Statistica Si	nica Preprint 1 2-2018-0005
Title	Spatial Factor Mode for High-Dimen and Large Spatial Data: An Appnear Mapping
Manuscript ID URL	SS-2018-0005 http://www.stat.sinica.edu.tw/stationi
DOI	10.5705/ss.202018.0005
Complete List of Authors	Daniel Taylor-Franciquez Andrew O. Finle Abhirup atta Thad a thock Ha zrn, en ce Cook cook and b Banerjee
Correspond	Andrew O. Finley
	finleya@msu.edu
vera subje	ct to English editing.

Spatial Factor Models for High-Dimensional and Large Spatial Data: An Application in Forest Variable Map

Daniel Taylor-Rodriguez¹, Andrew O. Finley^{2,*}, Abhirup Datta³, Chad Bebcock⁴ Hans-Erik Andersen⁵, Bruce D. Cook⁶, Douglas G. Hans-Erik Andersen⁵, Sudipto I. Vrjeg

¹Portland State University, ²Michigan State University, ³Johns H

⁴University of Washington, ⁵United States 10.

⁶National Aeronautics and Space Administration.

⁷University of California Los Angeles, * Co.

Abstract: Gathering information about forest expensive and arduous activity. Therefore, directly colled in d to produce high-resolution e data xt-generation collection initiatives maps over large spatial do ng (LiDAR) data are specifically aimed for remotely sensed light d na over large spatial domains. Given that Liat produci compl DAR data an s are often strongly correlated, it is possible to del, dict, and map forest variables over regions of interuse the with high-dimensional ($\sim 10^2$) spatially dependent LiDAR large number of locations ($\sim 10^5 - 10^6$). With this in mind, we spanal factor nearest neighbor Gaussian process (SF-NNGP) model, which we ed in a two-stage approach that connects the spatial structure found gnals with forest variables. We provide a simulation experiment that demonstrates the inferential and predictive performance of the SF-NNGP, and use the two-stage modeling strategy to generate complete-coverage maps of the forest variables, with associated uncertainty, over a large region of boreal forests in

interior Alaska.

Key words and phrases: LiDAR data, forest outcomes, nearest neighbor Gaussian processes, spatial prediction.

1. Introduction

Strong relationships between remotely sensed ection and rang DAR) data and forest variables have been do umented in the ner et al., 2009; Babcock et al., 2013; Næsset, rested settings, LiDAR data provide a high-dimensional signal that cterizes the vertical structure of the forest canopy at point-refe ons. Traditionally, LiDAR data acquisition camp s have sought complete-coverage at a high spatial resolution over relatively s al domains, resulting in a fine grid of point-referenced LiDNR ch settings, the link between gnals the LiDAR d ta and the f easurements on sparsely sampled forest inventory plots has lond to create high-resolution completecoverage predic forest variables. Commonly, this link is first the relevant features of the high-dimensional established l rac Machine Resion - reduction step (Babcock et al., 2015; Junt-LiDA (). Then the LiDAR features are used as predictors in a to explain the variability in the spatially coinciding forest regression variable es. Lastly, the model is applied to predict the forest outcomes at all locations across the domain where LiDAR signals have been observed.

Considerably more ambitious next-generation LiDAR collection initia-

tives, such as ICESAT-2 (ICESat-2, 2015), Global Ecosystem Dynamics Investigation LiDAR (GEDI) (GEDI, 2014), and NASA Goddard's LiDAR, Hyper-Spectral, and Thermal imager (G-LiHT) (G-LiHT, 2016), seek to quantify and map forest variables over vast spatial extent. To fulfill their goals in a cost-effective manner, these data-gathering programs do no colhey sparsely sai lect LiDAR data over the entire domain. In cations across the domain extent and over for st inventory p forest variables have been measured). While verage high-resolution maps of forest outcomes remains the primary interuse of these data, there is also interest in creating maps of ver nonsampled locations and assessing the spatial dependence within and among LiDAR signals.

variable prediction and map-Our motivating application foct s on fping in the real forests l using sparsely sampled LiDAR and forest variable measur m these regions, acquiring completecoverage LiDA data ive from a cost perspective (Andersen et al., Vet al., 2012). Because generating complete-**43**. 2011; Bolton et al. cover**o**re maj rariables (and perhaps LiDAR signals) is still the sampled LiDAR must be leveraged to inform the forest goal, the ions. One attractive solution is to move the LiDAR predictor varia to the left-hand side of the regression and then to model with the forest outcomes. When the number of LiDAR and forest variables is small, such joint models are possible via linear models of coregionalization; for example, see, Babcock et al. (2017) and Finley et al.

(2014a). Alternatively, if the LiDAR signal is high-dimensional, but observed at a small number of locations, reduced-rank models can be employed. For example, Banerjee et al. (2008), Ren and Banerjee (2013), and Finley et al. (2017) applied a reduced-rank predictive process modeling trategy to analyze similar high-dimensional data. However, such approaches cannot cale to data sets with tens of thousands of location and yield poor provide very performance (Stein, 2014).

