



## Efficient Algorithms for Bayesian Nearest Neighbor Gaussian Processes

Andrew O. Finley, Abhirup Datta, Bruce D. Cook, Douglas C. Morton, Hans E. Andersen & Sudipto Banerjee

To cite this article: Andrew O. Finley, Abhirup Datta, Bruce D. Cook, Douglas C. Morton, Hans E. Andersen & Sudipto Banerjee (2019): Efficient Algorithms for Bayesian Nearest Neighbor Gaussian Processes, Journal of Computational and Graphical Statistics, DOI: [10.1080/10618600.2018.1537924](https://doi.org/10.1080/10618600.2018.1537924)

To link to this article: <https://doi.org/10.1080/10618600.2018.1537924>



View supplementary material [↗](#)



Accepted author version posted online: 28 Nov 2018.  
Published online: 01 Apr 2019.



Submit your article to this journal [↗](#)



Article views: 187



View Crossmark data [↗](#)



Citing articles: 2 View citing articles [↗](#)



# Efficient Algorithms for Bayesian Nearest Neighbor Gaussian Processes

Andrew O. Finley<sup>a</sup>, Abhirup Datta<sup>b</sup>, Bruce D. Cook<sup>c</sup>, Douglas C. Morton<sup>c</sup>, Hans E. Andersen<sup>d</sup>, and Sudipto Banerjee<sup>e</sup>

<sup>a</sup>Michigan State University, East Lansing, MI; <sup>b</sup>Johns Hopkins University, Baltimore, MD; <sup>c</sup>National Aeronautics and Space Administration, Washington, DC; <sup>d</sup>United States Forest Service, Washington, DC; <sup>e</sup>University of California, Los Angeles, Los Angeles, CA

## ABSTRACT

We consider alternate formulations of recently proposed hierarchical nearest neighbor Gaussian process (NNGP) models for improved convergence, faster computing time, and more robust and reproducible Bayesian inference. Algorithms are defined that improve CPU memory management and exploit existing high-performance numerical linear algebra libraries. Computational and inferential benefits are assessed for alternate NNGP specifications using simulated datasets and remotely sensed light detection and ranging data collected over the U.S. Forest Service Tanana Inventory Unit (TIU) in a remote portion of Interior Alaska. The resulting data product is the first statistically robust map of forest canopy for the TIU. Supplemental materials for this article are available online.

## ARTICLE HISTORY

Received June 2017  
Revised July 2018

## KEYWORDS

Bayesian methods;  
Computationally intensive  
methods; Spatial analysis;  
Statistical computing;  
Stochastic processes

## 1. Introduction

As spatial statisticians confront massive datasets with locations  $\sim 10^6$  and increasingly demanding inferential questions, several existing approaches that once seemed attractive for locations in the order of  $10^4$  become impractical. Recent methodological developments within the burgeoning literature on this subject aim to deliver massively scalable spatial processes. Sun, Li, and Genton (2011) and Banerjee (2017) provided background and more current work (also see references therein), respectively, in this area. A recent contribution by Heaton et al. (2017) is particularly useful as it provides an overview of modeling approaches for large spatial data that are under active 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 nearest neighbor Gaussian process (NNGP: Datta et al. 2016a) models, the comparison presented by Heaton et al. (2017) considered reduced rank predictive processes (Banerjee et al. 2008; Finley et al. 2009), covariance tapering (Furrer and Sain 2010; Furrer 2016), gap filling (Gerber 2017), metakriging (Guhaniyogi and Banerjee 2018), spatial partitioning (Sang, Jun, and Huang 2011; Barbian and Assunção 2017), fixed rank kriging (Cressie and Johannesson 2008; Zammit-Mangion and Cressie 2017), multiresolution approximation (Katzfuss 2017), stochastic partial differential equations (Rue et al. 2017), lattice kriging (Nychka et al. 2015), and local approximate Gaussian processes (Gramacy and Apley 2015; Gramacy 2016). 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.

With a few exceptions, for example, Furrer and Sain (2010) and Gramacy (2016), the literature on scalable spatial process models has focused primarily on theoretical and methodological developments with little attention to the algorithmic details needed for effectively applying them. For example, Datta et al. (2016a) implemented a “sequential” Gibbs sampler that involves updating a high-dimensional latent random effect vector and is prone to high autocorrelations and slow convergence. Most of the aforementioned articles do not discuss how researchers can, in practice, exploit high-performance computing libraries to obviate expensive numerical linear algebra (e.g., expensive matrix multiplications and factorizations) and deliver full Bayesian inference for massive spatial datasets. We address this gap for the NNGP models here by outlining three alternate formulations that are significantly more efficient for practical implementation than Datta et al. (2016a). Along with the accompanying code supplied with this manuscript, our intended contribution is well aligned with recent emphasis on reproducible research for challenging data analysis in the context of massive spatial datasets.

Our motivating scientific application concerns forest resource monitoring efforts and, in particular, to create fine resolution canopy height predictions using remotely sensed data collected at over 5 million locations. Spatially explicit estimates of forest canopy height are key inputs to a variety of 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 model parameters and predictions that can be propagated through subsequent computer models of ecosystem function to yield more robust error quantification. Bayesian inference is attractive here as it supplies full posterior predictive

distributions for the outcomes and for the latent process at arbitrary locations in the region of interest.

The remainder of this article proceeds as follows. [Section 2](#) provides a brief overview of NNGP models and their computational aspects. This is followed by three distinct and efficient alternate formulations: the collapsed NNGP model, a NNGP model for the outcomes themselves (with no latent process), and a conjugate NNGP model that allows MCMC-free inference. [Section 3](#) offers detailed simulation experiments on model performance and assessment and also presents a detailed analysis of the U.S. Forest Service Tanana Inventory Unit (TIU) dataset. Finally, [Section 4](#) concludes the manuscript with a summary and an eye toward future work.

## 2. Nearest Neighbor Gaussian Processes

Let  $y(s_i)$  and  $\mathbf{x}(s_i)$  denote the response and the predictors observed at location  $s_i$ ,  $i = 1, 2, \dots, n$ . A spatial linear mixed model posits  $y(s_i) = \mathbf{x}(s_i)^\top \boldsymbol{\beta} + w(s_i) + \epsilon(s_i)$ , where the random effect  $w(s_i)$  sums up the effect of unknown or unobserved spatial covariates, and  $\epsilon(s_i)$  denotes the independent and identically observed noise. Gaussian processes (GP) are commonly used for modeling the unknown surface  $w(s)$ . In particular,  $w(s) \sim GP(0, C(\cdot, \cdot | \boldsymbol{\theta}))$  implies that  $\mathbf{w} = (w(s_1), w(s_2), \dots, w(s_n))^\top$  is Gaussian with mean zero and covariance  $\mathbf{C} = (c_{ij})$ , where  $c_{ij} = C(s_i, s_j | \boldsymbol{\theta})$  and  $\boldsymbol{\theta}$  denotes the GP covariance parameters. A popular choice for  $C(\cdot, \cdot | \boldsymbol{\theta})$  is the Matérn covariance function specified as:

$$C(s_i, s_j; \sigma^2, \phi, \nu) = \frac{\sigma^2}{2^{\nu-1} \Gamma(\nu)} (||s_i - s_j|| \phi)^\nu \mathcal{K}_\nu(||s_i - s_j|| \phi);$$

$$\phi > 0, \nu > 0, \quad (1)$$

where  $\boldsymbol{\theta} = \{\sigma^2, \phi, \nu\}$  and  $\mathcal{K}$  denotes the Bessel function of second kind. Customary Bayesian hierarchical models are constructed as

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

where  $p(\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2)$  is specified by assigning priors to  $\boldsymbol{\beta}$ ,  $\boldsymbol{\theta}$ , and  $\tau^2$ . When  $n$  is very large, implementing (2) poses multiple computational roadblocks. Firstly, storing the matrix  $\mathbf{C}$  requires  $O(n^2)$  dynamic memory. Furthermore, evaluating  $N(\mathbf{w} | \mathbf{0}, \mathbf{C})$  involves factorizations (e.g., Cholesky) that require  $O(n^3)$  floating point operations (flops) to solve linear systems involving  $\mathbf{C}$  and computing  $\det(\mathbf{C})$ . Finally, predicting the response at  $K$  new locations require an additional  $O(Kn^2)$  flops. Alternative parameterizations such as integrating  $\mathbf{w}$  out of (2) shrinks the size of the parameter space, but does not obviate these computational bottlenecks. Even for moderately large spatial datasets, say with  $\sim 10^4$ – $10^5$  locations, these memory and storage demands become prohibitive. For the TIU dataset with  $5 \times 10^6$  locations, implementing (2) is practically impossible.

As mentioned in [Section 1](#), we pursue massive scalability for full Bayesian inference exploiting the NNGP. The underlying idea is familiar in graphical models (see, e.g., Lauritzen 1996; Murphy 2012). The joint distribution for a random vector  $\mathbf{w}$  can be looked upon as a directed acyclic graph (DAG). We write  $p(\mathbf{w}) = p(w_1) \prod_{i=2}^n p(w_i | \text{Pa}[i])$ , where  $w_i \equiv w(s_i)$  and  $\text{Pa}[i] = \{w_1, w_2, \dots, w_{i-1}\}$  is the set of parents of  $w_i$ . We can construct

sparse models for  $\mathbf{w}$  by shrinking the size of  $\text{Pa}[i]$ . In spatial contexts, this can be done by defining  $\text{Pa}[i]$  to be the set of  $w(s_j)$ 's corresponding to a small number  $m$  of nearest neighboring locations of  $s_i$ . Approximations resulting from such shrinkage were originally proposed by Vecchia (1988) and studied and exploited by Stein, Chi, and Welty (2004), Stroud, Stein, and Lysen (2017), Datta et al. (2016a, 2016c), and Huang and Sun (2018). The NNGP builds upon previous ideas and extends finite-dimensional likelihood approximations to well-defined sparsity-inducing Gaussian processes for estimating (2).

Working with multivariate Gaussian densities makes the connection between conditional independence in DAGs and sparsity abundantly clear. We can write the multivariate Gaussian density  $N(\mathbf{w} | \mathbf{0}, \mathbf{C})$  as a linear model,

$$w_1 = 0 + \eta_1 \quad \text{and}$$

$$w_i = a_{i1}w_1 + a_{i2}w_2 + \dots + a_{i,i-1}w_{i-1} + \eta_i \quad \text{for } i = 2, \dots, n,$$

or, more compactly, simply as  $\mathbf{w} = \mathbf{A}\mathbf{w} + \boldsymbol{\eta}$ , where  $\mathbf{A}$  is  $n \times n$  strictly lower-triangular with elements  $a_{ij} = 0$ , whenever  $j \geq i$  and  $\boldsymbol{\eta} \sim N(\mathbf{0}, \mathbf{D})$  and  $\mathbf{D}$  is diagonal with entries  $d_{11} = \text{var}(w_1)$  and  $d_{ii} = \text{var}(w_i | \{w_j : j < i\})$  for  $i = 2, \dots, n$ .

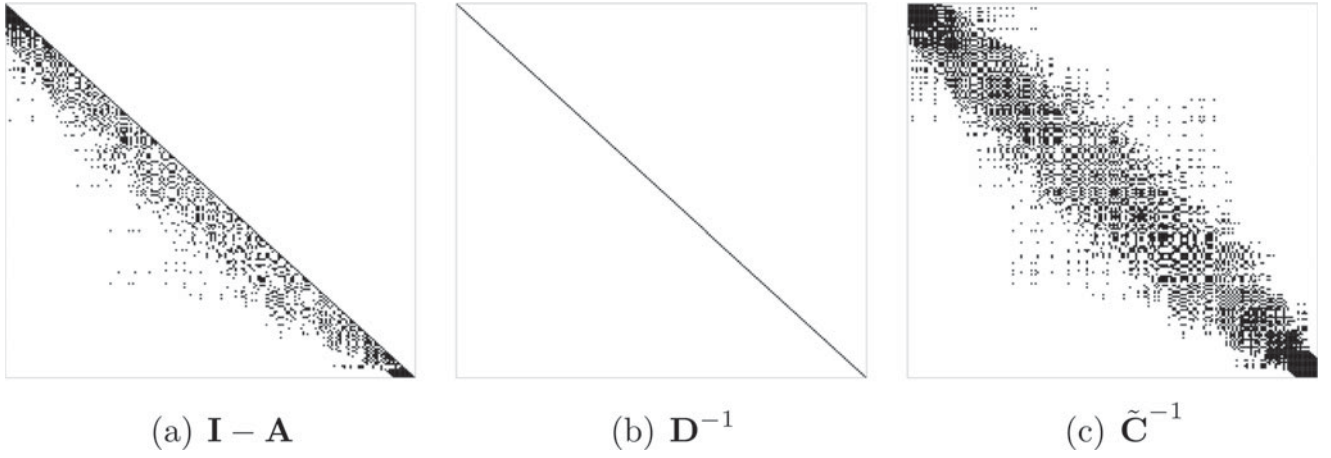
From the structure of  $\mathbf{A}$  it is evident that  $\mathbf{I} - \mathbf{A}$  is nonsingular and  $\mathbf{C} = (\mathbf{I} - \mathbf{A})^{-1} \mathbf{D} (\mathbf{I} - \mathbf{A})^{-\top}$ , where for any matrix  $\mathbf{M}$ ,  $\mathbf{M}^{-\top}$  refers to the inverse of its transpose. For any matrix  $\mathbf{M}$  and set of indices  $I_1, I_2 \subseteq \{1, 2, \dots, n\}$ , let  $\mathbf{M}[I_1, I_2]$  denote the submatrix of  $\mathbf{M}$  formed by the rows indexed by  $I_1$  and columns indexed by  $I_2$ . With the addition of  $\mathbf{D}[1, 1] = \mathbf{C}[1, 1]$  and the first row of  $\mathbf{A} = \mathbf{0}$ , the calculation of  $\mathbf{A}$  and  $\mathbf{D}$  is given in [Pseudocode 1](#), where  $1:i$  denotes the set  $\{1, 2, \dots, i\}$ ,  $\text{solve}(\mathbf{B}, \mathbf{b})$  computes the solution  $\mathbf{x}$  for the linear system  $\mathbf{B}\mathbf{x} = \mathbf{b}$ , and  $\text{dot}(\mathbf{u}, \mathbf{v})$  denotes the inner-product between two vectors  $\mathbf{u}$  and  $\mathbf{v}$ .

