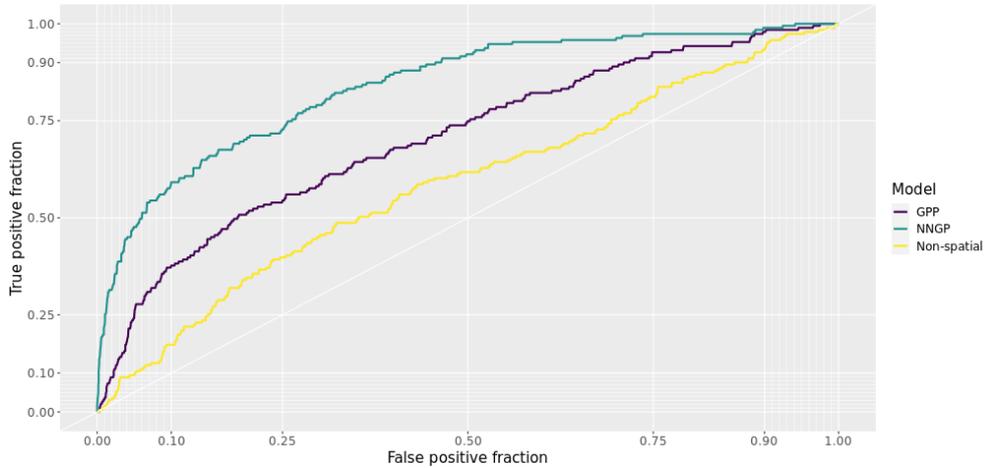


Figure 9: Michigan Eastern Hemlock SDM analysis out-of-sample posterior predictive mean ROC curves for the candidate models. Models with better prediction have curves closer to the top left corner.

the reduced rank predictive process model. The computing times per 1000 samples for the non-spatial, GPP, and NNGP models are 3, 13, and 9 seconds, respectively, using 10 cores.

### 3.5. Statistical gap-filling of a massive remotely sensed data

As reviewed in Section 1, methods and software are now emerging that can readily fit geostatistical models to data sets comprising locations in the 10s to 100s of thousands. However, as the number of observations climb into the millions, inferential and software options are quite limited. It is common to encounter data sets of such size in a variety of fields. For example, remotely sensed data of this magnitude either as gridded or scattered data products is now ubiquitous. Here we consider a massive-scale “gap-filling” exercise of missing normalized difference vegetation index (NDVI) – a measure of vegetation greenness — data from around 39 million locations. The NDVI data for a LandSat 8 sensor image was taken over Limpopo National Park, Mozambique, Africa on 2015-07-17. The image shown in Figure 10(a) has  $n_0 = 778,644$  pixels that are missing NDVI (denoted in gray) due to cloud cover during image acquisition. The red box in Figure 10(a) delineates a region with a large number of missing

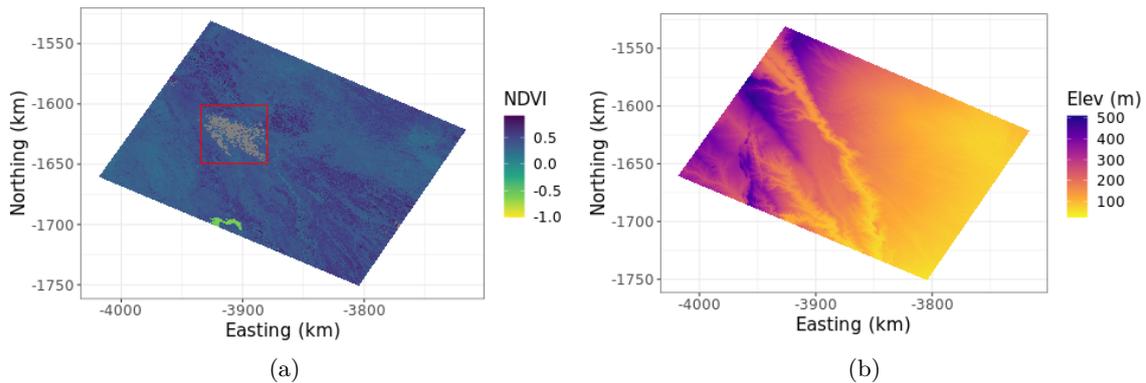


Figure 10: (a) Normalized Difference Vegetation Index (NDVI) LandSat 8 sensor image over Limpopo National Park, Mozambique, Africa on 2015-07-17, with gray pixels indicate missing NDVI data. (b) Land surface elevation used to help explain variation in NDVI.

pixels. Filling in missing NDVI values for this and other images in a time series over Limpopo was a step in a larger study conducted by [Desanker, Dahlin, and Finley \(2020\)](#) that looked at environmental drivers in vegetation phenology change. Our aim is to build a geostatistical model to predict the  $n_0$  missing NDVI pixels given the  $n = 38825052$  observed NDVI pixels and complete coverage land surface elevation predictor variable shown in Figure 10(b).