Models able to handle high-dimensional signature e number of locations and capable of estimating within and among locations dependence structures are needed. Recent modeling dev lewed in Heaton et al. (2017) and Banerjee (2017) highlight several options for robust and practical approximation of univaria n process (GP) models. A subset of these models can be e ly ext d to accommodate relatively ss) for example, see (Datta et al., small multivariate respons ticu application, we require an approach 2016a). Nevertheless, for that can cope with b gh-dimensional LiDAR measurements, ~ 50 the large collection of observed locations. outcomes at a loc

The nearest stable Gaussian process (NNGP) developed in Datta et al. (2016a), that is left (2016b), and Datta et al. (2016c) can be used with a large of locations, because its scalability is not mediated by the number of converged ved locations, but rather by the size of the nearest neighbor see. These models belong to the class of methods that induce sparsity on the spatial precision matrix, and exploit the natural representation

of sparsity provided by graphical models (Lauritzen, 1996; Murphy, 2012) to build a sparse GP that accurately approximates the original dense GP.

To tackle the high-dimensional LiDAR data set, we develop a Bayesian NNGP spatial factor model (SFM), referred to as the SHNNGP. Following Christensen and Amemiya (2002), Hogan and Thernis (2004), an and Banerjee (2013), the SFM structure en proximating the dence between multivariate (spatially depenit) outcomes t dimensional set of spatial factors, alleviating irectly with high-dimensional outcomes. The SF-NNGP allows us to del and map the LiDAR signals on both observed and unob ns, and, conditioning on the LiDAR spatial signatures, we can similarly map the forest variables over the entire spatial domai st. Furthermore, using a Bayesian approach for model fittin nables p equip the derived estimates and predictions with asso uncertainty, an essential requirement of many high-profile ves. ur methods are fully implemented in C++, using B (B) al., 2001; Zhang, 2016) to leverage efficient op s and openMP (Dagum and Menon, 1998) to multiprocessor ma the algorithm through parallelization. mpre e key

The section 2 introduced by the remainder of this paper is as follows. Section 2 introduced by the control of the Section 3, we formulate the proposed hierarchical design modeling strategy. Section 4 presents an analysis of a symmetric section 5 we develop and the available LiDAR and forest inventory data, in Section 5, we develop and validate a predictive model for the forest variables. We close by providing

insights, recommendations, and directions for future in Section 6.

2. Data Description

The Bonanza Creek Experimental Forest (BCEF) is a Long-Term Ecological Research (LTER) site consisting of vegetation and landforms to call of interior Alaska. The BCEF is 21,000 has a section Tanana River floodplain along the southeatern borders (LTER, 2016). Figure 1 shows the location and the location and the detailed in this section.

Forest variables were collected on 197 plots in 2 me USDA Forest Service Forest Inventory and Anal Program protocol (Bechtold and Patterson, 2005). We consider three forest a commonly used by forest professionals to make management cisio ove-ground biomass (AGB); tree density (D); and ba he AGB for individual trees was Rath Method described in Woodall et al. estimated using the Com (2015). The T pressed in thousands of trees per hectare. The BA for lot he s of the individual trees' cross-sectional areas in Ad to a per hectare basis.

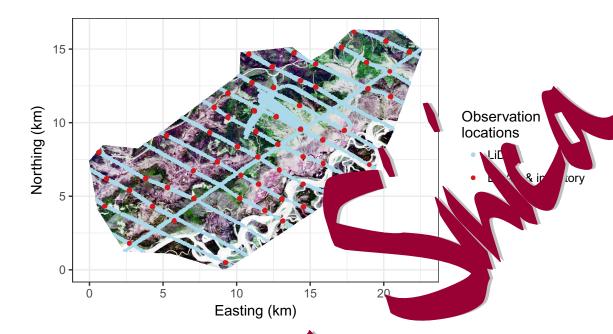


Figure 1: Bonanza Creek Experimental construct extent with color enhanced Landsat image and locations where the LiDAR signals were measured (Li-DAR in the legend) and location where the LiDAR signals and forest variables were measured LiDAR E envery P_g . In the legend).

ta were collected using a flight-line In the summer of 20 ASA Goddard's G-LiHT sensor (Cook et al., strip sampling 2013), which is a h abl unisensor system that accurately characterizes rtical distribution of canopy elements (Jakubowski terra et al., 2013). Point cloud information was summarized to et al., cell size to approximate field plot areas. Over each grid cell, $a 15 \times$ ns were generated by calculating the LiDAR return count psuedo-w densities for .5 m height bins between 0 and 28.5 m (i.e., 57 LiDAR outcomes per location). The LiDAR return count density for height bin l is defined as the number of returns in height bin l divided by the total number of LiDAR returns over the grid cell. Identical LiDAR psuedo-waveforms were obtained using point clouds extracted over each field plot. G-LiHT data for the study area are available online at https://gliht.gsfc.nasa.gov. For this analysis, 50,197 LiDAR observations were used for model-fitting.

A Landsat 8 top of atmosphere (TOA) reflectance product was product for the BCEF area for June of 2015. The Land 15 image was product to the June 2014 image owing to the excess a cloud cover it age. A tasseled cap transformation was applied as 3 TOA reflectance bands to obtain brightness, greenness, and wetness to cled cap indices (Baig et al., 2014). These indices are used as a green analysis.

Further details on the data set and our answer provided in Section 5.