**Pseudocode 1.** Computation of  $\mathbf{A}$  and  $\mathbf{D}$ .

```
for(i in 1:(n-1)) {
  A[i+1, 1:i] = solve(C[1:i, 1:i],
                    C[1:i, i+1])
  D[i+1, i+1] = C[i+1, i+1]
               - dot(C[i+1, 1:i],
                    A[i+1, 1:i])
}
```

While [Pseudocode 1](#) computes the Cholesky decomposition of  $\mathbf{C}$ , there is no apparent gain to be had from the preceding computations since, as the loop runs into higher values of  $i$  closer to  $n$ , the dimension of  $\mathbf{C}[1:i, 1:i]$  increases. Consequently, one will need to solve larger and larger linear systems and the computational complexity remains  $O(n^3)$ . Nevertheless, it immediately shows how to exploit sparsity, if we set some elements in the lower triangular part of  $\mathbf{A}$  to be zero. For example, suppose, we permit no more than  $m$  elements in each row of  $\mathbf{A}$  to be nonzero. Let  $N[i]$  be the set of indices  $j < i$  such that  $\mathbf{A}[i, j] \neq 0$ . One can then compute the elements of  $\mathbf{A}$  and  $\mathbf{D}$  following [Pseudocode 2](#).

In [Pseudocode 2](#), we solve  $n-1$  linear systems of size at most  $m \times m$ , where  $m = \max_i |N(i)|$ . This can be performed in



**Figure 1.** Structure of the factors making up the sparse  $\tilde{\mathbf{C}}^{-1}$  matrix for  $n = 200$  and  $m = 10$ .

**Pseudocode 2.** Sparsity inducing computation of  $\mathbf{A}$  and  $\mathbf{D}$ .

```

for(i in 1:(n-1)) {
  A[i+1,N[i+1]]
    = solve(C[N[i+1],N[i+1]],
           C[N[i+1],i+1])
  D[i+1,i+1]
    = C[i+1,i+1] - dot(C[i+1, N[i+1]],
                       A[i+1,N[i+1]])
}

```

$O(nm^3)$  flops. Furthermore, these computations can be performed in parallel as each iteration of the loop is independent of the others. The above discussion provides a very useful strategy for constructing a sparse precision matrix. Starting with a dense  $n \times n$  matrix  $\mathbf{C}$ , we construct a sparse strictly lower-triangular matrix  $\mathbf{A}$  with no more than  $m(\ll n)$  nonzero entries in each row, and the diagonal matrix  $\mathbf{D}$  using **Pseudocode 2**, such that the matrix  $\tilde{\mathbf{C}} = (\mathbf{I} - \mathbf{A})^{-1}\mathbf{D}(\mathbf{I} - \mathbf{A})^{-\top}$  is a covariance matrix whose inverse  $\tilde{\mathbf{C}}^{-1} = (\mathbf{I} - \mathbf{A})^{\top}\mathbf{D}^{-1}(\mathbf{I} - \mathbf{A})$  is sparse. **Figure 1** presents a visual representation of the sparsity.

The factorization of  $\tilde{\mathbf{C}}^{-1}$  facilitates cheap computation of quadratic forms  $\mathbf{u}^{\top}\tilde{\mathbf{C}}^{-1}\mathbf{v}$  in terms  $\mathbf{A}$  and  $\mathbf{D}$ . The algorithm to evaluate such quadratic forms  $\text{qf}(\mathbf{u}, \mathbf{v}, \mathbf{A}, \mathbf{D})$  is provided in **Pseudocode 3**, where  $*$  and  $/$  denote multiplication and division by scalars, respectively.

**Pseudocode 3.** Computation of quadratic form.

```

qf(u,v,A,D) = u[1] * v[1] / D[1,1]
for(i in 2:n) {
  qf(u,v,A,D) = qf(u,v,A,D)
                + (u[i] - dot(A[i,N(i)],
                             u[N(i)]))
                * (v[i] - dot(A[i,N(i)],
                             v[N(i)])) / D[i,i]
}

```

Observe the algorithm in **Pseudocode 3** only involves inner products of  $m \times 1$  vectors. So, the entire `for` loop can be computed using  $O(nm)$  flops as compared to  $O(n^2)$  flops typically required to evaluate quadratic forms involving an  $n \times n$  dense

matrix. Also, importantly, the determinant of  $\tilde{\mathbf{C}}$  is obtained with almost no additional cost—it is simply  $\prod_{i=1}^n \mathbf{D}[i, i]$ .

Hence, while  $\tilde{\mathbf{C}}$  need not be sparse, the density  $N(\mathbf{w} | \mathbf{0}, \tilde{\mathbf{C}})$  is cheap to compute requiring only  $O(n)$  flops. This was exploited by Datta et al. (2016a), where the neighbor sets were constructed based on  $m$  nearest neighbors and the traditional GP prior for  $\mathbf{w}$  in (2) was replaced with an NNGP prior  $N(\mathbf{w} | \mathbf{0}, \tilde{\mathbf{C}})$ . The Markov chain Monte Carlo (MCMC) implementation of the NNGP model in Datta et al. (2016a) requires updating the  $n$  latent spatial effects  $\mathbf{w}$  sequentially, in addition to the regression and covariance parameters. While this ensures substantial computational scalability in terms of evaluating the likelihood, the behavior of MCMC convergence for such a high-dimensional model is difficult to study and may well prove unreliable.

We observed that, for very large spatial datasets, sequential updating of the random effects often leads to very poor mixing in the MCMC (see Figures S2 and S3). The computational gains per MCMC iteration is thus offset by a slow converging MCMC. Liu, Wong, and Kong (1994) showed that MCMC algorithms, where one or more variables are marginalized out tend to have lower autocorrelation and improved convergence behavior. Here, we explore NNGP models that drastically reduce the parameter dimensionality of the NNGP models by marginalizing over the entire vector of spatial random effects. Three different variants are developed, including an MCMC free conjugate model, and their relative merits and demerits are assessed both in terms of computational burden as well as model prediction and inference. Simulation experiments using spatial datasets of up to 10 million locations are conducted to assess the models' performance. Finally, we use the NNGP models to analyze the TIU dataset comprising over 5 million locations. To our knowledge, fully Bayesian analysis of spatial data at such scales is unprecedented.

## 2.1. Collapsed NNGP

The hierarchical model (2) or its NNGP analog impart a nice interpretation to the spatial random effects. The latent surface  $w(\mathbf{s})$  can provide a lot of information about the effect of missing covariates or unobserved physical processes. Hence, inference about  $\mathbf{w}$  is often critical for the researchers to improve the

understanding of the underlying scientific phenomenon. Here, we provide a collapsed NNGP model that enjoys the frugality of a low-dimensional MCMC chain but allows for full recovery of the latent random effects. We begin with the two-stage hierarchical specification  $N(\mathbf{y} | \mathbf{X}\boldsymbol{\beta} + \mathbf{w}, \tau^2 \mathbf{I}) \times N(\mathbf{w} | \mathbf{0}, \tilde{\mathbf{C}})$  and avoid sampling  $\mathbf{w}$  in the Gibbs' sampler by integrating out  $\mathbf{w}$  to obtain the collapsed NNGP model

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Lambda}) \text{ where } \boldsymbol{\Lambda} = \tilde{\mathbf{C}} + \tau^2 \mathbf{I}. \quad (3)$$

This model has only  $p + 4$  parameters compared to  $n + p + 4$  parameters in the hierarchical model. We use a conjugate prior  $N(\boldsymbol{\mu}_\beta, \mathbf{V}_\beta)$  for  $\boldsymbol{\beta}$ , inverse gamma priors for the spatial and noise variances, and uniform priors for the range and smoothness parameters. We use the  $u | \cdot$  notation to denote the full conditional distribution of any random variable  $u$  in the Gibbs' sampler. Let  $N(i)$  denote the set of indices corresponding to neighbor set of  $s_i$ . Observe that, although from Section 2, we know  $\tilde{\mathbf{C}} = (\mathbf{I} - \mathbf{A})^{-1} \mathbf{D} (\mathbf{I} - \mathbf{A})^{-\top}$ ,  $\boldsymbol{\Lambda}$  does not enjoy any such convenient factorization. In fact,  $\boldsymbol{\Lambda}^{-1}$  is also not guaranteed to be sparse, but exploiting the Sherman Woodbury Morrison (SWM) identity, we can write  $\boldsymbol{\Lambda}^{-1} = \tau^{-2} \mathbf{I} - \tau^{-4} \boldsymbol{\Omega}^{-1}$ , where  $\boldsymbol{\Omega} = (\tilde{\mathbf{C}}^{-1} + \tau^{-2} \mathbf{I})$  enjoys the same sparsity as  $\tilde{\mathbf{C}}^{-1}$ . Also, using a familiar determinant identity, we have  $\det(\boldsymbol{\Lambda}) = \tau^{2n} \det(\tilde{\mathbf{C}}) \det(\boldsymbol{\Omega})$ .

We exploit these matrix identities in conjunction with sparse matrix algorithms to obtain posterior distributions of the parameters  $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$ . In fact, the necessary computations can

be done by entirely avoiding expensive matrix computations and is described in detail in Algorithm 1. In addition to the inner product function  $\text{dot}(\cdot, \cdot)$  introduced earlier, we require a fill-reducing permutation matrix and a sparse Cholesky factorization ( $\text{sparsechol}(\cdot)$ ) for a sparse positive-definite matrix (note,  $\text{dot}(\cdot, \cdot)$ ,  $\text{sparsechol}(\cdot)$ , and subsequent functions that share this font are pseudocode). Large matrix-matrix and matrix-vector multiplications either involve at least one triangular matrix ( $\text{trmm}(\cdot, \cdot)$  or  $\text{trmv}(\cdot, \cdot)$ , where  $\text{mm}$  and  $\text{mv}$  denote matrix-matrix and matrix-vector operations), or at least one sparse matrix ( $\text{sparsemm}(\cdot, \cdot)$ , or  $\text{sparsemv}(\cdot, \cdot)$ ). We also use  $\text{diagsolve}(\cdot, \cdot)$  and  $\text{trsolve}(\cdot, \cdot)$  to solve linear systems with a diagonal or triangular coefficient matrix, respectively. We perform Cholesky decompositions, matrix-vector multiplications and solve linear equations involving general unstructured matrices using  $\text{chol}(\cdot)$ ,  $\text{gemv}(\cdot, \cdot)$ , and  $\text{solve}(\cdot, \cdot)$ , respectively, only for small  $p \times p$  or  $m \times m$  matrices, where both  $p$  and  $m$  are much less than  $n$ . Other utilities used in Algorithm 1 are  $\text{diag}(\cdot)$  to extract the diagonal elements of a matrix,  $\text{prod}(\cdot)$  to compute the product of the elements in a vector and  $\text{rnorm}(\cdot)$  to generate a specified number of random variables (as an integer argument) from a standard  $N(0, 1)$  distribution.