While a MCMC solution is feasible via a call to `spNNGP` the runtime would be on the order of days, even on a multiprocessor computer. Therefore we opt for a MCMC-free approach using `spConjNNGP` – trading some inference about the spatial process parameters for substantially decreased computing time. An initial variogram analysis using a sample of residuals from NDVI regressed on elevation helped define `spConjNNGP`'s `theta.alpha` search grid of 100  $\phi$  and  $\alpha$  combinations, and for setting the prior for  $\sigma^2$ . As noted in the code below, we used minimum CRPS to select the optimal set of covariance parameters (i.e., specified via the `score.rule` argument).

```
R> theta.alpha <- as.matrix(expand.grid(seq(1.5, 0.01, length.out = 10),
+   seq(3 / 200, 3 / 25, length.out = 10)))
R> colnames(theta.alpha) <- c("alpha", "phi")
R> ser.c <- spConjNNGP(ndvi ~ elev, coords = c("x", "y"), data = ser.mod,
+   cov.model = "exponential", sigma.sq.IG = c(2, 0.01),
+   n.neighbors = 10, theta.alpha = theta.alpha, k.fold = 2,
+   score.rule = "crps", X.0 = cbind(1, ser.ho$elev),
+   coords.0 = as.matrix(ser.ho[, c("x", "y")]), n.omp.threads = 40)
```

```
## -----
##           Starting k-fold
## -----
##
|
|                                     | 0%
|
| *****                          | 50%
```

```

|
| *****| 100%
## -----
##           Model description
## -----
## NNGP Conjugate model fit with 38825052 observations.
##
## Number of covariates 2 (including intercept if specified).
##
## Using the exponential spatial correlation model.
##
## Using 10 nearest neighbors.
##
## Source compiled with OpenMP support and model fit using 40 thread(s).
## -----
## Priors and hyperpriors:
##           beta flat.
##           sigma.sq IG hyperpriors shape=2.00000 and scale=0.01000
## -----
## Predicting at 778644 locations.
## -----
##           Estimation for parameter set(s)
## Set phi=0.12000 and alpha=0.01000

```

Given the size of this data set, the search time to identify nearest neighbor sets can take a considerable amount of time, even when using the [Ra and Kim \(1993\)](#) fast code book algorithm invoked by `search.type = "cb"`. In this case the search time was 8.86 minutes (`search.proc.time()` is held in `neighbor.info$nn.indx.run.time` in the model object). Once the optimal  $\phi$  and  $\alpha$  was found, the runtime for parameter estimation and prediction for the missing pixels was 31.13 minutes. We have experimented with a variety of tree-based data structures (e.g., *kd*-trees modified to accommodate the nearest neighbor search constraints) and associated search algorithms, and believe there are substantial gains in search time efficiency to be had with additional development.

As shown below, a subsequent call to `summary` returns the optimal parameter set as well as point estimates for the other model parameters. Given the large size of this data set, parameter variance estimates are remarkably small (for some parameters the variance might be smaller than machine precision). The message at the end of the `summary` output reminds us that posterior samples were not requested in the initial call to `spConjNNGP` or `summary`. Generating a reasonable number of posterior samples for parameters would take a few hours for a data set of this size.

```
R> summary(ser.c, digits = 8)
```

```
## Call:
## spConjNNGP(formula = ndvi ~ elev, data = ser.mod, coords = c("x", "y"),
##           n.neighbors = 10, theta.alpha = theta.alpha, sigma.sq.IG = c(2,
##           0.01), cov.model = "exponential", k.fold = 2, score.rule = "crps",

```

```

##      X.0 = cbind(1, ser.ho$elev), coords.0 = as.matrix(ser.ho[, c("x",
##              "y")]), n.omp.threads = 40)
##
## Model class is NNGP, method conjugate, family gaussian.
##
##              Estimate    Variance
## (Intercept) 0.14544817 0.00002334
## elev        0.00155903 0.00000000
## sigma.sq    0.06520584 0.00000000
## phi         0.12000000 0.00000000
## alpha       0.01000000 0.00000000

## If posterior summaries are desired, then either rerun spConjNNGP with
## fit.rep=TRUE, or specify the summary argument sub.sample to indicate the
## number of fitted and replicated samples to collect.