3. Modeling Strateg

unce tainty-equipped predictions of forest Our goal is to model and variables using vined in LiDAR signals. Consider a LiDAR signal, $\mathbf{z}(\cdot)$, note collection of locations, $\mathcal{T}_z = \{\mathbf{s}_1, \dots, \mathbf{s}_{n_z}\},\$ serv Smes, $\mathbf{y}(\cdot)$, observed at locations in the set $\mathcal{T}_y =$ and a Furthermore, let $\mathcal{T}_{\emptyset} = \left\{\mathbf{t}_1, \dots, \mathbf{t}_{n_{\emptyset}}\right\}$ denote a set of loither LiDAR signals nor forest outcomes are available, but cations wr where r as are of interest. Thus, the set of locations where both Li-DAR and forest outcomes are mapped corresponds to $\mathcal{T} = (\mathcal{T}_z \cup \mathcal{T}_{\emptyset})$, with $\mathcal{T} \subset \mathcal{D} \subset \mathbb{R}^2$, where \mathcal{D} is the spatial domain of interest. Note that although $\mathbf{z}(\cdot)$ and $\mathbf{y}(\cdot)$ are "observed" at locations in \mathcal{T}_z and \mathcal{T}_y , respectively, we allow for missing values that are to be imputed in these sets. We make this distinction because locations where imputation is performed are part of the model fitting, whereas for locations in \mathcal{T}_{\emptyset} , predictions are driven expost facto from the posterior predictive distribution; see Section 3.4.

The LiDAR signals are high-dimensional variables of measurements h_z , whereas the forest outcomes are relatively mall-dimensional variables, $h_y \ll h_z$, assumed to have support on \mathbb{R}^{h_y} . The lipade LiDAR signals are strongly dependent on each other; LiDAR signals variable that the composition of a forest, and, as a plethora of example in the large demonstrated (Ene et al., 2018; Finley et al., 2014b; Nelson et al., 2017), variability in forest outcome variables can be writially explained by LiDAR characteristics.

3.1 Linking the LiDAR 11 Inventory Data

We seek to consect the form accomes and LiDAR signals as a two-step process. First, we form at the relative model to extract the spatial signature from the LiD. First a propositions in \mathcal{T}_z , which can also be used to interpolate LiDAR smalls \mathcal{T}_{\emptyset} . Along with other spatially referenced predictors, the LiDAR smalls region in \mathcal{T}_y are used as predictors to build the model from the forest outcomes. Moreover, a component that captures the model from exclusive to the forest outcomes can also be specified, if required. For $\mathbf{s} \in \mathcal{D}$, this two-stage model is given by

Stage 1:
$$\mathbf{z}(\mathbf{s}) = \mathbf{X}_z(\mathbf{s})'\boldsymbol{\beta}_z + \mathbf{w}^*(\mathbf{s}) + \boldsymbol{\varepsilon}_z(\mathbf{s}),$$
 (3.1)

Stage 2:
$$\mathbf{y}(\mathbf{s}) = \mathbf{X}_y(\mathbf{s})'\boldsymbol{\beta}_y + \boldsymbol{\Upsilon}\mathbf{w}^*(\mathbf{s}) + \mathbf{v}^*(\mathbf{s}) + \mathbf{v}_y(\mathbf{s}).$$
 (3.2)

Note that the influence of $\mathbf{z}(\mathbf{s})$ over $\mathbf{y}(\mathbf{s})$ in (3.2) in exerted solely \mathbf{v} ight its spatial component, $\mathbf{w}^*(\mathbf{s})$. There are set all arguments in factors approach, as opposed to substituting $\mathbf{z}(\mathbf{s})$ of $\mathbf{v}^*(\mathbf{s})'$ $\mathbf{v}^*(\mathbf{s})'$ in $\mathbf{v}^*(\mathbf{s})$ as covariates directly into (3.2). Among these, and most important for our setting, $\mathbf{z}(\mathbf{s})$, $\boldsymbol{\mu}_z(\mathbf{s})$ and $\mathbf{w}^*(\mathbf{s})$ are all high-dimensional objective $\mathbf{z}^*(\mathbf{s})$ reduces the dimensionality of the problem by casting it under the factor model structure, as shown in Section 3.2 and ddition, the elements within $\mathbf{z}(\mathbf{s})$ are strongly correlated, hence multic \mathbf{v}^* early issues would arise if it was included directly in (2).

and $\mathbf{X}_y(\mathbf{s})'\boldsymbol{\beta}_y$ capture large-scale In (3.1) and (3.2), the variation. For $\mathbf{k} \in \{1, 1\}$ presents a fixed $h_{\kappa} \times p_{\kappa}$ block-diagonal redictors, where $p_{\kappa} = \sum_{j=1}^{h_{\kappa}} p_{\kappa,j}$, having as matrix of spatial, lock r length- $p_{\kappa,j}$ vector $\mathbf{x}_j^{\kappa}(\mathbf{s})'$. The length- p_{κ} vector $\boldsymbol{\beta}_{\kappa}$ ts j diago regression coefficients associated with $X_{\kappa}(s)'$. The vectors corresp are h_z - and h_y -dimensional zero-centered stochastic processes over \mathcal{D} , respectively. ely. The process $\mathbf{w}^{\star}(\mathbf{s})$ captures the spatial variation of $\mathbf{z}(\mathbf{s})$, mesizes additional spatial variation in the forest outcomes. The $h_y \times h_z$ matrix Υ connects the spatial information extracted from the LiDAR model into the forest outcomes model. The vectors $\boldsymbol{\varepsilon}_z(\mathbf{s}) \sim N_{h_z}(\mathbf{0}, \boldsymbol{\Psi}_z)$ and

 $\varepsilon_y(\mathbf{s}) \sim N_{h_y}(\mathbf{0}, \Psi_y)$ represent uncorrelated random errors (i.e., Ψ_z and Ψ_y are diagonal) at finer scales.

Implementing this modeling strategy directly is challenging owing to the high-dimensionality of the LiDAR signals ($h_z \sim 50$) and the massive number of spatially dependent observations ($n \sim 10^5$). Thus, it is impossibly to attempt using common computing resources the following section we formulate a viable alternative to models (3.1 and (3.2).