Observe that the entire Algorithm 1 is devoid of any expensive operations like  $\text{solve}$ ,  $\text{chol}$ , or  $\text{gemv}$  on dense  $n \times n$  matrices. All such operations are limited to  $m \times m$  or  $p \times p$  matrices, where both  $m$  and  $p$  are small. The computational costs in terms of flops of all such steps are listed in the algorithm and

---

#### Algorithm 1 Collapsed NNGP: Sampling from the posterior.

---

##### MCMC steps for updating $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$

---

##### 1: Metropolis-Hastings (MH) update for $\{\boldsymbol{\theta}, \tau^2\}$ :

$$p(\boldsymbol{\theta}, \tau^2 | \cdot) \propto p(\boldsymbol{\theta}, \tau^2) \times \frac{1}{\sqrt{\det(\boldsymbol{\Lambda})}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Lambda}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)$$

(a) Use Pseudocode 1 to obtain  $\mathbf{A}$  and  $\mathbf{D}$  using  $\mathbf{C}$  and  $\{N(i) | i = 1, 2, \dots, n\}$

$O(nm^3)$  flops

(b)  $\boldsymbol{\Omega} = \text{trmm}((\mathbf{I} - \mathbf{A})^\top, \text{diagsolve}(\mathbf{D}, \mathbf{I} - \mathbf{A})) + \tau^{-2} * \mathbf{I}$

$O(nm^2)$  flops

(c) Find a fill reducing permutation matrix  $\mathbf{P}$  for  $\boldsymbol{\Omega}$

(d)  $\mathbf{L} = \text{sparsechol}(\text{sparsemm}(\text{sparsemm}(\mathbf{P}, \boldsymbol{\Omega}), \mathbf{P}^\top))$

(e)  $\mathbf{r} = \mathbf{y} - \text{gemv}(\mathbf{X}, \boldsymbol{\beta})$ ;  $\mathbf{u} = \text{trsolve}(\mathbf{L}, \text{sparsemv}(\mathbf{P}, \mathbf{r}))$ ;  $\mathbf{v} = \text{trsolve}(\mathbf{L}^\top, \mathbf{u})$

$O(np)$  flops

(f)  $q = \text{dot}(\mathbf{r}, \mathbf{r})/\tau^2 - \text{dot}(\mathbf{r}, \text{sparsemv}(\mathbf{P}, \mathbf{v}))/\tau^4$

(g)  $d = \tau^{2*n} * \text{prod}(\text{diag}(\mathbf{D})) * \text{prod}(\text{diag}(\mathbf{L}))^2$

$O(n)$  flops

(h) Generate  $p(\boldsymbol{\theta}, \tau^2 | \cdot) \propto \frac{\exp(-q/2) * p(\boldsymbol{\theta}, \tau^2)}{\text{sqrt}(d)}$

##### 2: Gibb's sampler update for $\boldsymbol{\beta}$ :

$\boldsymbol{\beta} | \cdot \sim N(\mathbf{B}^{-1} \mathbf{b}, \mathbf{B}^{-1})$ , where  $\mathbf{B} = \mathbf{X}^\top \boldsymbol{\Lambda}^{-1} \mathbf{X} + \mathbf{V}_\beta^{-1}$  and  $\mathbf{b} = \mathbf{X}^\top \boldsymbol{\Lambda}^{-1} \mathbf{y} + \mathbf{V}_\beta^{-1} \boldsymbol{\mu}_\beta$

(a) for (j in 1:n) {  
     $\mathbf{u}_j = \text{trsolve}(\mathbf{L}, \text{sparsemv}(\mathbf{P}, \mathbf{X}[, j]))$ ;  $\mathbf{v}_j = \text{trsolve}(\mathbf{L}^\top, \mathbf{u}_j)$   
}

(b)  $\mathbf{F} = \text{solve}(\mathbf{V}_\beta, \mathbf{I})$ ;  $\mathbf{f} = \text{solve}(\mathbf{V}_\beta, \boldsymbol{\mu}_\beta)$

$O(p^3)$  flops

(c) Solve for  $p \times p$  matrix  $\mathbf{B}$  and  $p \times 1$  vector  $\mathbf{b}$ :

$O(np^2)$  flops

```
for (j in 1:p) {
     $\mathbf{b}[j] = \text{dot}(\mathbf{y}, \mathbf{X}[, j])/\tau^2 - \text{dot}(\mathbf{y}, \text{sparsemv}(\mathbf{P}, \mathbf{v}_j))/\tau^4 + \mathbf{f}[j]$ 
    for (i in 1:p) {
         $\mathbf{B}[i, j] = \text{dot}(\mathbf{X}[, i], \mathbf{X}[, j])/\tau^2 - \text{dot}(\mathbf{X}[, i], \text{sparsemv}(\mathbf{P}, \mathbf{v}_j))/\tau^4 + \mathbf{F}[i, j]$ 
    }
}
```

(d)  $\boldsymbol{\beta} = \text{solve}(\mathbf{B}, \mathbf{b}) + \text{trsolve}(\text{chol}(\mathbf{B}), \text{rnorm}(p))$

$O(p^3)$  flops

##### 3: Repeat Steps (1) and (2) (except Step 1(c)) $N$ times to obtain $N$ MCMC samples for $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$

---



are linear in  $n$ . However, the exact cost of the steps involving  $\mathbf{L}$  in Algorithm 1 (Steps 1(d)–(e)) depends on the data design. Although  $\mathbf{\Omega}$  is sparse  $O(nm^2)$  nonzero entries, the sparsity of its Cholesky factor  $\mathbf{L}$  actually depends on the location of the nonzero entries. Hence, we used a fill reducing permutation  $\mathbf{P}$  that increases the sparsity of the Cholesky factor. Although  $\mathbf{P}$  needs to be evaluated only once before the MCMC, finding the optimal  $\mathbf{P}$  yielding the least fill-in is an NP-complete problem. Hence, algorithms have been proposed to improve sparsity patterns based on a variety of fill-in minimizing heuristics, see, for example, Amestoy, Davis, and Duff (1996), Karypis and Kumar (1998), and Hager (2002) (also see Section 3).

When flops per iteration of MCMC are considered, computational requirements for the collapsed NNGP model is data dependent and may exceed the exact linear flops usage for the hierarchical NNGP Algorithm. We also observed this in simulation experiments described in Section 3. However, the improved MCMC convergence for the collapsed NNGP, as observed in Figures S2 and S5, implies that substantial computational gains accrue by truncating the MCMC run. Furthermore, all the `for` loops in Algorithm 1 can be evaluated independent of each other using parallel computing resources.

The collapsed model nicely separates the MCMC sampler for parameter estimation from posterior estimation of spatial random effects and subsequent predictions. Computational benefits accrue from using the quantities  $\mathbf{L}$  and  $\mathbf{u}$  already computed in Steps 1(d)–(e) of Algorithm 1 corresponding to the post-convergence samples of  $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$ . This is presented in Algorithm 2.

---

#### Algorithm 2 Collapsed NNGP: Posterior predictive inference

---

Post-MCMC steps using  $\mathbf{L}$  and  $\mathbf{u}$  from Steps 1(d)–(e) of Algorithm 1 for post-convergence samples of  $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$

---

- 1: **Sample from  $p(\mathbf{w} | \cdot)$  one-for-one for each post-convergence sample of  $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$**   
 $\mathbf{w} | \cdot \sim N(\mathbf{B}^{-1}\mathbf{b}, \mathbf{B}^{-1})$ , where  $\mathbf{B} = \tilde{\mathbf{C}}^{-1} + \tau^{-2}\mathbf{I}$  and  $\mathbf{b} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})/\tau^2$ 
    - (a)  $\mathbf{z} = \text{rnorm}(n)$   $O(n)$  flops
    - (b)  $\mathbf{w} = \text{sparsemv}(\mathbf{P}^\top, \text{trsolve}(\mathbf{L}^\top, \mathbf{u}/\tau^2 + \mathbf{z}))$
  - 2: **Prediction at a new location  $\mathbf{s}_0$ :**  
 $y(\mathbf{s}_0) | \cdot \sim N(\mathbf{x}(\mathbf{s}_0)^\top \boldsymbol{\beta} + w(\mathbf{s}_0), \tau^2)$ 
    - (a) Find  $N_0$  — set of  $m$  nearest neighbors of  $\mathbf{s}_0$  among  $\{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n\}$   $O(n)$  flops
    - (b)  $\mathbf{c} = \mathbf{C}(\mathbf{s}_0, N_0; \boldsymbol{\theta})$   $O(m)$  flops
    - (c)  $\mathbf{m} = \text{dot}(\mathbf{c}, \text{solve}(\mathbf{C}(N_0, N_0), \mathbf{w}[N_0]))$   $O(m^3)$  flops  
 $\mathbf{v} = \mathbf{C}(\mathbf{s}_0, \mathbf{s}_0; \boldsymbol{\theta}) - \text{dot}(\mathbf{c}, \text{solve}(\mathbf{C}(N_0, N_0), \mathbf{c}))$
    - (d)  $w(\mathbf{s}_0) = \mathbf{m} + \text{sqrt}(\mathbf{v}) * \text{rnorm}(1)$   $O(p)$  flops  
 $y(\mathbf{s}_0) = \text{dot}(\mathbf{x}(\mathbf{s}_0), \boldsymbol{\beta}) + w(\mathbf{s}_0) + \tau * \text{rnorm}(1)$   $O(p)$  flops
- 

Algorithm 2 demonstrates how inference on  $w(\mathbf{s})$  and  $y(\mathbf{s})$  can be easily achieved for any spatial location using the post burn-in samples of  $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$ . We first sample the spatial random effects  $p(\mathbf{w} | \mathbf{y})$  for the observed locations, use them to sample from  $p(w(\mathbf{s}_0) | \mathbf{y})$  and then from  $p(y(\mathbf{s}_0) | \mathbf{y})$ .

## 2.2. NNGP for the Response

Both the sequential NNGP Algorithm in Datta et al. (2016a) or the collapsed version in Section 2.1 accomplishes prediction

at a new location via recovering the spatial random effects first, proceeded by kriging at the new location. This differed from Vecchia's (1988) original approach, which applied nearest neighbor approximation directly to the marginal likelihood of  $\mathbf{y}$ . The recovery of the spatial random effects becomes necessary if inference on the latent process is of interest. Although recovering  $\mathbf{w}$ , as discussed earlier, has its own importance, if spatial interpolation of the response is the primary objective, this intermediate step is often a computational burden. In this Section, we propose a NNGP model for the response  $\mathbf{y}$  that sacrifices the ability to recover  $\mathbf{w}$  and directly predicts the response at new locations.

Datta et al. (2016a) demonstrated that a NNGP model can be derived from any GP. If  $w(\mathbf{s}) \sim \text{GP}(0, C(\cdot, \cdot))$ , then the response  $y(\mathbf{s}) \sim \text{GP}(\mathbf{x}(\mathbf{s})^\top \boldsymbol{\beta}, \Sigma(\cdot, \cdot))$  is also a GP, where  $\Sigma(\mathbf{s}_i, \mathbf{s}_j) = C(\mathbf{s}_i, \mathbf{s}_j) + \tau^2 I(\mathbf{s}_i = \mathbf{s}_j)$ . Hence, we can directly derive an NNGP for the response process  $y(\mathbf{s})$ . For finite dimensional realizations  $\mathbf{y}$ , likelihood under the response NNGP model is identical to Vecchia's composite likelihood. Datta et al. (2016a) extended this notion to a fully Bayesian setup. The key observation is that Vecchia's approximation corresponds to a proper multivariate Gaussian distribution obtained by simply replacing the covariance matrix  $\Sigma = \mathbf{C} + \tau^2 \mathbf{I}$  with its nearest-neighbor approximation  $\tilde{\Sigma}$  as described in Section 2. The sparsity properties documented in Section 2 apply to  $\tilde{\Sigma}$  as well. MCMC steps for parameter estimation and prediction using this response NNGP model are provided in Algorithm 3.