```

For prediction, as illustrated in the previous analyses, we can pass the `spConjNNGP` object `ser.c` to `predict`. Alternatively, a slightly more efficient option (which avoids computing observed location covariance with their neighbor set twice, i.e., once for estimating parameters in `spCongNNGP` and again in subsequent call to `predict`) is to specify the prediction locations and associated design matrix via `coords.0` and `X.0` in the initial call to `spConjNNGP`. In this case, `ser.c` includes mean (`y.0.hat`) and variance (`y.0.hat.var`) estimates for the prediction locations. These predicted mean and variance of the mean estimates are shown for the missing pixels delineated by the red box in Figure 10(a) (see Figure 11(a) for a zoomed in view) in Figures 11(b) and (c), respectively. As expected, the prediction variance surface shows that pixels close to observed pixels have higher precision due to borrowing of information through the spatial correlation structure.

## 4. Summary

The `spNNGP` R package provides a suite of NNGP-based (Datta *et al.* 2016) spatial regression models for both Gaussian and non-Gaussian point-referenced outcomes that are spatially indexed. The package implements the MCMC and MCMC-free algorithms detailed Finley *et al.* (2019) with the addition of the Pólya-Gamma latent variable model for binomial outcomes. Special care was taken to design algorithms that take advantage of multiprocessor computer via **OpenMP** and those with threaded **BLAS** and **LAPACK** libraries. Our future aim is to add functionality to accommodate multivariate outcomes, where we envisage two settings, first, where a limited number of outcomes (e.g., fewer than 10) might be handled using a NNGP linear model of coregionalization (Gelfand, Schmidt, Banerjee, and Sirmans 2004), second, where the number of outcomes is larger and requires some dimension reduction, e.g., via a NNGP spatial factor model, akin to the model detailed in Taylor-Rodriguez *et al.* (2019). For such highly multivariate data, besides factor models, we will also explore new multivariate covariance functions using sparse inter-variable graphical models that parsimoniously capture the relationship between multiple variables. Future releases will also provide the flexibility to specify a general space-varying coefficient model like the `spBayes spSVC` function (Finley and Banerjee 2020). Finally, our current functions for delivering exact Bayesian inference using conjugate models for the response will be expanded to accommodate latent