3.2 The Spatial Factor NNGP Model

To make models (3.1) and (3.2) tractable with limited power, we combine a dimension-reduction approach and a sparsity-inducing technique. In particular, we introduce the SF-NNGI in the which brings together the SFM structure (Schmidt and Gelfan), 2003 level et al., 2008; Zhang, 2007; Ren and Bar rjee, 2013) and the structure (Schmidt and Gelfan), 2005 level et al., 2016b,c,a).

While the SFM strug able the analysis of high-dimensional reproblemations of a relatively small number of sponse vectors v usi c p. independent stock s, NNGPs make it possible to fit spatial promber of spatial observations is particularly large. cess odels NNGPs te the parent (dense) GP using the natural representation ided by graphical models (Lauritzen, 1996; Murphy, 2012), by assumi ditional independence—where conditioning is on the nearest with locations outside of the neighbor set. The result is a proper (but sparse) GP that accurately approximates the original dense GP. In contrast to other sparsity-inducing approaches, NNGPs allow for interpolation at unobserved locations and can be used to make full inference on model parameters, including the latent processes. Combining the SFM structure with NNGPs provides a methodology capable of coping simultaneously with high-dimensional response vectors and a large number of spatially dependent observations.

Under the traditional SFM structure, spatially indence is introd by defining the spatial process as $\mathbf{w}^*(\mathbf{s}) = \Lambda \mathbf{w} / \sim \mathrm{GP}(\mathbf{0}, \mathcal{H}(\cdot))$, is a factor loadings matrix (commonly tall an small-dimensional vector of independent spatial GPs, providing the non-parable multivariate cross-covariance function given by

$$\mathcal{H}(\mathbf{h} \ \phi) = \operatorname{cov}(\mathbf{\Lambda} \mathbf{w} \ \mathbf{s}, \quad \mathbf{w}(\mathbf{s} + \mathbf{h}))$$
$$= \sum_{q_w} C_k(\mathbf{h} \ | \quad \mathbf{k} \lambda_k', \tag{3.3}$$

for locations $\mathbf{s}, \mathbf{s} + \mathbf{h} \in \mathcal{D}$ is $\mathcal{C}_k (\varphi_k)$ denotes a univariate parametric correlation function, \mathcal{A} denotes \mathbf{h} the kth column of $\mathbf{\Lambda}$. This cross-covariance matrix is induced \mathcal{A} -vertex $\mathbf{A} \leq l$) spatial factors $\mathbf{w}(\mathbf{s})$ with independent components $\mathbf{a} = \mathbf{v} \cdot (0, \mathcal{C}_k(\cdot \mid \phi_k))$.

As so m, \mathbf{v} 's (3.1) and (3.2) can be reformulated as SF-NNGPs by he spatial processes $\mathbf{w}^*(\mathbf{s})$ and $\mathbf{v}^*(\mathbf{s})$ as

$$\mathbf{w}^{\star}(\mathbf{s}) = \mathbf{\Lambda}_z \mathbf{w}(\mathbf{s}) \text{ and } \mathbf{v}^{\star}(\mathbf{s}) = \mathbf{\Gamma} \mathbf{v}(\mathbf{s}),$$
 (3.4)

where the matrices $\Lambda_z = ((\lambda_{hk}^{(z)}))_{h_z \times q_w}$ and $\Gamma = ((\gamma_{lr}))_{h_y \times q_v}$ correspond to the

factor loadings matrices, and the new spatial factors for $\mathbf{s} \in \mathcal{D}$ are given by

$$\mathbf{w}(\mathbf{s}) \sim \prod_{k=1}^{q_w} \text{NNGP}\left(0, \tilde{\mathcal{C}}(\cdot \mid \phi_k^w)\right), \text{ and}$$

$$\mathbf{v}(\mathbf{s}) \sim \prod_{r=1}^{q_v} \text{NNGP}\left(0, \tilde{\mathcal{C}}(\cdot \mid \phi_r^v)\right).$$

The expressions NNGP $\left(0, \tilde{\mathcal{C}}(\cdot \mid \phi_k^w)\right)$ and NN $\left(0, \tilde{\mathcal{C}}(\cdot \mid \phi_r^v)\right)$ do to derived from the parent processes $\operatorname{GP}(0, \mathcal{C}(\cdot \mid \phi_r^w))$ and $\operatorname{GP}(0, \cdot v^v)$, respectively. Here, $\mathcal{C}(\cdot \mid \phi)$ represents the spatial correlation function the spatial decay parameter ϕ . The factor model representation in f(x) and so a significant reduction in the dimensionality of the problem because the spatial factors $\mathbf{w}(\mathbf{s}) = (w_k(\mathbf{s}) \mid 1 \leq k \leq q_w)$ and $\mathbf{v} = (v_r(\mathbf{s}) : 1 \leq r \leq q_v)$ have dimensions $q_w << h_z$ and $q_v \leq h_y$, respect to

Combining these elements, we letting $\Lambda_y = \Upsilon \Lambda_z = ((\lambda_{lk}^{(y)}))_{h_y \times q_w}$, a computationally viable version of (a. (3.2) is

Stage
$$\mathbf{s})'\boldsymbol{\beta}_z + \boldsymbol{\Lambda}_z \mathbf{w}(\mathbf{s}) + \boldsymbol{\varepsilon}_z(\mathbf{s})$$
 (3.5)

St
$$\mathbf{X}_y(\mathbf{s})'\boldsymbol{\beta}_y + \boldsymbol{\Lambda}_y \mathbf{w}(\mathbf{s}) + \boldsymbol{\Gamma} \mathbf{v}(\mathbf{s}) + \boldsymbol{\varepsilon}_y(\mathbf{s}).$$
 (3.6)

identifiable terson, 2003). Identifiability for SFMs can be achieved either upper triangle of the loadings matrix equal to zero and its diagonal elements all equal to one (Geweke and Zhou, 1996; Lopes and West, 2004; Aguilar and West, 2010), or, as in Ren and Banerjee (2013), by fixing

the sign of one element in each column of the factor loadings matrix, while enforcing an ordering constraint among the spatial decay parameters of the univariate correlation functions. We choose to ensure rotation and scale identifiability by using the former approach.