---

#### Algorithm 3 Response NNGP model: Sampling from the posterior

---

MCMC steps for updating  $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$

---

- 1: **Metropolis–Hastings (MH) update for  $\{\boldsymbol{\theta}, \tau^2\}$ :**  
 $p(\boldsymbol{\theta}, \tau^2 | \cdot) \propto p(\boldsymbol{\theta}, \tau^2) \times \frac{1}{\sqrt{\det(\tilde{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \tilde{\Sigma}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)$ 
    - (a) Use Pseudocode 1 to obtain  $\mathbf{A}$  and  $\mathbf{D}$  using  $\Sigma$  and  $\{N(i) | i = 1, 2, \dots, n\}$   $O(nm^3)$  flops
    - (b)  $\mathbf{e} = \mathbf{y} - \text{gemv}(\mathbf{X}, \boldsymbol{\beta})$ ; Using Pseudocode 3,  $\mathbf{q} = \mathbf{qf}(\mathbf{e}, \mathbf{e}, \mathbf{A}, \mathbf{D})$   $O(n(p+m))$  flops
    - (c)  $\mathbf{d} = \text{prod}(\text{diag}(\mathbf{D}))$   $O(n)$  flops
    - (d) Generate  $p(\boldsymbol{\theta}, \tau^2 | \cdot) \propto \frac{\exp(-\mathbf{q}/2) * p(\boldsymbol{\theta}, \tau^2)}{\text{sqrt}(\mathbf{d})}$   $O(1)$  flops
  - 2: **Gibb's sampler update for  $\boldsymbol{\beta}$ :**  
 $\boldsymbol{\beta} | \cdot \sim N(\mathbf{B}^{-1}\mathbf{b}, \mathbf{B}^{-1})$ , where  $\mathbf{B} = \mathbf{X}^\top \tilde{\Sigma}^{-1} \mathbf{X} + \mathbf{V}_\beta^{-1}$  and  $\mathbf{b} = \mathbf{X}^\top \tilde{\Sigma}^{-1} \mathbf{y} + \mathbf{V}_\beta^{-1} \boldsymbol{\mu}_\beta$ 
    - (a)  $\mathbf{F} = \text{solve}(\mathbf{V}_\beta, \mathbf{I})$ ;  $\mathbf{f} = \text{solve}(\mathbf{V}_\beta, \boldsymbol{\mu}_\beta)$   $O(p^3)$  flops
    - (b) Solve for  $p \times p$  matrix  $\mathbf{B}$  and  $p \times 1$  vector  $\mathbf{b}$  using Pseudocode 3:  $O(nmp^2)$  flops  

```

for (i in 1:p) {
  b[i] = qf(X[, i], Y, A, D) + f[i]
  for (j in 1:p) {
    B[1, j] = qf(X[, i], X[, j], A, D) + F[1, j]
  }
}

```
    - (c)  $\boldsymbol{\beta} = \text{solve}(\mathbf{B}, \mathbf{b}) + \text{trsolve}(\text{chol}(\mathbf{B}), \text{rnorm}(p))$   $O(p^3)$  flops
  - 3: Repeat Steps (1) and (2)  $N$  times to obtain  $N$  MCMC samples for  $\{\boldsymbol{\beta}, \boldsymbol{\theta}, \tau^2\}$
-

Unlike the collapsed NNGP model, the computational cost for each step of [Algorithm 3](#) does not depend on the spatial design of the data and is exactly linear in  $n$ . This is a result of the complete absence of the latent spatial effects  $\mathbf{w}$  in the model. Once again, parallel computing can be leveraged to evaluate all the `for` loops. A caveat with the response model is that recovery of  $\mathbf{w}$  is not possible as highlighted in [Datta et al. \(2016a\)](#). However, if that is of peripheral concern, the response model offers a computationally parsimonious solution for fully Bayesian analysis of massive spatial datasets. Posterior predictive inference, therefore, consists only of predicting the outcome  $y(\mathbf{s})$  at any arbitrary location  $\mathbf{s}$ . This is achieved easily through [Algorithm 4](#), where  $\mathbf{y}_{N(\mathbf{s}_0)}$  represents the subvector of  $\mathbf{y}$  corresponding to the points in  $N(\mathbf{s}_0)$ ,  $\mathbf{X}_{N(\mathbf{s}_0)}$  is the corresponding design matrix, and  $\Sigma_0$  is the  $m \times m$  covariance matrix for  $\mathbf{y}_{N(\mathbf{s}_0)}$ .

---

**Algorithm 4** Response NNGP model: Posterior predictive inference

---

Post-MCMC steps using post-convergence samples of  $\{\beta, \theta, \tau^2\}$

---

- 1: **Sample from  $p(y(\mathbf{s}_0) | \cdot)$  one-for-one for each post-convergence sample of  $\{\beta, \theta, \tau^2\}$**   
 $y(\mathbf{s}_0) | \cdot \sim N(\mathbf{x}(\mathbf{s}_0)^\top \beta + \mathbf{c}_0^\top \Sigma_0^{-1} (\mathbf{y}_{N(\mathbf{s}_0)} - \mathbf{X}_{N(\mathbf{s}_0)} \beta), \Sigma(\mathbf{s}_0, \mathbf{s}_0) - \mathbf{c}_0^\top \Sigma_0^{-1} \mathbf{c}_0)$ 
    - (a) Find  $N_0$  — set of  $m$  nearest neighbors of  $\mathbf{s}_0$  among  $\{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n\}$   
 $O(n)$  flops
    - (b)  $\mathbf{c} = \Sigma(\mathbf{s}_0, N_0; \theta)$   $O(m)$  flops
    - (c)  $\mathbf{m} = \text{dot}(\mathbf{c}, \text{solve}(\Sigma[N_0, N_0], \mathbf{y}[N_0] - \text{dot}(\mathbf{X}[N_0, ], \beta)))$   $O(m^3)$  flops  
 $\mathbf{v} = \Sigma(\mathbf{s}_0, \mathbf{s}_0) - \text{dot}(\mathbf{c}, \text{solve}(\Sigma[N_0, N_0], \mathbf{c}))$
    - (d)  $y(\mathbf{s}_0) = \text{dot}(\mathbf{x}(\mathbf{s}_0), \beta) + \mathbf{m} + \text{sqrt}(\mathbf{v}) * \text{rnorm}(1)$   $O(p)$  flops
- 

### 2.3. MCMC-Free Exact Bayesian Inference Using Conjugate NNGP

The fully Bayesian approaches developed in [Datta et al. \(2016a\)](#) and in [Sections 2.1](#) and [2.2](#) provide complete posterior distributions for all parameters. However, for massive spatial datasets containing millions of observations, running the Gibbs' samplers for several thousand iterations may still be prohibitively slow. One advantage of NNGP over similar scalable statistical approaches for large spatial data is that it offers a probability model. Here, we exploit this fact to achieve exact Bayesian inference.

We define  $\alpha = \tau^2 / \sigma^2$  and rewrite the marginal model from [Section 2.2](#) as  $N(\mathbf{y} | \mathbf{X}\beta, \sigma^2 \mathbf{M})$ , where  $\mathbf{M} = \mathbf{G} + \alpha \mathbf{I}$  and  $\mathbf{G}$  denotes the Matern correlation matrix corresponding to the covariance matrix  $\mathbf{C}$ , that is,  $\mathbf{G}[i, j] = C(\mathbf{s}_i, \mathbf{s}_j, (1, \nu, \phi)^\top)$ . Once again, the analogous NNGP model can be obtained by replacing the dense matrix  $\mathbf{M}$  with its nearest-neighbor approximation  $\tilde{\mathbf{M}}$ . Note that  $\tilde{\mathbf{M}}$  depends on  $\alpha$ , the spatial range  $\phi$ , and smoothness  $\nu$ . Empirically, in spatial regression models, the spatial process parameters  $\phi$  and  $\nu$  are often not well estimated due to multimodality issues. In fixed domain asymptotic settings (see, e.g., [Zhang 2004](#)) it is impossible to jointly identify the spatial covariance parameters. Consequently, if inference for the covariance parameters is not of interest, it might be possible to

fix them at reasonable values with minimal effect on prediction or point estimates of other model parameters. For example, the smoothness parameter  $\nu$  could be fixed at 0.5, which reduces (1) to the exponential covariance function, and  $\phi$  and  $\alpha$  could be estimated using  $K$ -fold cross-validation.

For fixed  $\alpha$  and  $\phi$ , we obtain the familiar conjugate Bayesian linear regression model  $IG(\sigma^2 | a_\sigma, b_\sigma) \times N(\beta | \mu_\beta, \sigma^2 \mathbf{V}_\beta) \times N(\mathbf{y} | \mathbf{X}\beta, \sigma^2 \tilde{\mathbf{M}})$  with joint posterior distribution

$$p(\beta, \sigma^2 | \mathbf{y}) \propto \underbrace{IG(\sigma^2 | a_\sigma^*, b_\sigma^*)}_{p(\sigma^2 | \mathbf{y})} \times \underbrace{N(\beta | \mathbf{B}^{-1} \mathbf{b}, \sigma^2 \mathbf{B}^{-1})}_{p(\beta | \sigma^2, \mathbf{y})},$$

where  $a_\sigma^* = a_\sigma + n/2$ ,  $b_\sigma^* = b_\sigma + \frac{1}{2}(\mu_\beta^\top \mathbf{V}_\beta^{-1} \mu_\beta + \mathbf{y}^\top \tilde{\mathbf{M}}^{-1} \mathbf{y} - \mathbf{b}^\top \mathbf{B}^{-1} \mathbf{b})$ ,  $\mathbf{B} = \mathbf{V}_\beta^{-1} + \mathbf{X}^\top \tilde{\mathbf{M}}^{-1} \mathbf{X}$  and  $\mathbf{b} = \mathbf{V}_\beta^{-1} \mu_\beta + \mathbf{X}^\top \tilde{\mathbf{M}}^{-1} \mathbf{y}$ .

It is easy to directly sample  $\sigma^2 \sim IG(a_\sigma^*, b_\sigma^*)$  and then sample  $\beta \sim N(\mathbf{B}^{-1} \mathbf{b}, \sigma^2 \mathbf{B}^{-1})$  one-for-one for each drawn  $\sigma^2$ . This produces samples from the marginal posterior distributions  $\beta | \mathbf{y} \sim \text{MVS-}t_{2a_\sigma^*}(\mathbf{B}^{-1} \mathbf{b}, \frac{b_\sigma^*}{a_\sigma^*} \mathbf{B}^{-1})$  and  $\sigma^2 | \mathbf{y} \sim IG(a_\sigma^*, b_\sigma^*)$ , where  $\text{MVS-}t_\kappa(\mathbf{B}^{-1} \mathbf{b}, (b/a) \mathbf{B}^{-1})$  denotes the multivariate noncentral Student's  $t$  distribution with degrees of freedom  $\kappa$ , mean  $\mathbf{B}^{-1} \mathbf{b}$  and variance  $b \mathbf{B}^{-1} / (\kappa - 1)$ . The marginal posterior mean and variance for  $\sigma^2$  are  $b_\sigma^* / (a_\sigma^* - 1)$  and  $b_\sigma^{*2} / (a_\sigma^* - 1)^2 (a_\sigma^* - 2)$ , respectively.

Instead of sampling from the posterior directly, we prefer a fast evaluation of the marginal posterior distributions to effectively implement the aforementioned cross-validated approach. Steps for efficiently evaluating the above is provided in [Algorithm 5](#). The marginal posterior predictive distribution at a new location  $\mathbf{s}_0$  is given by  $y(\mathbf{s}_0) | \mathbf{y} \sim t_{2a_\sigma^*}(m_0, b_\sigma^* v_0 / a_\sigma^*)$ , where expressions for  $m_0$  and  $v_0$  are provided in [Step 3](#) of [Algorithm 5](#). We deploy hyper-parameter tuning based on  $K$ -fold cross-validation to choose the optimal  $\alpha$  and  $\phi$  from a grid of possible values. In our data analysis, we have chosen broad endpoints of the grid using exploratory variograms. However, as suggested by one reviewer, reparametrizing  $\alpha^* = \alpha / (1 + \alpha)$  and  $\phi^* = \phi / (1 + \phi)$  would ensure that the new hyper-parameters are within  $[0, 1]$  and can facilitate a more automated grid-search. In applications, where the exploratory variograms are inaccurate, the latter parameterization will possibly be more useful.

We denote the indices and locations corresponding to the  $k$ th fold of the data by  $I(k)$  and  $S(k)$ , respectively, whereas  $I(-k)$  and  $S(-k)$ , respectively, denote the analogous quantities when the  $k$ th fold is excluded from the data. Also, let  $N(i, k)$  denote the neighbor set for a location  $\mathbf{s}_i$  constructed from the locations in  $S(-k)$ . Details of the cross-validation procedure are also provided in [Algorithm 5](#).