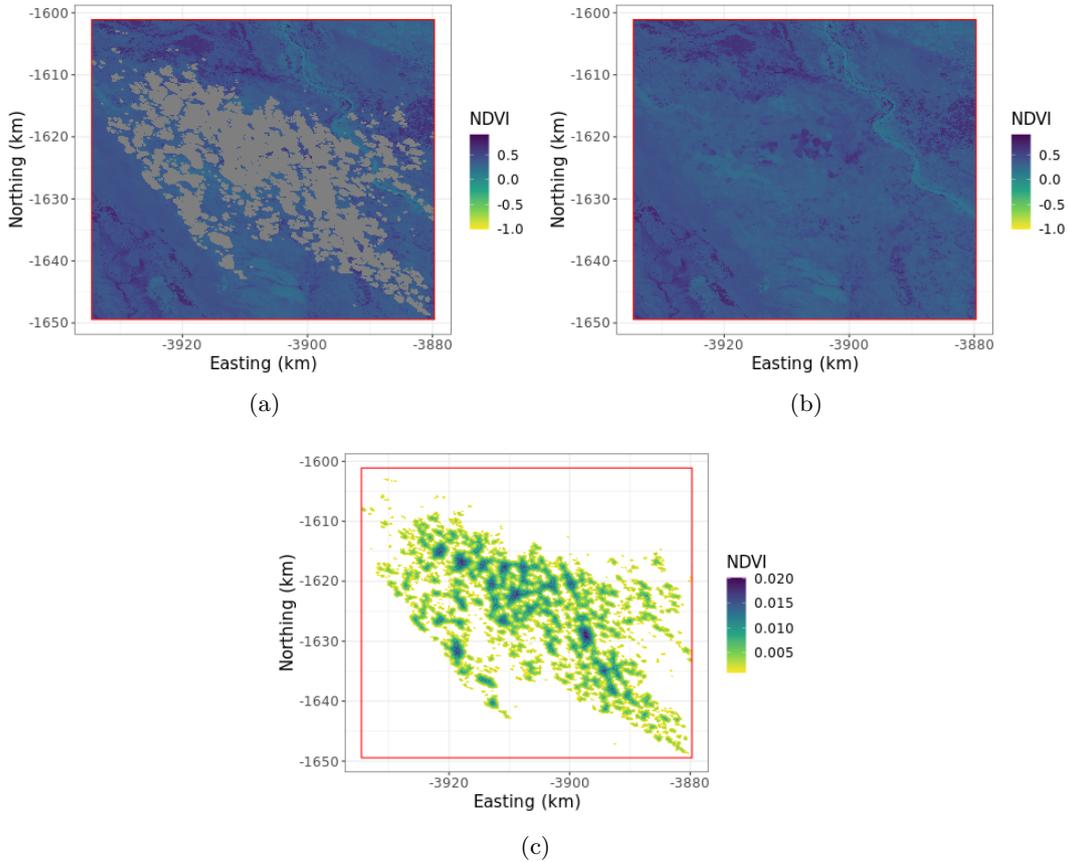


Figure 11: (a) NDVI image subset identified by the red box in Figure 10(a), with gray pixels indicate missing NDVI data. (b) Missing pixel posterior predictive distribution mean. (c) Missing pixel posterior predictive distribution variance.

spatial effects (Zhang, Banerjee, and Finley 2021) and multivariate spatial regression using the matrix-variate normal and inverse-Wishart families.

## Acknowledgments

Finley was supported by National Science Foundation (NSF) EF-1253225, DMS-1916395, DEB-2213565, and National Aeronautics and Space Administration’s Carbon Monitoring System project. Banerjee was supported by NSF DMS-1513654, DMS-2113778, and DMS-1916349. Datta was supported by NSF DMS-1915803. The authors thank Michele Peruzzi for fruitful conversations about computing.

## References

Abdalati W, Zwally HJ, Bindschadler R, Csatho B, Farrell SL, Fricker HA, Harding D, Kwok R, Lefsky M, Markus T, Marshak A, Neumann T, Palm S, Schutz B, Smith B, Spinhirne J, Webb C (2010). “The ICESat-2 Laser Altimetry Mission.” *Proceedings of the IEEE*, **98**(5), 735–751. doi:10.1109/jproc.2009.2034765.

- Abdulah S, Ltaief H, Sun Y, Genton MG, Keyes DE (2018). “**ExaGeoStat**: A High Performance Unified Software for Geostatistics on Manycore Systems.” *IEEE Transactions on Parallel and Distributed Systems*, **29**(12), 2771–2784. doi:10.1109/tpds.2018.2850749.
- Banerjee S (2017). “High-Dimensional Bayesian Geostatistics.” *Bayesian Analysis*, **12**, 583–614. doi:10.1214/17-ba1056r.
- Banerjee S, Carlin BP, Gelfand AE (2014). *Hierarchical Modeling and Analysis for Spatial Data*. Monographs on Statistics & Applied Probability, 2nd edition. Chapman & Hall/CRC.
- Banerjee S, Gelfand AE, Finley AO, Sang H (2008). “Gaussian Predictive Process Models for Large Spatial Data Sets.” *Journal of the Royal Statistical Society B*, **70**(4), 825–848. doi:10.1111/j.1467-9868.2008.00663.x.
- Barbian MH, Assunção RM (2017). “Spatial Subsemble Estimator for Large Geostatistical Data.” *Spatial Statistics*, **22**, 68–88. doi:10.1016/j.spasta.2017.08.004.
- Bivand R, Nowosad J (2022). *CRAN Task View: Analysis of Spatial Data*. Version 2022-04-14, URL <https://CRAN.R-project.org/view=Spatial>.
- Bonanza Creek LTER (2019). “Bonanza Creek LTER, Institute of Arctic Biology, University of Alaska Fairbanks.” Accessed: 2022-1-04, URL <http://www.lter.uaf.edu/research/study-sites-bcef>.
- Cook BD, Corp LA, Nelson RF, Middleton EM, Morton DC, McCorkel JT, Masek JG, Ranson KJ, Ly V, Montesano PM (2013). “NASA Goddard’s LiDAR, Hyperspectral and Thermal (G-LiHT) Airborne Imager.” *Remote Sensing*, **5**(8), 4045–4066. doi:10.3390/rs5084045.
- Cressie N, Johannesson G (2008). “Fixed Rank Kriging for Very Large Spatial Data Sets.” *Journal of the Royal Statistical Society B*, **70**, 209–226. doi:10.1111/j.1467-9868.2007.00633.x.
- Dagum L, Menon R (1998). “**OpenMP**: An Industry Standard API for Shared-Memory Programming.” *IEEE Computational Science & Engineering*, **5**(1), 46–55. doi:10.1109/99.660313.
- Datta A, Banerjee S, Finley AO, Gelfand AE (2016). “Hierarchical Nearest-Neighbor Gaussian Process Models for Large Geostatistical Datasets.” *Journal of the American Statistical Association*, **111**(514), 800–812. doi:10.1080/01621459.2015.1044091.
- Desanker G, Dahlin KM, Finley AO (2020). “Environmental Controls on Landsat-Derived Phenoregions Across an East African Megatransect.” *Ecosphere*, **11**(5), e03143. doi:10.1002/ecs2.3143.
- Emery X (2009). “The Kriging Update Equations and Their Application to the Selection of Neighboring Data.” *Computational Geosciences*, **13**(3), 269–280. doi:10.1007/s10596-008-9116-8.
- Finley AO, Banerjee S (2020). “Bayesian Spatially Varying Coefficient Models in the **spBayes** R Package.” *Environmental Modelling & Software*, **125**, 104608. doi:10.1016/j.envsoft.2019.104608.