With the SFM structure in place, introducing the NNGP reduce the expensive ($\sim n_z^3 q_w$ and $\sim n_y^3 q_v$) calculation require to invert the decree-variance matrices from the parent GPs by n_z , and $n_y q_v$ parallely, each of order m^3 . Here, m is the number of side of the NNGP, with $m \ll n_y \leq n_z$. In simulations, Datta et al. (20.17) found that, in most cases, $10 \leq m \leq 20$ provides an excellent on to the parent process; thus, the number of operations required is nearly linear in n.

For completeness, additional details of Scattering NGPs, and the sampling algorithm are included in the online upple of . For a more thorough treatment of SFM's, refer to Navarti term (2013) and Genton and Kleiber (2015), and for NNGPs, refer to Date et al. (2016c).

3.3 Prior Special Tal Hierarchical Formulation

Importantly, \mathbf{u} and \mathbf{v} and \mathbf{v} and \mathbf{v} are fitted separately such that $\mathbf{w}(\mathbf{s})$ exclusive approximates for $\mathbf{w}(\mathbf{s})$ (e.g., the posterior means) in (3.6) disregards the uncertainty present in the LiDAR spatial signal. Thus, to propagate this uncertainty through the forest outcome predictions, at each iteration of the Markov Chain Monte Carlo (MCMC) algorithm for $\mathbf{y}(\mathbf{s})$, we draw a sample for $\mathbf{w}(\mathbf{s})$ ($\mathbf{s} \in \mathcal{T}_y$) MCMC samples obtained when fitting model (3.5).

As mentioned in the previous section, the stochastic processes that capture the spatial structure are assumed to follow NNGPs. Given that an NNGP is a proper GP, at a finite collection of locations, the NNGPs induce zero-centered multivariate normal priors, with covariance natrices given by $\tilde{\mathbf{C}}^{(w)}$ and $\tilde{\mathbf{C}}^{(v)}$, respectively. Additionally, we use suitably noninform tive priors for all other parameters, thus providing that it is sampling strategies.

at or conjugat In particular, we assume that β is eithe matrices Γ and Λ_z are constrained as describ below the diagonal assumed to be standard normal. All elements in are also assumed to follow a standard normal distribution. ntries in Ψ_z and Ψ_y are assigned half-t priors. Lastly, we assume uniform priors for the elements of the spatial decay vector (), $\phi_{w,k}: 1 \leq k \leq q_w$) and $\phi_v = (\phi_{v,r} : 1 \le r \le q_v)$ in the $\log 0.05/\zeta_{\rm max}, -\log 0.01/\zeta_{\rm min}),$ erval where ζ_{\min} and ζ_{\max} are the naximum distances, respectively, ϕ_y are not conjugate with their across all locations. Given corresponding Veliho are sampled using random walk Metropolis steps.

To joint more insities for the first and second stages of the algo-

rithm are proportional to

Stage 1:

$$\pi(\boldsymbol{\phi}_{w}) \operatorname{N}_{n_{z}q_{w}}(\mathbf{w}_{\mathcal{T}_{z}}|\mathbf{0}, \tilde{\mathbf{C}}^{(w)}) \left(\prod_{k=1}^{q_{w}} \prod_{j>k}^{h_{z}} \operatorname{N}(\lambda_{jk}^{(z)}|0, \mathbf{0}) \right) \times \pi(\boldsymbol{\beta}_{z}) \left(\prod_{j=1}^{h_{z}} \mathcal{IG}(\psi_{j}^{z}|\nu/2, \nu/2, \nu/2, \mathbf{0}) \right) \times \left(\prod_{\mathbf{s}_{i} \in \mathcal{T}_{z}} \operatorname{N}_{h_{z}}(\mathbf{z}(\mathbf{s}_{i})|\mathbf{X}_{z}(\mathbf{s}_{i}) \right) \times \left(\prod_{\mathbf{s}_{i} \in \mathcal{T}_{z}} \operatorname{N}_{h_{z}}(\mathbf{z}(\mathbf{s}_{i})|\mathbf{X}_{z}(\mathbf{s}_{i}) \right)$$

$$(5.7)$$

Stage 2:

$$\pi(\boldsymbol{\phi}_{v}) \operatorname{N}_{n_{y}q_{v}}(\nabla \tau_{y}|\mathbf{0}, \tilde{\mathbf{C}}^{(v)}) \left(\prod_{k=1}^{q_{w}} \sum_{j=1}^{h_{v}} (\mathbf{v}_{ik}^{(y)}|\mathbf{0}, 1) \right) \left(\prod_{r=1}^{q_{v}} \prod_{j>r} \operatorname{N}(\gamma_{jr}|\mathbf{0}, 1) \right) \times \pi(\boldsymbol{\beta}_{y}) \left(\prod_{j=1}^{h_{y}} \mathcal{I}\mathcal{G}(\mathbf{v}_{i}^{(y)}|\mathbf{v}_{j}^{(y)}|\mathbf{v}_{j}^{(y)}) \mathcal{J}(a_{y,j}|1/2, 1/A^{2}) \right) \times \left(\prod_{\mathbf{s}_{i}} \prod_{j=1}^{h_{y}} \mathcal{I}(\mathbf{s}_{i}) \mathbf{Y}_{y}(\mathbf{s}_{i})^{j} \boldsymbol{\beta}_{y} + \boldsymbol{\Lambda}_{y} \mathbf{w}(\mathbf{s}_{i}) + \boldsymbol{\Gamma} \mathbf{v}(\mathbf{s}_{i}), \boldsymbol{\Psi}_{y}) \right), \quad (3.8)$$

where $\mathbf{w}_{\mathcal{T}_z} = \mathbf{v}(\mathbf{s}_i \mid \mathbf{s}_i \in \mathcal{T}_y)'$ and $\mathbf{v}_{\mathcal{T}} = (\mathbf{v}(\mathbf{s}_i)' : \mathbf{s}_i \in \mathcal{T}_y)'$, such that

$$\mathbf{N}_{q_w}(\mathbf{w}|\tilde{\mathbf{C}}^{(w)}) = \prod_{\mathbf{s}_i \in \mathcal{T}_z} \mathbf{N}_{q_w}(\mathbf{w}(\mathbf{s}_i) | \mathbf{B}_i^{(w)} \mathbf{w}_{N(i)}, \mathbf{F}_i^{(w)}), \text{ and}$$

$$\mathbf{N}_{n_w o} = \mathbf{T} |\mathbf{0}, \tilde{\mathbf{C}}^{(v)}) = \prod_{\mathbf{s}_i \in \mathcal{T}_y} \mathbf{N}_{q_v}(\mathbf{v}(\mathbf{s}_i) | \mathbf{B}_i^{(v)} \mathbf{v}_{N(i)}, \mathbf{F}_i^{(v)}). \tag{3.9}$$

The expressions on the right-hand side of (3.9) result from the construction of the NNGP (see online supplement). For an m-neighbor NNGP, let

 $m_i = \min\{m, i-1\}$ denote the number of neighbors for location \mathbf{s}_i . The index set N(i) for location $\mathbf{s}_i \in \mathcal{T}_z$ contains its m_i nearest neighbors; thus, $\mathbf{w}_{N(i)}$ corresponds to the vector $(\mathbf{w}(\mathbf{s}_j)': \mathbf{s}_j \in N(i) \subset \mathcal{T}_z)'$. The neighbor set for $\mathbf{v}(\mathbf{s}_i)$ is defined analogously. Letting $u \in \{w, v\}$, $\mathbf{B}_i^{(u)}$ depotes a $q_u \times m_i q_u$ block matrix, with the $q_u \times q_u$ diagonal blocks containing the kriging with on, $\mathbf{F}_{i}^{(u)}$ corresp for the q_u spatial factors for each neighbor. or the q_u spatia the $q_u \times q_u$ diagonal matrix with the variance details tioned on the neighbor set N(i) (see Section 8 on $\mathbf{B}_{i}^{(u)}$ and $\mathbf{F}_{i}^{(u)}$). Lastly, the parameters $\{a_{y,j}\}_{j=1}^{h_{y}}$ and $\{a_{z,k}\}_{k\neq j}^{h_{y}}$ mplete the hierarchical representation of the half-t prior distr d for ψ_i^y and ψ_k^z , respectively, and the hyperparameter A is simply chosen to be some large value (say, 100).

Owing to prior conjugacy, the following conditional densities for all parameters except ϕ_w and ϕ_v can be sampling algorithms be for the online supplement.

3.4 Imputation a relation

As mationed of the AR signals are collected over the large spatial region \mathcal{T}_z , where for outcome observations are confined to the smaller subset of rock. \mathcal{T}_y . Additionally, there are relevant out-of-sample locations where neither iDAR nor forest outcomes are observed, \mathcal{T}_{\emptyset} . Finally, there are no cations within the corresponding reference sets \mathcal{T}_z and \mathcal{T}_y that have some or all missing outcomes. It is thus essential for this modeling effort to provide the means to accurately impute the missing values in \mathcal{T}_z or

 \mathcal{T}_y . This enables us to generate LiDAR predictions in \mathcal{T}_{\emptyset} and forest outcome predictions within $\mathcal{T}_{\emptyset} \cup (\mathcal{T}_z \setminus \mathcal{T}_y)$. Given the NNGP formulation, both the imputation and the out-of-sample prediction are remarkably inexpensive.