[Algorithm 5](#) completely circumvents MCMC based iterative sampling and only requires at most  $O(n)$  flops per step. Although the calculations need to be replicated for every  $(\phi, \alpha)$  combination, unlike the MCMC based algorithms that run serially, this step can be run in parallel. Moreover, kriging is often less sensitive to the choice of the covariance parameters so cross-validation can be done at a moderately crude resolution on the  $(\phi, \alpha)$  domain. Hence, the Algorithm remains extremely fast. This incredible scalability makes the conjugate NNGP model an attractive choice for ultra high-dimensional spatial

**Algorithm 5** MCMC free posterior sampling for conjugate NNGP model**Hyper parameter tuning**

- 1: Fix  $\alpha$  and  $\phi$ , split the data into  $K$  folds.
  - (a) Find the collection of neighbor sets  $\mathcal{N} = \{N(i, k) : i = 1, 2, \dots, n; k = 1, 2, \dots, K\}$
- 2: Obtain posterior means for  $\beta$  and  $\sigma^2$  after removing the  $k^{th}$  fold of the data:
  - (a) Use [Pseudocode 1](#) to obtain  $\mathbf{A}(k)$  and  $\mathbf{D}(k)$  from  $\mathbf{M}[S(-k), S(-k)]$  and  $\mathcal{N}$   $O(nm^3)$  flops
  - (b)  $\mathbf{F} = \text{solve}(\mathbf{V}_\beta, \mathbf{I})$ ;  $\mathbf{f} = \text{solve}(\mathbf{V}_\beta, \mu_\beta)$   $O(p^3)$  flops
  - (c) Solve for  $p \times p$  matrix  $\mathbf{B}(k)$  and  $p \times 1$  vector  $\mathbf{b}(k)$  using [Pseudocode 3](#):  $O(nmp^2)$  flops

```

for (i in 1:p) {
  b(k)[i] = qf(X[S(-k), i], Y[S(-k)], A(k), D(k)) + f[i]
  for (j in 1:p) {
    B(k)[i, j] = qf(X[S(-k), i], X[S(-k), j], A(k), D(k)) + F[i, j]
  }
}

```
  - (d)  $\mathbf{V}(k) = \text{solve}(\mathbf{B}(k), \mathbf{I})$ ;  $\mathbf{g}(k) = \text{gemv}(\mathbf{V}(k), \mathbf{b}(k))$   $O(p^3)$  flops

$$\hat{a}_\sigma^*(k) = a_\sigma + (n - n/K) / 2$$

$$\hat{b}_\sigma^*(k) = b_\sigma + (\text{dot}(\mu_\beta, \mathbf{f}) + \text{qf}(\mathbf{Y}[S(-k)], \mathbf{Y}[S(-k)], \mathbf{A}(k), \mathbf{D}(k)) - \text{dot}(\mathbf{g}(k), \mathbf{b}(k))) / 2$$
  - (e)  $\hat{\beta} = \mathbf{g}(k)$ ;  $\hat{\sigma}^2 = \hat{b}_\sigma^*(k) / (\hat{a}_\sigma^*(k) - 1)$
- 3: Predicting posterior means of  $y[S(k)]$ :  $O(nm^3/K)$  flops

```

for (s in S(k)) {
  N(s, k) = m-nearest neighbors of s from S(-k)
  z = M(s, N(s, k))
  w = solve(M[N(s, k), N(s, k)], z)
  m0 = y(s) = dot(x(s), g(k)) + dot(w, (y[N(s, k)] - dot(X[N(s, k)], g(k))))
  u = x(s) - dot(X[N(s, k)], w)
  v0 = dot(u, gemv(V(k), u)) + 1 + alpha - dot(w, z)
  var(y(s)) = b_sigma^*(k) v0 / (a_sigma^*(k) - 1)
}

```
- 4: Root Mean Square Predictive Error (RMSPE) over  $K$  folds:  $O(n)$  flops
  - (a) Initialize  $e = 0$ 

```

for (k in 1:K) for (s_i in S[k]) {
  e = e + (y(s_i) - y_hat(s_i))^2
}

```
- 5: Cross-validation for choosing  $\alpha$  and  $\phi$ 
  - (a) Repeat steps (2) and (3) for  $G$  values of  $\alpha$  and  $\phi$   $O(GKnm(p^2 + m^2))$  flops
  - (b) Choose  $\alpha_0$  and  $\phi_0$  as the value that minimizes the average RMSPE  $O(G)$  flops

**Parameter estimation and prediction**

- 6: Repeat step (2) with  $(\alpha_0, \phi_0)^\top$  and the full data to get  $(\beta, \sigma^2) | \mathbf{y}$   $O(nmp^2 + nm^3)$  flops
- 7: Repeat step (3) with  $(\alpha_0, \phi_0)^\top$  and the full data to predict at a new location  $\mathbf{s}_0$  to obtain the mean and variance of  $y(\mathbf{s}_0) | \mathbf{y}$   $O(m^3)$  flops

data. Although this approach philosophically departs from the true Bayesian paradigm, often inference about covariance parameters is of little interest and this hybrid cross-validation approach offers a pragmatic compromise.

**3. Illustrations****3.1. Implementation**

This section details two simulation experiments and the analysis of a large remotely sensed dataset. In the analyses, we consider the candidate models labeled: *Sequential* defined in Datta et al. (2016a), *Collapsed* defined in [Section 2.1](#), *Response* defined in [Section 2.2](#), and *Conjugate* defined in [Section 2.3](#).

Two additional analyses are provided in the web supplement. First, [Section S3](#), compares full GP and NNGP model parameter estimates and predictive performance. Second,

[Section S4](#), moves beyond the typical geostatistical setting, where  $\mathbf{s}$  indexes data in two-dimensions, for example, latitude and longitude, to a more general settings, where data are indexed in  $N$ -dimensions. Such data are common in computer experiments, where  $\mathbf{s}$  indexes outcomes associated with a set of values on  $N$  computer model inputs. Here too, we apply a Matérn covariance function. Response and conjugate model out-of-sample predictive performance is shown to be comparable with that achieved using a local approximate Gaussian processes as implemented in the `laGP` R package (Gramacy and Sun 2017; Gramacy 2016).

Samplers were programmed in C++ and used `openBLAS` (Zhang 2016) and Linear Algebra Package (LAPACK; [www.netlib.org/lapack](http://www.netlib.org/lapack)) for efficient matrix computations. `openBLAS` is an implementation of Basic Linear Algebra Subprograms (BLAS; [www.netlib.org/blas](http://www.netlib.org/blas)) capable of exploiting multiple processors. Additional multiprocessor parallelization used



openMP (Dagum and Menon 1998) to improve performance of key steps within the samplers. In particular, substantial gains were realized by distributing the calculation of NNGP precision matrix components using the openMP `omp for directive`. Updating these matrices is necessary for each MCMC iteration in the sequential, response, and collapsed models, and for each conjugate model cross-validation iteration. An `omp for directive` with `reduction` clause was also effectively used to evaluate the [Pseudocode 3](#) found in all models.

For the collapsed model, SuiteSparse version 4.4.5 (Davis 2016a) provided an interface to: fill-in minimizing algorithms, for example, AMD (Amestoy, Davis, and Duff 2004), METIS (Karypis and Kumar 1998), and CHOLMOD (Chen et al. 2008) version 3.0.6 used for supernodal openBLAS-based Cholesky factorization to obtain  $\mathbf{L}$  of  $\mathbf{P}(\tilde{\mathbf{C}}^{-1} + \tau^{-2}\mathbf{I})\mathbf{P}^\top$ , and solvers for sparse triangular systems. Also see the text by Davis (2006).

For each analysis using the collapsed model, nine fill-in algorithms were considered (for details see Chen et al. 2008; Davis 2016b, pp. 4 and 16, respectively) for formation of the permutation matrix  $\mathbf{P}$ . Assessment of the various fill-in algorithms is based on the resulting pattern of nonzero matrix elements. This is important for our setting because the initial pattern of the NNGP precision matrix is determined by the neighbor set and, hence, discovery of an *optimal* permutation matrix need only be done once prior to sampling.

Implementing NNGP models requires a neighbor set for each observed location. For a given location  $\mathbf{s}_i$ , a brute force approach to finding the neighbor set calculates Euclidean distances to  $\mathbf{s}_1$ ,  $\mathbf{s}_2$ , and  $\mathbf{s}_{i-1}$ , sorts these distances while keeping track of locations' indexes, then selects the  $m$  minimum distance neighbors. This brute force approach is computationally demanding. Subsequent analyses use a relatively simple to implement fast nearest neighbor search algorithm proposed by Ra and Kim (1993) that provides substantial efficiency gains over the brute force search (see supplemental material for details).

All subsequent analyses were conducted on a Linux workstation with two 18-core Intel processors and 512 GB of memory. Unless otherwise noted, posterior inference used the last  $1 \times 10^4$  iterations from each of three chains of  $2.5 \times 10^4$  iterations. Chains run for a given model were initiated at different values and each chain was given a unique random number generator seed. Following Datta et al. (2016a), all models were fit using  $m = 15$  neighbors unless noted otherwise.

*Code and data needed to reproduce the analyses are provided in the web supplement.*

### 3.2. Experiment #1

The aim of this experiment was to assess NNGP model run time. To achieve this, we selected data subsets for a range of  $n$  from the TIU dataset described in [Sections 1](#) and [3.4](#). The posited model follows (2) and includes an intercept and slope regression coefficients, and an exponential covariance function with parameters  $\sigma^2$ ,  $\phi$ , and residual variance  $\tau^2$ . A “flat” improper prior distribution was assigned to each regression coefficient,  $\beta$ 's, which places equal weight on all possible values of the parameter. The variance components  $\tau^2$  and  $\sigma^2$  were assigned inverse-Gamma  $IG(2, 10)$  priors, and a uniform  $U(0.1, 10)$  prior for the

decay parameter  $\phi$ . The support on the decay corresponds to an effective spatial range (i.e., the distance where the spatial correlation is 0.05) between 0.3 and 30 km (see [Section 3.4](#) for specifics on the TIU domain and dataset).

[Figure 2\(a\)](#) shows run time for a dataset of  $n = 5 \times 10^4$  and number of CPUs used to complete one MCMC iteration (not including the initial nearest neighbor set search time, which is common across models). Two versions of the collapsed model are shown, one assumes the permutation matrix  $\mathbf{P}$  is diagonal (labeled *no perm*) and the other allows CHOLMOD to select an approximately optimal permutation matrix (labeled *perm*). Here, and in other experiments, using a fill-in reducing permutation matrix provides substantial time efficiency gains. The response model provides full posterior inference on all parameters, with the exception of  $\mathbf{w}$ , and dramatically faster run time compared to the collapsed model. Inference for the conjugate model, including  $\hat{\boldsymbol{\beta}}$  and  $\hat{\sigma}^2$  ([Algorithm 5](#)), requires about the same amount of time as one response model MCMC iteration. Explicitly updating  $\mathbf{w}$  is relatively slow; hence, the sequential model's computing time falls somewhere between that of the collapsed and response models.

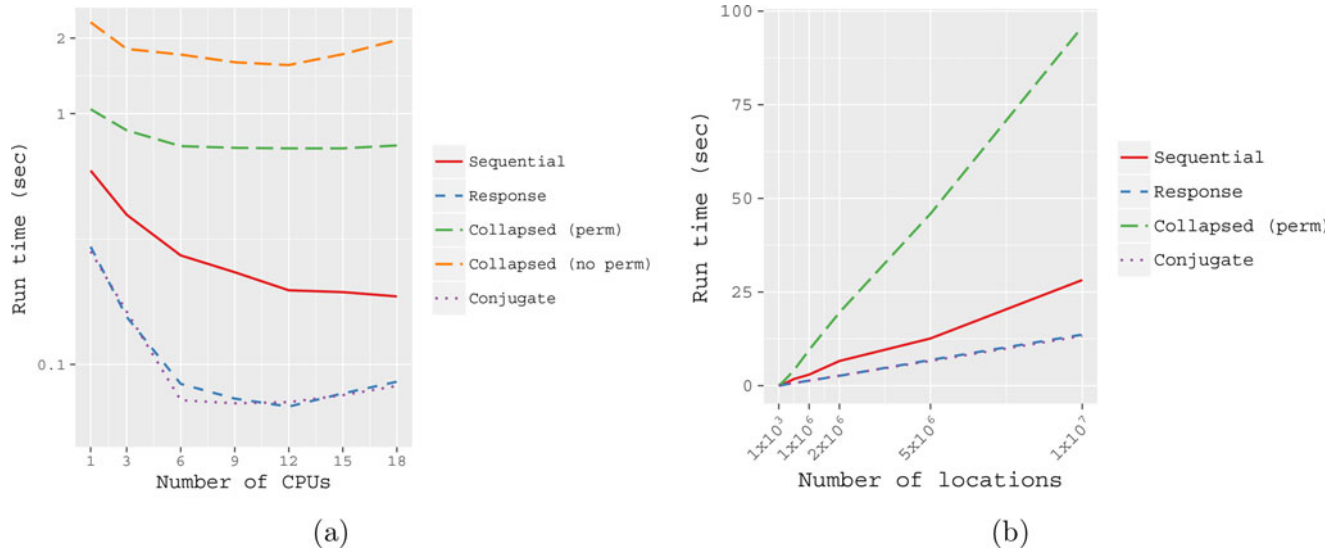
For all models, [Figure 2\(a\)](#) show marginal improvement in run time beyond  $\sim 6$  CPUs and negligible improvement beyond  $\sim 12$  CPUs. We attribute the slight increase in run time beyond  $\sim 12$  CPU seen in some models to communication overhead. Run time is actual execution time, or “wall clock time,” of the specified number of MCMC iterations. Points of diminishing return on number of CPUs used will change with  $n$ ; however, exploratory analysis across the range of  $n$  considered here suggested 12 CPUs is the bound for substantial gains (clearly this also depends on computing environment and programming decisions).