- Finley AO, Banerjee S, Gelfand A (2015). “**SpBayes** for Large Univariate and Multivariate Point-Referenced Spatio-Temporal Data Models.” *Journal of Statistical Software*, **63**(13), 1–28. doi:10.18637/jss.v063.i13.
- Finley AO, Datta A, Banerjee S (2022). *spNNGP: Spatial Regression Models for Large Datasets using Nearest Neighbor Gaussian Processes*. R package version 1.0.0, URL <https://CRAN.R-project.org/package=spNNGP>.
- Finley AO, Datta A, Cook BD, Morton DC, Andersen HE, Banerjee S (2019). “Efficient Algorithms for Bayesian Nearest Neighbor Gaussian Processes.” *Journal of Computational and Graphical Statistics*, **28**, 401–414. doi:10.1080/10618600.2018.1537924.
- Finley AO, Sang H, Banerjee S, Gelfand AE (2009). “Improving the Performance of Predictive Process Modeling for Large Datasets.” *Computational Statistics & Data Analysis*, **53**(8), 2873–2884. doi:10.1016/j.csda.2008.09.008.
- Finney MA (2004). “FARSITE: Fire Area Simulator – Model Development and Evaluation.” *Research Paper RMRS-RP-4*, U.S. Department of Agriculture, Forest Service, Rocky Mountain Research Station.
- Furrer R, Flury R, Gerber F (2021). *spam: SParse Matrix*. R package version 2.7-0, URL <https://CRAN.R-project.org/package=spam>.
- Furrer R, Sain SR (2010). “**spam**: A Sparse Matrix R Package with Emphasis on MCMC Methods for Gaussian Markov Random Fields.” *Journal of Statistical Software*, **36**(10), 1–25. doi:10.18637/jss.v036.i10.
- G-LiHT (2019). “G-LiHT: Goddard’s LiDAR, Hyperspectral & Thermal Imager.” Accessed: 2022-11-05, URL <https://gliht.gsfc.nasa.gov/>.
- GEDI (2014). “Global Ecosystem Dynamics Investigation LiDAR.” Accessed: 2022-01-04, URL <https://gedi.umd.edu/>.
- Gelfand AE, Ghosh SK (1998). “Model Choice: A Minimum Posterior Predictive Loss Approach.” *Biometrika*, **85**(1), 1–11. doi:10.1093/biomet/85.1.1.
- Gelfand AE, Schmidt AM, Banerjee S, Sirmans CF (2004). “Nonstationary Multivariate Process Modeling through Spatially Varying Coregionalization.” *Test*, **13**(2), 263–312. doi:10.1007/bf02595775.
- Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB (2013). *Bayesian Data Analysis*. Texts in Statistical Science, 3rd edition. Chapman & Hall/CRC. doi:10.1201/b16018.
- Gerber F (2021). *gapfill: Fill Missing Values in Satellite Data*. R package version 0.9.6-1, URL <https://CRAN.R-project.org/package=gapfill>.
- Gneiting T, Raftery AE (2007). “Strictly Proper Scoring Rules, Prediction, and Estimation.” *Journal of the American Statistical Association*, **102**(477), 359–378. doi:10.1198/016214506000001437.