Imputation is straightforward. Let $\mathbf{s}_{\bullet} \in \mathcal{T}_z$ be a location where $\mathbf{z}(\mathbf{s}_{\bullet})$ is missing. Then, $\mathbf{z}(\mathbf{s}_{\bullet})$ is drawn as part of the sampling algorithm from $N_{h_z}(\mathbf{X}_z(\mathbf{s}_{\bullet})'\boldsymbol{\beta}_z + \boldsymbol{\Lambda}_z\mathbf{w}(\mathbf{s}_{\bullet}), \boldsymbol{\Psi}_z)$, where $\mathbf{w}(\mathbf{s}_{\bullet})$ is the online supplied from the full ditional posterior density in Equation (S3.1) the online supplied in the full conditional posterior for $\mathbf{v}(\mathbf{s}_{\bullet})$, where $\mathbf{s}_{\bullet} \in \mathcal{T}_y$, the procedure $\mathbf{v}(\mathbf{s}_{\bullet})$ is a location where $\mathbf{z}(\mathbf{s}_{\bullet})$ is a loca

The procedure to predict a new LiDAR observation \mathcal{T}_{\emptyset} , begins by sampling the spatial factor $\mathbf{w}(\mathbf{s}_z)$ from $N_{q_w}(\mathbf{B}_{\circ}^{(w)}\mathbf{w}_{N(\mathbf{s}_{\circ})}, \mathbf{F}_{\circ}^{(w)})$, with $\mathbf{B}_{\circ}^{(w)}$ and $\mathbf{F}_{\circ}^{(w)}$ defined a before. Note that \mathbf{w} carest neighbor set $N(\mathbf{s}_{\circ})$ is assumed to be in \mathcal{T}_z . Then, \mathbf{w} draw $\mathbf{v} = \mathbf{z}_{\mathcal{T}_z}$ from $N_{h_z}(\mathbf{X}_z(\mathbf{s}_{\circ})'\boldsymbol{\beta}_z + \mathbf{\Lambda}_z\mathbf{w}(\mathbf{s}_{\circ}), \mathbf{\Psi}_z)$. This is decreased from the posterior samples of $\{\boldsymbol{\beta}_z, \mathbf{\Lambda}_z, \mathbf{\Psi}_z, \boldsymbol{\phi}_w\}$ obtained from the long algorithm.

To predict the force of $\mathbf{v}(\mathbf{s}_{\circ})$ at $\mathbf{s}_{\circ} \in \mathcal{T}_{\emptyset} \cup (\mathcal{T}_{z} \setminus \mathcal{T}_{y})$, we first generate samples of $\mathbf{v}(\mathbf{s}_{\circ}) \sim (\mathbf{E}_{\bullet} \mathbf{v}_{\Lambda_{\circ} \mathbf{s}_{\circ}}), \mathbf{F}_{\circ}^{(v)})$. Given that $\mathbf{y}(\mathbf{s}_{\circ})$ depends on $\mathbf{w}(\mathbf{s}_{\circ})$, we combine the laster draws of $\{\boldsymbol{\beta}_{y}, \boldsymbol{\Lambda}_{y}, \boldsymbol{\Gamma}, \boldsymbol{\Psi}_{y}, \boldsymbol{\phi}_{v}\}$ with those of $\mathbf{w}(\mathbf{s}_{\circ})$, obtained the dicting $\mathbf{z}(\mathbf{s}_{\circ})$, and draw predicted values for $\mathbf{y}(\mathbf{s}_{\circ}) \mid \mathbf{y}_{\mathcal{T}_{y}}$ from $\mathbf{v}(\mathbf{s}_{\circ}) \mid \mathbf{v}_{y} \mid \mathbf{v}_{y}$

4. Simulation: Recovering Low-dimensional Structure

In the following simulation exercise we focus exclusively on the high-dimensional component (i.e., the first stage) of the model described above. The simulation below was devised to illustrate the ability of our approach to recover the true low-dimensional structure when the data are generated from a low-dimensional SFM with dense spatial factors

We generate a synthetic data set for $h_z =$ outcomes in nObelio cations from the spatial factor model $\mathbf{z}(\mathbf{s}) = \mathbf{X}_z(\mathbf{s}) \beta_z + \mathbf{\Lambda}_z$ Here, $\mathbf{X}_z(\mathbf{s})'$ is a 50 × 150 block-diagonal matrix of predictors, and β vector of regression coefficients, both defined as before. We const rs $\tilde{\mathbf{w}}(\mathbf{s}) \sim \prod_{k=1}^{8} \mathrm{GP}(0, \mathcal{C}(\cdot \mid \tilde{\phi}_{k}^{z}),$ predictors for all outcomes. The spatial where $\mathcal{C}(\cdot | \tilde{\phi}_k^z)$ is an exponential correlation. n with decay parameter $\tilde{\phi}_k^z$. Additionally, for identifiable we as \blacksquare that the 50 × 8 factor loadings matrix A has zeros in ngle and ones along the diagonal. Finally, $\tilde{\boldsymbol{\varepsilon}}_z \sim N_{h_z} \left(\mathbf{0}, \tilde{\boldsymbol{\Psi}}_z \right)$ $= \operatorname{diag}(\psi_k^z : k = 1, \dots, 8).$

We assess the lemma (3.5) to recover the model parameters from the true data ones, or process, impute missing outcomes, and predict at out-of-wave (a 1.5). The SF-NNGP model was fitted for $q_w \in \{3, 5, 8, 10\}$ at n = 10 neighbors. Of the 10,000 locations, we assume all 50 outcomes to be missing in 200 locations chosen at random. These out-of-sample prediction and model validation.