[Figure 2\(b\)](#) shows time required to execute one sampler iteration by  $n$ . The response and conjugate models deliver inference across  $n$  in  $\sim 1/3$  and  $\sim 1/10$  the time required by the Sequential and Collapsed models, respectively. For  $n=1 \times 10^7$  the run time is approximately 28, 13, 13, and 95 sec for the sequential, response, conjugate, and collapsed, respectively.

### 3.3. Experiment #2

This experiment compared parameters estimates and predictive performance among the NNGP models for a large dataset. Also, the potential to identify *optimal* values of  $\phi$  and  $\alpha$  via cross-validation was assessed for the conjugate model. We generated observations at  $6 \times 10^4$  locations within a unit square domain from model (2), the  $n \times n$  spatial covariance matrix  $\mathbf{C}$  was formed using (1) with  $\nu$  fixed at 0.5, and the mean comprised an intercept and covariate  $\mathbf{x}_1$  drawn from independent  $N(0, 1)$ . Observations were then generated using the parameter values given in the column labeled *True* in [Table 1](#). Observations at  $n = 5 \times 10^4$  of these locations, selected at random, were used to estimate model parameters. Observations at the remaining  $1 \times 10^4$  holdout locations were used to assess model predictive performance.

Following [Section 2.3](#), five-fold cross-validation aimed at minimizing RMSPE and continuous rank probability score (CRPS; Gneiting and Raftery 2007) for the conjugate model



**Figure 2.** (a) Run time required for one sampler iteration using  $n = 5 \times 10^4$  by number of CPUs (y-axis is on the log scale). (b) Run time required for one sampler iteration by number of locations.

**Table 1.** Simulated dataset, parameter credible intervals 50% (2.5%, 97.5%) and predictive validation.

Parameter	True	Sequential (metrop)	Sequential (slice)	Response	Collapsed	Conjugate
$\beta_0$	1	<b>0.64 (0.53, 0.75)</b>	<b>0.56 (0.44, 0.79)</b>	<b>0.84 (0.70, 0.99)</b>	1.10 (0.51, 1.79)	0.84
$\beta_1$	5	5.00 (5.00, 5.01)	5.00 (5.00, 5.01)	5.01 (5.00, 5.01)	5.00 (5.00, 5.01)	5.01
$\sigma^2$	1	<b>1.95 (1.44, 2.21)</b>	<b>1.68 (1.11, 2.19)</b>	1.03 (0.91, 1.21)	<b>1.69 (1.16, 2.24)</b>	0.98
$\tau^2$	1	1.00 (0.98, 1.01)	1.00 (0.98, 1.01)	1.00 (0.98, 1.01)	1.00 (0.98, 1.01)	1.02
$\phi$	6	<b>3.39 (3.03, 4.54)</b>	3.98 (3.04, 6.05)	6.26 (4.88, 7.78)	<b>3.95 (3.01, 5.83)</b>	4.05
CRPS		0.59	0.59	0.6	0.59	0.59
RMSPE		1.04	1.04	1.05	1.04	1.05
95% PIC		93.13	92.63	93.08	92.77	94.94
95% PIW		3.87	3.85	3.93	3.84	4.11

NOTE: Bold entries indicate where the true value is not within the 95% credible interval.

are given in Figure 3. We observe that a broad range of  $\phi$  and  $\alpha$  values deliver comparable predictive performance, and minimization of RMSPE and CRPS yield approximately the same estimates of  $\phi$  and  $\alpha$ .

In addition to RMSPE and CRPS, percent of holdout observations covered by their corresponding predictive distribution 95% credible interval (PCI), and mean width of the predictive distributions' 95% credible interval (PIW) were used to assess NNGP model predictive performance. Results given in Table 1 show the NNGP models yield comparable parameter estimates and prediction. Here, the conjugate model's  $\phi$  and  $\alpha$  were selected to minimize RMSPE (results are comparable for minimization of CRPS).

Candidate models' Gelman–Rubin (Gelman and Rubin 1992) potential scale reduction factor figures and MCMC chain trace plots are given in Figures S2–S5 of the web supplement. These figures show the response and collapsed models provide faster chain convergence for the intercept and spatial covariance parameters compared to sequential model. Additional analysis in Section S3 of the web supplement reveal that for a smaller dataset generated using the same model, the Sequential model parameter posteriors do not match well that of the full GP.

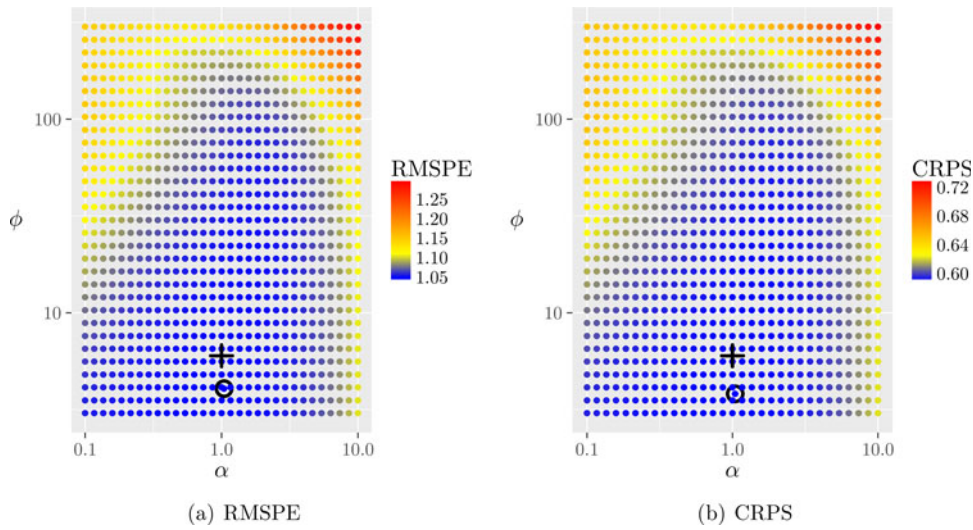
### 3.4. Tanana Inventory Unit Forest Canopy Height

Our goal is to create a high-resolution forest canopy height data product, with accompanying uncertainty estimates for

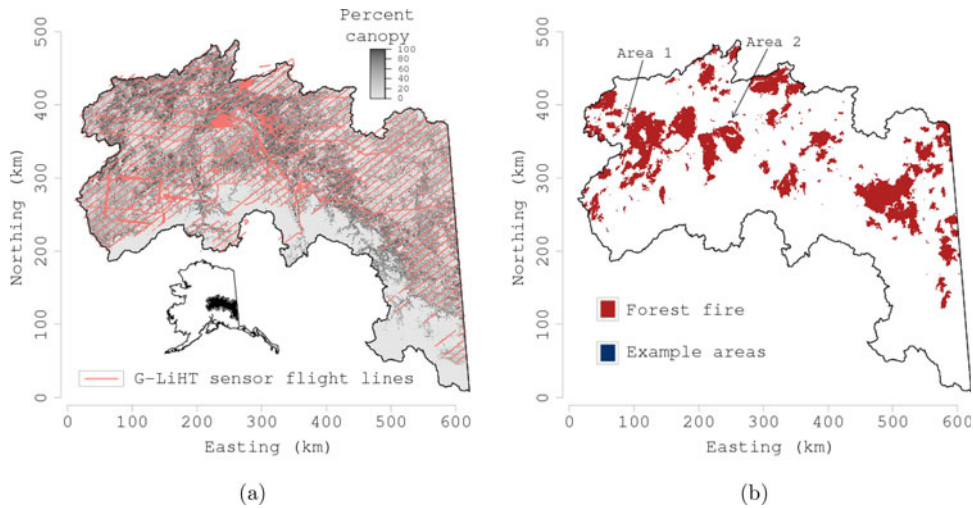
prediction and spatial correlation parameters, for the U.S. Forest Service Tanana Inventory Unit (TIU) that covers a large portion of Interior Alaska using a sparse sample of light detection and ranging (LiDAR) data from NASA Goddard's LiDAR, hyperspectral, and thermal (G-LiHT) Airborne Imager (Cook et al. 2013).

For remote forested regions, combining sparse airborne LiDAR data with a sparse network of forest inventory data provides a cost-effective means to deliver predictive maps of forest canopy height. In this study, LiDAR data were acquired across the US Forest Service TIU in Interior Alaska, approximately 140,000 km<sup>2</sup>, using the NASA Goddard's LiDAR, hyperspectral, and thermal (G-LiHT) Airborne Imager (Cook et al. 2013). The G-LiHT instrument package simultaneously acquires data from a suite of remote sensing instruments to collect complementary information on forest structure (LiDAR), vegetation composition (hyperspectral), and forest health (hyperspectral and thermal).

Here, we consider G-LiHT LiDAR data collected during a 2014 TIU flight campaign. The campaign collected a systematic sample covering ~8% of the TIU, with 78 parallel flight lines spaced ~9 km apart, Figure 4(a), along with incidental measurements to-and-from the transects. The nominal flying altitude of data collection in the TIU was 335 m above ground level, resulting in a sample swath width of ~180 m (30° field of view) and sample density of three laser pulses m<sup>2</sup>. Point cloud data were classified and used to generate bare earth elevation



**Figure 3.** Conjugate model cross-validation results for selection of  $\alpha$  and  $\phi$  using the simulated dataset. Parameter combination with minimum scoring rule indicated with open circle symbol  $\circ$  and true combination used to generate the data indicated with a plus symbol  $+$ .



**Figure 4.** TIU, Alaska, study region. (a) G-LiHT flight lines where canopy height was measured at  $5 \times 10^6$  locations and percent tree cover predictor variable. (b) Occurrence of forest fire within the past 20 years predictor variable and two example areas for prediction illustration.

and canopy height models at 1 m ground sample distance, as described in Cook et al. (2013). G-LiHT point cloud data and derived products are available online at <http://gliht.gsfc.nasa.gov>. The data was processed following methods in Cook et al. (2013), such that 28,751,400 LiDAR-based estimates of forest canopy height were available on a  $15 \times 15$  m grid along the flight lines. Each grid cell yielded an estimate of canopy height calculated as the height below, which 95% of the pulse data was recorded. The subsequent analysis uses a random sample of  $5.025 \times 10^6$  observations from the larger LiDAR dataset.

Two predictors that completely cover the TIU were considered. First, a Landsat derived percent tree cover data product developed by Hansen et al. (2013), shown as the gray scale surface in Figure 4(a). This product provides percent tree cover estimates for peak growing season in 2010 (most recent year available) and was created using a regression tree model applied to Landsat 7 ETM+ annual composites. These data are provided by the United States Geological Survey (USGS) on an approximate 30 m grid covering the entire globe (Hansen et al. 2013). Second, the perimeters of past fire events from 1947

to 2014 were obtained from the Alaska Interagency Coordination Center Alaska fire history data product (AICC 2016). Forest recovery/regrowth following fire is very slow in Interior Alaska. Hence, we discretized the fire history data to 1 if the fire occurred within the past 20 years and 0 otherwise, Figure 4(b).