- Gramacy RB (2016). “**laGP**: Large-Scale Spatial Modeling Via Local Approximate Gaussian Processes in R.” *Journal of Statistical Software*, **72**(1), 1–46. doi:[10.18637/jss.v072.i01](https://doi.org/10.18637/jss.v072.i01).
- Gramacy RB, Apley DW (2015). “Local Gaussian Process Approximation for Large Computer Experiments.” *Journal of Computational and Graphical Statistics*, **24**(2), 561–578. doi:[10.1080/10618600.2014.914442](https://doi.org/10.1080/10618600.2014.914442).
- Guhaniyogi R, Banerjee S (2018). “Meta-Kriging: Scalable Bayesian Modeling and Inference for Massive Spatial Datasets.” *Technometrics*, **60**(4), 430–444. doi:[10.1080/00401706.2018.1437474](https://doi.org/10.1080/00401706.2018.1437474).
- Guhaniyogi R, Finley AO, Banerjee S, Gelfand AE (2011). “Adaptive Gaussian Predictive Process Models for Large Spatial Datasets.” *Environmetrics*, **22**(8), 997–1007. doi:[10.1002/env.1131](https://doi.org/10.1002/env.1131).
- Guinness J (2018). “Permutation and Grouping Methods for Sharpening Gaussian Process Approximations.” *Technometrics*, **60**(4), 415–429. doi:[10.1080/00401706.2018.1437476](https://doi.org/10.1080/00401706.2018.1437476).
- Guisan A, Zimmermann NE (2000). “Predictive Habitat Distribution Models in Ecology.” *Ecological Modelling*, **135**, 147–186. doi:[10.1016/S0304-3800\(00\)00354-9](https://doi.org/10.1016/S0304-3800(00)00354-9).
- Hanks E, Schliep E, Hooten M, Hoeting J (2015). “Restricted Spatial Regression in Practice: Geostatistical Models, Confounding, and Robustness Under Model Misspecification.” *Environmetrics*, **26**, 243–254. doi:[10.1002/env.2331](https://doi.org/10.1002/env.2331).
- Hansen MC, Potapov PV, Moore R, Hancher M, Turubanova SA, Tyukavina A, Thau D, Stehman SV, Goetz SJ, Loveland TR, Kommareddy A, Egorov A, Chini L, Justice CO, Townshend JRG (2013). “High-Resolution Global Maps of 21st-Century Forest Cover Change.” *Science*, **342**(6160), 850–853. doi:[10.1126/science.1244693](https://doi.org/10.1126/science.1244693).
- Heaton MJ, Datta A, Finley AO, Furrer R, Guinness J, Guhaniyogi R, Gerber F, Gramacy RB, Hammerling D, Katzfuss M, Lindgren F, Nychka DW, Sun F, Zammit-Mangion A (2019). “A Case Study Competition Among Methods for Analyzing Large Spatial Data.” *Journal of Agricultural, Biological and Environmental Statistics*, **24**(3), 398–425. doi:[10.1007/s13253-018-00348-w](https://doi.org/10.1007/s13253-018-00348-w).
- Hurt GC, Dubayah R, Drake J, Moorcroft PR, Pacala SW, Blair JB, Fearon MG (2004). “Beyond Potential Vegetation: Combining Lidar Data and a Height-Structured Model for Carbon Studies.” *Ecological Applications*, **14**(3), 873–883. doi:[10.1890/02-5317](https://doi.org/10.1890/02-5317).
- ICESat-2 (2015). “Ice, Cloud, and Land Elevation Satellite-2.” Accessed: 2022-01-04, URL <http://icesat.gsfc.nasa.gov/>.
- Intel Corp (2019). “Intel Math Kernel Library.” Santa Clara.
- Katzfuss M (2017). “A Multi-Resolution Approximation for Massive Spatial Datasets.” *Journal of the American Statistical Association*, **112**(517), 201–214. doi:[10.1080/01621459.2015.1123632](https://doi.org/10.1080/01621459.2015.1123632).
- Klein T, Randin C, Korner C (2015). “Water Availability Predicts Forest Canopy Height at the Global Scale.” *Ecology Letters*, **18**(12), 1311–1320. doi:[10.1111/ele.12525](https://doi.org/10.1111/ele.12525).

- Lany NK, Zarnetske PL, Finley AO, McCullough DG (2020). “Complementary Strengths of Spatially-Explicit and Multi-Species Distribution Models.” *Ecography*, **43**(3), 456–466. doi:10.1111/ecog.04728.
- Lefsky MA (2010). “A Global Forest Canopy Height Map from the Moderate Resolution Imaging Spectroradiometer and the Geoscience Laser Altimeter System.” *Geophysical Research Letters*, **37**(15). doi:10.1029/2010gl1043622.
- Lindgren F, Rue H, Lindström J (2011). “An Explicit Link Between Gaussian Fields and Gaussian Markov Random Fields: The Stochastic Partial Differential Equation Approach.” *Journal of the Royal Statistical Society B*, **73**(4), 423–498. doi:10.1111/j.1467-9868.2011.00777.x.
- Nychka D, Bandyopadhyay S, Hammerling D, Lindgren F, Sain S (2015). “A Multiresolution Gaussian Process Model for the Analysis of Large Spatial Datasets.” *Journal of Computational and Graphical Statistics*, **24**(2), 579–599. doi:10.1080/10618600.2014.914946.
- Nychka D, Hammerling D, Sain S, Lenssen N, Smirniotis C, Iverson M (2019). **LatticeKrig: Multi-Resolution Kriging Based on Markov Random Fields**. R package version 8.4, URL <https://CRAN.R-project.org/package=LatticeKrig>.
- Plummer M, Best N, Cowles K, Vines K (2006). “**coda**: Convergence Diagnosis and Output Analysis for MCMC.” *R News*, **6**(1), 7–11.
- Polson NG, Scott JG, Windle J (2013). “Bayesian Inference for Logistic Models Using Pólya-Gamma Latent Variables.” *Journal of the American Statistical Association*, **108**(504), 1339–1349. doi:10.1080/01621459.2013.829001.
- Pulliam HR (2000). “On the Relationship Between Niche and Distribution.” *Ecology Letters*, **3**(4), 349–361. doi:10.1046/j.1461-0248.2000.00143.x.
- Ra S, Kim J (1993). “A Fast Mean-Distance-Ordered Partial Codebook Search Algorithm for Image Vector Quantization.” *IEEE Transactions on Circuits and Systems II: Analog and Digital Signal Processing*, **40**(9), 576–579. doi:10.1109/82.257335.
- Rasmussen CE (2003). “Gaussian Processes in Machine Learning.” In *Summer School on Machine Learning*, pp. 63–71. Springer-Verlag.
- R Core Team (2022). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Ribeiro Jr PJ, Diggle PJ (2020). **geoR: Analysis of Geostatistical Data**. R package version 1.8-1, URL <https://CRAN.R-project.org/package=geoR>.
- Risser MD, Turek D (2020). “Bayesian Inference for High-Dimensional Nonstationary Gaussian Processes.” *Journal of Statistical Computation and Simulation*, **90**(16), 2902–2928. doi:10.1080/00949655.2020.1792472.
- Roberts GO, Rosenthal JS (2009). “Examples of Adaptive MCMC.” *Journal of Computational and Graphical Statistics*, **18**(2), 349–367. doi:10.1198/jcgs.2009.06134.