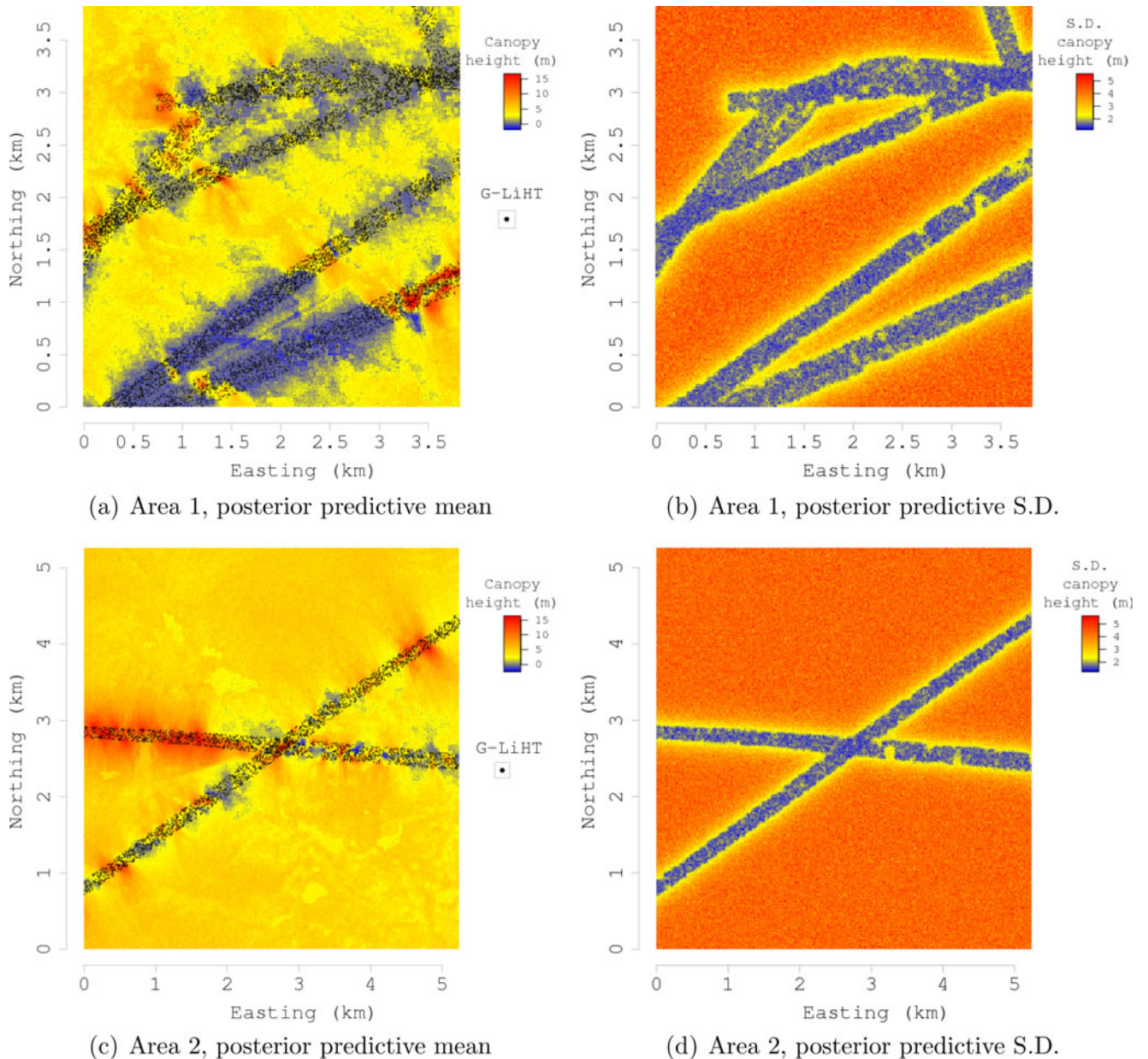
We explored the relationship between canopy height, tree cover, and fire history using a nonspatial regression model and NNGP response, collapsed, and conjugate models. We did not consider the sequential model here because of the convergence issues seen in the preceding experiments. Exploratory analysis using the nonspatial regression suggested both predictors explain a substantial portion of variability in canopy height (Table 2), with a positive association between canopy height and tree cover (TC) and negative association between canopy height and recent fire occurrence (fire). These results are consistent with our understanding of the TIU forest system. The tree cover variable captures forest canopy sparseness—with sparser canopies resulting in LiDAR height percentiles shifted toward the ground. Recently, burned areas are typically replaced with regenerating, shorter stature, forests.



**Table 2.** TIU dataset results.

Parameter	Nonspatial regression	Response	Collapsed	Conjugate minimize RMSPE
$\beta_0$	-2.46 (-2.47, -2.45)	2.37 (2.31, 2.42)	2.41 (2.35, 2.47)	2.51
$\beta_{TC}$	0.13 (0.13, 0.13)	0.02 (0.02, 0.02)	0.02 (0.02, 0.02)	0.02
$\beta_{Fire}$	-0.13 (-0.14, -0.12)	0.43 (0.39, 0.48)	0.39 (0.34, 0.43)	0.35
$\sigma^2$	—	17.29 (17.13, 17.41)	18.67 (18.50, 18.81)	23.21
$\tau^2$	17.39 (17.37, 17.41)	1.55 (1.54, 1.55)	1.56 (1.55, 1.56)	1.21
$\phi$	—	4.15 (4.13, 4.19)	3.73 (3.70, 3.77)	3.83
$\alpha$	—	—	—	0.052
CRPS	2.3	0.86	0.86	0.84
RMSPE	4.19	1.72	1.73	1.71
95% PIC	93.43	94.29	94.25	94.85
95% PIW	16.27	6.58	6.56	6.73
Run time (hours)	—	38.29	318.81	0.002

NOTE: Parameter credible intervals, 50% (2.5%, 97.5%), predictive validation, and run time for  $25 \times 10^3$  MCMC iterations.



**Figure 5.** About 95th LiDAR percentile height posterior predictive distribution summary at a 30 m pixel resolution for the two example areas identified in Figure 4(b).



For all models, the intercept and slope regression parameters were given flat prior distributions. The variance components  $\tau^2$  and  $\sigma^2$  were assigned inverse-Gamma  $IG(2, 10)$  priors. We assumed an exponential spatial correlation function with a uniform  $U(0.1, 10)$  prior on the decay parameter. The support on the decay corresponds to an effective spatial range between 0.3 and 30 km. Observations at  $n = 5 \times 10^6$  locations, selected at random, were used to estimate model parameters. Observations at the remaining  $2.5 \times 10^4$  holdout locations were used to assess model predictive performance. Parameter estimates and prediction performance summaries for candidate models are given in Table 2. Results for the  $m = 15$  and  $m = 25$  models were indistinguishable, hence, only  $m = 15$  results are presented. Here, NNGP models provide approximately the same predictive performance, and a substantial improvement over the nonspatial regression.

As suggested by Figure 2(b), and seen again here, the collapsed model using a full reducing permutation and 12 CPU requires an excessively long run time, that is, about 2 weeks to generate  $25 \times 10^3$  MCMC samples. If one is willing to forgo estimates of spatial random effects, the response model offers greatly improved run time, that is, about 1.5 days, and parameter and prediction inference comparable to the collapsed model. The conjugate model delivers the shortest run time and predictive inference comparable to the other NNGP models.

Figure 4(b) identifies two example areas selected to illustrate how LiDAR and the other data inform forest canopy height prediction. As suggested by the prediction metrics in Table 2, all three NNGP models delivered nearly identical prediction map products. Figure 5 shows the posterior predictive distribution mean and standard deviation from the response model with  $m = 15$  for the two areas. Here, the left subplots identify LiDAR data locations as black points along the flight lines. The presence of strong residual spatial autocorrelation results in fine-scale prediction within, and adjacent to, the flight lines (Figures 5(a) and (c)) and more precise posterior predictive distributions as reflected in the standard deviation maps (Figures 5(b) and (d)). Predictions more than a km from the flight lines are informed primarily by tree cover and fire occurrence predictors.

The TIU forest's vertical and horizontal structure is highly heterogeneous due, in large part to topography, hydrology, and disturbance history, for example, fire. This heterogeneity is reflected in the relatively short estimated effective range of just over 1 km (Table 2).

These results provide key input needed for planning future LiDAR campaigns to collect data to inform canopy height models. Using more informative predictor variables would certainly improve prediction across the TIU; however, few complete-coverage high spatial resolution data layers exist, other than those produced using moderate spatial resolution remote sensing products, for example, the Landsat based tree cover predictor used here.

As seen here, high spatial resolution wall-to-wall map predictions can be achieved with sufficient LiDAR coverage and use of fine-scale residual spatial structure. The G-LiHT LiDAR data—spatially dense along the 180 m swath widths—could better inform canopy height prediction across the TIU if it covered a larger swath width. This could be accomplished by increasing the flight altitude. While a higher nominal flying altitude will

increase the swath width, it will also decrease the spatial density of LiDAR observations. Our results suggest that LiDAR density is less important than coverage width, given models were fit using only  $\sim 17\%$  ( $5 \times 10^6 / 28,751,400$ ) of available data and even then it appears we had ample information to inform prediction within flight lines. This observation has implications for the other LiDAR collection campaigns, for example, ICESat-2 (Abdalati et al. 2010; ICESat-2 2015) and Global Ecosystem Dynamics Investigation LiDAR GEDI (2014), when they choose between pulse density and swath width.

#### 4. Summary

Our aim has been to propose alternate formulations and derivatives of Bayesian NNGP models developed by Datta et al. (2016a) to substantially improve computational efficiency for fully process-based inference. These improvements make it feasible to bring a rich set of hierarchical spatial Gaussian process models to bear on data intensive analyses such as the TIU forest canopy mapping effort. Analysis of simulated data shows that compared with the sequential specification of Datta et al. (2016a), the response and collapsed models offer improved MCMC chain behavior for the intercept and spatial covariance parameters. If full inference about the spatial random effects is of interest, then the response or conjugate models are not appropriate. So while the collapsed model can be computationally intensive, depending on the burden imposed by the sparse Cholesky decomposition, it is the only fully Bayesian alternative to the sequential Gibbs sampler developed in Datta et al. (2016a) and should generally be selected over the latter due to its significantly improved chain convergence. Furthermore, recent work by Katzfuss and Guinness (2017) shows that the collapsed model provides a better approximation of the full GP than the response model in the sense of Kullback–Leibler divergence from the full GP model. If model parameter estimation and/or spatial interpolation of the response is the primary objective, the response model offers substantial computational gains over the collapsed model. Finally, relative to the other NNGP models, the conjugate model delivers massive gains in computational efficiency and seemingly uncompromised predictive inference, but requires specification of the models' spatial decay and  $\alpha$  parameters. However, as demonstrated in the simulation and TIU analyses, these parameters can be effectively selected via cross-validation. The response and conjugate NNGP models are available for public use in the `spNNGP` package (Finley, Randin, and Korner 2017) in R.

The response model emerges a viable option for obtaining full Bayesian inference about spatial covariance parameters and prediction units. A fully Bayesian kriging model capable of handling  $5 \times 10^6$  observations on standard computing architectures is an exciting advancement and opens the door to using a rich set of process models to tackle complex problems in big data settings. For example, the response and collapsed NNGP models can seamlessly replace GP within multivariate, space-varying coefficients, and space-time settings (see, e.g., Datta et al. 2016a, 2016c, 2016b). The conjugate model provides a new tool for delivering fast interpolation with few inferential concessions. Extension of the conjugate model to some of the more complex

hierarchical frameworks noted above provides an additional avenue for development.

The TIU analysis shows the advantage of embedding the NNGP as a sparsity-inducing prior within a hierarchical modeling framework. The proposed NNGP specifications yield complete coverage forest canopy height prediction maps with associated uncertainty estimates using sparsely sampled but locally dense  $n = 5 \times 10^6$  LiDAR data. The resulting data product is the first statistically robust map of forest canopy for the TIU. Insight into residual spatial dependence will help guide planning for upcoming LiDAR data collection campaigns at global and local scales to improve prediction by leveraging information in more optimally located canopy height observations.

There remains much to be explored in NNGP models. Recent investigations by Guinness (2018) suggest that the Kullback–Leibler divergence between full Gaussian process likelihoods and Vecchia-type nearest neighbor approximations can be sensitive to topological ordering. Our preliminary explorations seem to suggest that while the Kullback–Leibler divergence from the truth may be affected, substantive inference in the form of parameter estimates and predictive performance (based upon root-mean-square-predictions) are very robust. Guinness (2018) also demonstrated empirically that certain carefully chosen orderings of the locations lead to a better approximation of the full GP by NNGP, than what is achieved by the simple co-ordinate based ordering. All the algorithms, we propose here are flexible to the choice of ordering. While, we have continued to use co-ordinate based ordering for all the data analysis here, we could as easily use any of the orderings proposed by Guinness (2018). We are currently conducting further investigations with the ordering suggested by Guinness (2018) and intend to report on our findings in a subsequent work.

A limiting factor for the hybrid approach adopted in the conjugate NNGP model is the cross-validation procedure for selecting the hyper-parameters. For most spatial applications, the isotropic Matérn functions are often the preferred choice for the covariance kernel, and is convenient for implementing the conjugate model as it only involves two or three unknown parameters. Hence, cross-validation using a grid search on a three dimensional space is computationally feasible. However, as pointed out by one reviewer, many other GP-based applications use more complex covariance functions involving several parameters. For example, in computer model emulations, separable Gaussian covariance functions are commonly used, for which there is a co-ordinate specific range parameter. As with all cross-validation based procedures, the conjugate model will also suffer from the curse of dimensionality in such richly parametrized settings, as searching for optimal or near-optimal points in a high-dimensional space is highly inefficient. Newer strategies need to be conceived for hyper parameter estimation in such settings.

Another pertinent matter concerns the performance of NNGP models for nonstationary processes. Naive implementations using neighbor selection based on simple Euclidean metrics may not be desirable. Here, the dynamic neighbor-finding algorithms proposed by Datta et al. (2016c) in spatiotemporal contexts may offer a better starting point than finding suitable metrics to choose neighbors. Still, work needs to be done in developing and analyzing analogous algorithms

for nonstationary processes. Finally, there is scope to explore NNGP models for high-dimensional multivariate outcomes using spatial factor models (Taylor-Rodriguez et al. 2018) or graphical Gaussian models and assessing their efficiency for highly complex multivariate spatial datasets.

## Supplementary Materials

Supplementary information provides additional comparative analyses of NNGP specifications versus full GP models and alternative GP approximations.