- Rue H, Held L (2005). *Gaussian Markov Random Fields: Theory and Applications*. Monographs on Statistics & Applied Probability. Chapman & Hall/CRC.
- Rue H, Martino S, Lindgren F, Simpson D, Riebler A, Krainski ET, Fuglstad GA (2021). *INLA: Bayesian Analysis of Latent Gaussian Models Using Integrated Nested Laplace Approximations*. R package version 21.06.11, URL <http://r-inla.org/>.
- Saha A, Datta A (2018). “BRISC: Bootstrap for Rapid Inference on Spatial Covariances.” *Stat*, p. e184. doi:10.1002/sta4.184.
- Saha A, Datta A (2022). *BRISC: Fast Inference for Large Spatial Datasets Using BRISC*. R package version 1.0.5, URL <https://CRAN.R-project.org/package=BRISC>.
- Sang H, Jun M, Huang JZ (2011). “Covariance Approximation for Large Multivariate Spatial Data Sets with an Application to Multiple Climate Model Errors.” *The Annals of Applied Statistics*, pp. 2519–2548. doi:10.1214/11-aos478.
- Spiegelhalter DJ, Best NG, Carlin BP, Van Der Linde A (2002). “Bayesian Measures of Model Complexity and Fit.” *Journal of the Royal Statistical Society B*, **64**(4), 583–639. doi:10.1111/1467-9868.00353.
- Stein ML (2012). *Interpolation of Spatial Data: Some Theory for Kriging*. Springer-Verlag. doi:10.1007/978-1-4612-1494-6.
- Stein ML (2014). “Limitations on Low Rank Approximations for Covariance Matrices of Spatial Data.” *Spatial Statistics*, **8**, 1–19. doi:10.1016/j.spasta.2013.06.003. Spatial Statistics Miami.
- Stein ML, Chi Z, Welty LJ (2004). “Approximating Likelihoods for Large Spatial Data Sets.” *Journal of the Royal Statistical Society B*, **66**(2), 275–296. doi:10.1046/j.1369-7412.2003.05512.x.
- Stratton RD (2006). “Guidance on Spatial Wildland Fire Analysis: Models, Tools, and Techniques.” *General Technical Report RMRS-GTR-183*, U.S. Department of Agriculture, Forest Service, Rocky Mountain Research Station.
- Sun Y, Li B, Genton MG (2011). “Geostatistics for Large Datasets.” In JM Montero, E Porcu, M Schlather (eds.), *Advances and Challenges in Space-Time Modelling of Natural Events*, pp. 55–77. Springer-Verlag.
- Taylor-Rodriguez D, Finley AO, Datta A, Babcock C, Andersen HE, Cook BD, Morton DC, Banerjee S (2019). “Spatial Factor Models for High-Dimensional and Large Spatial Data: An Application in Forest Variable Mapping.” *Statistica Sinica*, **29**(3), 1155–1180. doi:10.5705/ss.202018.0005.
- Turek D, Risser M (2022). *BayesNSGP: Bayesian Analysis of Non-Stationary Gaussian Process Models*. R package version 0.1.2, URL <https://CRAN.R-project.org/package=BayesNSGP>.
- Vecchia AV (1988). “Estimation and Model Identification for Continuous Spatial Processes.” *Journal of the Royal Statistical Society B*, **50**, 297–312. doi:10.1111/j.2517-6161.1988.tb01729.x.

- Zammit-Mangion A, Cressie N (2021). “**FRK**: An R Package for Spatial and Spatio-Temporal Prediction with Large Datasets.” *Journal of Statistical Software*, **98**(4), 1–48. doi:10.18637/jss.v098.i04.
- Zhang L, Banerjee S, Finley AO (2021). “High-Dimensional Multivariate Geostatistics: A Bayesian Matrix-Normal Approach.” *Environmetrics*, **13**(3), 269–280. doi:10.1002/env.2675.
- Zhang X (2016). “An Optimized BLAS Library Based on GotoBLAS2.” <https://github.com/xianyi/OpenBLAS/>. Accessed: 2022-01-04.

**Affiliation:**

Andrew O. Finley  
Department of Forestry  
Michigan State University  
Natural Resources Building  
480 Wilson Road, Room 126  
East Lansing, MI 48824-6402, United States of America  
E-mail: [finleya@msu.edu](mailto:finleya@msu.edu)  
URL: <https://www.finley-lab.com/>

Abhirup Datta  
Department of Biostatistics  
Johns Hopkins Bloomberg  
School of Public Health  
615 N Wolfe Street, Room E3640  
Baltimore, MD 21205, United States of America  
E-mail: [abhidatta@jhu.edu](mailto:abhidatta@jhu.edu)  
URL: <http://abhidatta.com/>

Sudipto Banerjee  
Fielding School of Public Health  
University of California, Los Angeles  
650 Charles E. Young Dr. South  
Los Angeles, CA 90095-1772, United States of America  
E-mail: [sudipto@ucla.edu](mailto:sudipto@ucla.edu)  
URL: <https://ph.ucla.edu/faculty/banerjee>