## Funding

Finley was supported by National Science Foundation (NSF) DMS-1513481, EF-1137309, EF-1241874, and EF-1253225. Cook, Morton, and Finley were supported by NASA Carbon Monitoring System grants. Banerjee was supported by NSF DMS-1513654, NSF IIS-1562303, and NIH/NIEHS R01-ES027027.

## References

- Abdalati, W., Zwally, H., Bindschadler, R., Csatho, B., Farrell, S., Fricker, H., Harding, D., Kwok, R., Lefsky, M., Markus, T., Marshak, A., Neumann, T., Palm, S., Schutz, B., Smith, B., Spinhirne, J., and Webb, C. (2010), “The ICESat-2 Laser Altimetry Mission,” *Proceedings of the IEEE*, 98, 735–751. [12]
- AICC (2016), “Fire History in Alaska,” available at [http://afsmaps.blm.gov/imf\\_firehistory/imf.jsp?site=firehistory](http://afsmaps.blm.gov/imf_firehistory/imf.jsp?site=firehistory). [10]
- Amestoy, P. R., Davis, T. A., and Duff, I. S. (1996), “An Approximate Minimum Degree Ordering Algorithm,” *SIAM Journal on Matrix Analysis and Applications*, 17, 886–905. [5]
- (2004), “Algorithm 837: AMD, An Approximate Minimum Degree Ordering Algorithm,” *ACM Transactions on Mathematical Software*, 30, 381–388. [8]
- Banerjee, S. (2017), “High-Dimensional Bayesian Geostatistics,” *Bayesian Analysis*, 12, 583–614. [1]
- Banerjee, S., Gelfand, A. E., Finley, A. O., and Sang, H. (2008), “Gaussian Predictive Process Models for Large Spatial Datasets,” *Journal of the Royal Statistical Society, Series B*, 70, 825–848. [1]
- Barbian, M. H., and Assunção, R. M. (2017), “Spatial Subsample Estimator for Large Geostatistical Data,” *Spatial Statistics*, 22, 68–88. [1]
- Chen, Y., Davis, T. A., Hager, W. W., and Rajamanickam, S. (2008), “Algorithm 887: CHOLMOD, Supernodal Sparse Cholesky Factorization and Update/Downdate,” *ACM Transactions on Mathematical Software*, 35, 1–14. [8]
- Cook, B., Corp, L., Nelson, R., Middleton, E., Morton, D., McCorkel, J., Masek, J., Ranson, K., Ly, V., and Montesano, P. (2013), “NASA Goddard’s LiDAR, Hyperspectral and Thermal (G-LiHT) Airborne Imager,” *Remote Sensing*, 5, 4045–4066. [9,10]
- Cressie, N., and Johannesson, G. (2008), “Fixed Rank Kriging for Very Large Spatial Data Sets,” *Journal of the Royal Statistical Society, Series B*, 70, 209–226. [1]
- Dagum, L. and Menon, R. (1998), “OpenMP: An Industry Standard API for Shared-Memory Programming,” *Computational Science & Engineering, IEEE*, 5, 46–55. [8]
- Datta, A., Banerjee, S., Finley, A. O., and Gelfand, A. E. (2016a), “Hierarchical Nearest-Neighbor Gaussian Process Models for Large Geostatistical Datasets,” *Journal of the American Statistical Association*, 111, 800–812. [1,2,3,5,6,7,8,12]
- (2016b), “On Nearest-Neighbor Gaussian Process Models for Massive Spatial Data,” *Wiley Interdisciplinary Reviews: Computational Statistics*, 8, 162–171. [12]
- Datta, A., Banerjee, S., Finley, A. O., Hamm, N. A. S., and Schaap, M. (2016c), “Nonseparable Dynamic Nearest Neighbor Gaussian Process Models for Large Spatio-Temporal Data With an Application to

- Particulate Matter Analysis,” *Annals of Applied Statistics*, 10, 1286–1316. [2,12,13]
- Davis, T. A. (2006), *Direct Methods for Sparse Linear Systems*, Philadelphia, PA: Society for Industrial and Applied Mathematics. [8]
- (2016a), “A Suite of Sparse Matrix Software,” available at [www.suitesparse.com](http://www.suitesparse.com). [8]
- (2016b), “User Guide for CHOLMOD: A Sparse Cholesky Factorization and Modification Package,” available at [www.suitesparse.com](http://www.suitesparse.com). [8]
- Finley, A., Datta, A., and Banerjee, S. (2017), *spNNGP: Spatial Regression Models for Large Datasets Using Nearest Neighbor Gaussian Processes*, r package version 0.1.1. [12]
- Finley, A. O., Sang, H., Banerjee, S., and Gelfand, A. E. (2009), “Improving the Performance of Predictive Process Modeling for Large Datasets,” *Computational Statistics & Data Analysis*, 53, 2873–2884. [1]
- Finney, M. A. (2004), “FARSITE: Fire Area Simulator-Model Development and Evaluation,” Technical Report Research Paper RMRS-RP-4, U.S. Department of Agriculture, Forest Service, Rocky Mountain Research Station. [1]
- Furrer, R. (2016), *spam: SPArse Matrix*, r package version 1.4-0. [1]
- Furrer, R., and Sain, S. R. (2010), “spam: A Sparse Matrix R Package with Emphasis on MCMC Methods for Gaussian Markov Random Fields,” *Journal of Statistical Software*, 36, 1–25. [1]
- GEDI (2014), “Global Ecosystem Dynamics Investigation LiDAR,” available at <http://science.nasa.gov/missions/gedi/>. [12]
- Gelman, A., and Rubin, D. (1992), “Inference from Iterative Simulation Using Multiple Sequences,” *Statistical Science*, 7, 457–511. [9]
- Gerber, F. (2017), *gapfill: Fill Missing Values in Satellite Data*, r package version 0.9.5. [1]
- Gneiting, T., and Raftery, A. E. (2007), “Strictly Proper Scoring Rules, Prediction, and Estimation,” *Journal of the American Statistical Association*, 102, 359–378. [8]
- Gramacy, R. B. (2016), “laGP: Large-Scale Spatial Modeling via Local Approximate Gaussian Processes in R,” *Journal of Statistical Software*, 72, 1–46. [1,7]
- Gramacy, R. B., and Apley, D. W. (2015), “Local Gaussian Process Approximation for Large Computer Experiments,” *Journal of Computational and Graphical Statistics*, 24, 561–578. [1]
- Gramacy, R. B., and Sun, F. (2017), *laGP: Local Approximate Gaussian Process Regression*, r package version 1.5-1. [7]
- Guhaniyogi, R., and Banerjee, S. (2018), “Meta-Kriging: Scalable Bayesian Modeling and Inference for Massive Spatial Datasets,” *Technometrics*, 60, 430–444. [1]
- Guinness, J. (2018), “Permutation and Grouping Methods for Sharpening Gaussian Process Approximations,” *Technometrics*, 60, 415–429. [13]
- Hager, W. W. (2002), “Minimizing the Profile of a Symmetric Matrix,” *SIAM Journal on Scientific Computing*, 23, 1799–1816. [5]
- Hansen, M. C., Potapov, P. V., Moore, R., Hancher, M., Turubanova, S. A., Tyukavina, A., Thau, D., Stehman, S. V., Goetz, S. J., Loveland, T. R., Kommareddy, A., Egorov, A., Chini, L., Justice, C. O., and Townshend, J. R. G. (2013), “High-Resolution Global Maps of 21st-Century Forest Cover Change,” *Science*, 342, 850–853. [10]
- Heaton, M. J., Datta, A., Finley, A. O., Furrer, R., Guhaniyogi, R., Gerber, F., Gramacy, R. B., Hammerling, D., Katzfuss, M., Lindgren, F., Nychka, D. W., Sun, F., and Zammit-Mangion, A. (2017), “Methods for Analyzing Large Spatial Data: A Review and Comparison,” ArXiv e-prints, available at <https://arxiv.org/abs/1710.05013>. [1]
- Huang, H., and Sun, Y. (2018), “Hierarchical Low Rank Approximation of Likelihoods for Large Spatial Datasets,” *Journal of Computational and Graphical Statistics*, 27, 110–118. [2]
- Hurt, G. C., Dubayah, R., Drake, J., Moorcroft, P. R., Pacala, S. W., Blair, J. B., and Fearon, M. G. (2004), “Beyond Potential Vegetation: Combining Lidar Data and a Height-Structured Model for Carbon Studies,” *Ecological Applications*, 14, 873–883. [1]
- ICESat-2 (2015), “Ice, Cloud, and Land Elevation Satellite-2,” available at <http://icesat.gsfc.nasa.gov/>. [12]
- Karypis, G., and Kumar, V. (1998), “A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs,” *SIAM Journal on Scientific Computing*, 20, 359–392. [5,8]
- Katzfuss, M. (2017), “A Multi-Resolution Approximation for Massive Spatial Datasets,” *Journal of the American Statistical Association*, 112, 201–214. [1]
- Katzfuss, M., and Guinness, J. (2017), “A General Framework for Vecchia Approximations of Gaussian Processes,” available at <https://arxiv.org/pdf/1708.06302.pdf>. [12]
- Klein, T., Randin, C., and Korner, C. (2015), “Water Availability Predicts Forest Canopy Height at the Global Scale,” *Ecology Letters*, 18, 1311–1320. [1]
- Lauritzen, S. L. (1996), *Graphical Models*, Oxford: Clarendon Press. [2]
- Lefsky, M. A. (2010), “A Global Forest Canopy Height Map from the Moderate Resolution Imaging Spectroradiometer and the Geoscience Laser Altimeter System,” *Geophysical Research Letters*, 37, 115401. [1]
- Liu, J. S., Wong, W. H., and Kong, A. (1994), “Covariance Structure of the Gibbs Sampler with Applications to the Comparisons of Estimators and Augmentation Schemes,” *Biometrika*, 81, 27–40. [3]
- Murphy, K. (2012), *Machine Learning: A Probabilistic Perspective*, Cambridge, MA: The MIT Press. [2]
- Nychka, D., Bandyopadhyay, S., Hammerling, D., Lindgren, F., and Sain, S. (2015), “A Multiresolution Gaussian Process Model for the Analysis of Large Spatial Datasets,” *Journal of Computational and Graphical Statistics*, 24, 579–599. [1]
- Ra, S. W., and Kim, J. K. (1993), “Fast Mean-Distance-Ordered Partial Codebook Search Algorithm for Image Vector Quantization,” *IEEE Transactions on Circuits and Systems II*, 40, 576–579. [8]
- Rue, H., Martino, S., Lindgren, F., Simpson, D., Riebler, A., Krainski, E. T., and Fuglstad, G.-A. (2017), *INLA: Bayesian Analysis of Latent Gaussian Models Using Integrated Nested Laplace Approximations*, r package version 17.06.20. [1]
- Sang, H., Jun, M., and Huang, J. Z. (2011), “Covariance Approximation for Large Multivariate Spatial Data Sets with an Application to Multiple Climate Model Errors,” *The Annals of Applied Statistics*, 5, 2519–2548. [1]
- Stein, M. L., Chi, Z., and Welty, L. J. (2004), “Approximating Likelihoods for Large Spatial Data Sets,” *Journal of the Royal Statistical Society, Series B*, 66, 275–296. [2]
- Stratton, R. D. (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. [1]
- Stroud, J. R., Stein, M. L., and Lysen, S. (2017), “Bayesian and Maximum Likelihood Estimation for Gaussian Processes on an Incomplete Lattice,” *Journal of Computational and Graphical Statistics*, 26, 108–120. [2]
- Sun, Y., Li, B., and Genton, M. (2011), “Geostatistics for Large Datasets,” In *Advances and Challenges in Space-time Modelling of Natural Events*, eds. J. Montero, E. Porcu, and M. Schlather, Berlin: Springer-Verlag, pp. 55–77. [1]
- Taylor-Rodriguez, D., Finley, A. O., Datta, A., Babcock, C., Andersen, H.-E., Cook, B. D., Morton, D. C., and Banerjee, S. (2018), “Spatial Factor Models for High-Dimensional and Large Spatial Data: An Application in Forest Variable Mapping,” ArXiv:1801.02078. [13]
- Vecchia, A. V. (1988), “Estimation and Model Identification for Continuous Spatial Processes,” *Journal of the Royal Statistical Society, Series B*, 50, 297–312. [2,5]
- Zammit-Mangion, A., and Cressie, N. (2017), “FRK: An R Package for Spatial and Spatio-Temporal Prediction with Large Datasets,” arXiv preprint arXiv:1705.08105. [1]
- Zhang, H. (2004), “Inconsistent Estimation and Asymptotically Equal Interpolations in Model-Based Geostatistics,” *Journal of the American Statistical Association*, 99, 250–261. [6]
- Zhang, X. (2016), “An Optimized BLAS Library Based on GotoBLAS2,” available at <https://github.com/xianyi/OpenBLAS/>. [7]