The first result worth highlighting is the gains in computational efficiency

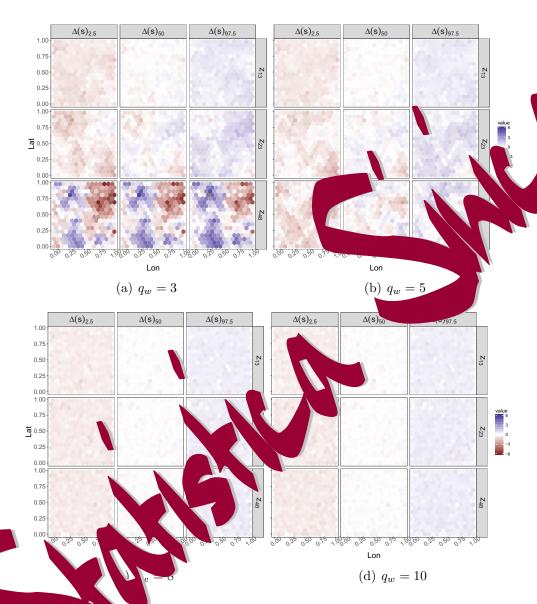
provided by the SF-NNGP. For this simulation exercise—a relatively computationally challenging problem—fitting the largest model considered (i.e., $q_w = 10$) with 50,000 MCMC iterations on a Linux server with an Intel i7 processor (two eight-core) and 16 GB of memory, the runtim was 4.88 hours. As shown below, the proposed approach is able to recover the true model parameters, accurately impute missing data, are tracely are precise precise precise all with suitable uncertainty estimates.

For all values of q_w , the SF-NNGP accul ression coefficients β_z for all predictors and responses (Figure 1 in the onsupplement). In contrast, the quality of the estimates for the variance components $\tilde{\psi}_k^z$'s was compromised when was lower than the true number of spatial factors. This shavior is expect of ver values of q_w , the ψ_k^z 's that the spatial component attempt to compensate for the admional with too few patial facto ture (Figure 2 in the supplemen- ψ_z , the coverage for $\tilde{\psi}_z$ was 88% and tary material). For q_w ose to $\tilde{\psi}_k^z$ with tight 95% credible sets. 84%, respective, with

When $q_w \neq 8$, to differ side s of the fitted Λ_z , ϕ_w , and $\mathbf{w}(\mathbf{s})$ do not match those of their \mathbf{w} of the true model. Therefore, to assess the quality of the fit for s_1 of s_2 of the true model. Therefore, to assess the quality of the fit for s_1 of s_2 of s_3 of the true model, give $\mathbf{w}^*(\mathbf{s}) = \mathbf{\Lambda}_z \mathbf{w}(\mathbf{s})$, for $\mathbf{s} \in \mathcal{T}_z$, to that of the true model, give $\mathbf{w}^*(\mathbf{s}) = \tilde{\mathbf{\Lambda}}_z \tilde{\mathbf{w}}(\mathbf{s})$.

cocations in \mathcal{T}_z , we calculate $\Delta(\mathbf{s}) = \mathbf{w}^*(\mathbf{s}) - \tilde{\mathbf{w}}^*(\mathbf{s})$ (fitted minus true spatial signal) for each MCMC draw of the parameters. For all $\mathbf{s} \in \mathcal{T}_z$, we obtained the median and 95% credible set for $\Delta(\mathbf{s})$. To facilitate visu-

alization, in Figure 2, we show the results for only three responses, selected at random, from the 50 considered. The columns of each panel map the quantiles 2.5, 50, and 97.5 for $\Delta(\mathbf{s})$, with three locations (13, 23, and 48) plotted by row. The fitted spatial signal when $q_w \in \{3, 1\}$ only partially recovers the true signal, with coverages of 26.13% and 42.06%, respectively, for $q_w = 3$ and $q_w = 5$. When $q_w \in \{8, 10\}$, the very of the spatial value at the extremely accurate: over all responses, the extrage is 94.78 in $q_w = 10$.



minus true spatial signal, $\Delta(\mathbf{s}) = \mathbf{w}^*(\mathbf{s}) - \tilde{\mathbf{w}}^*(\mathbf{s})$, for locations $\mathbf{s}_{13}, \mathbf{s}_{23}, \mathbf{s}_{48}$. In left to right, the columns in each panel show percentiles $\Delta(\mathbf{s})$, for $\Delta(\mathbf{s})$, respectively.

In addition to the previous results, it is also encouraging to find that when the dimension of the SF-NNGP model matches that of the true model, both the factor loadings $(\tilde{\Lambda}_z)$ and the spatial decay parameters $(\tilde{\phi}_z)$ from the true spatial process can be recovered accurately (Figures 3 and 4).

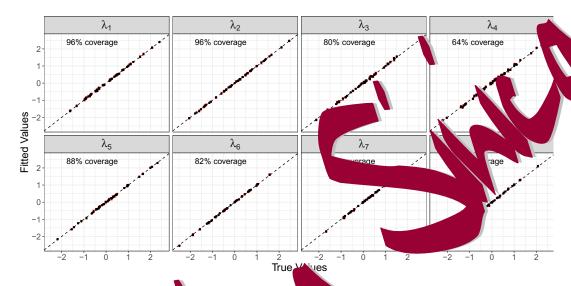


Figure 3: Fitted vs. true factor leadings ix parameters (95% credible sets and medians) for $q_w = 8$.

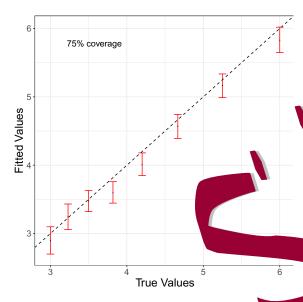


Figure 4: Fitted vs. true spatial decay parameters (see and medians) for $q_w = 8$.

Model performance in terms of the accuracy of imputation and prediction improves drastically as the number of forces approaches that of the true model; see Figures 5 and 6 in the organization applement.

Table 1 collapares v_w SE NNGP using different measures of outof-sample predictive energials. In particular, the continuous rank probability score (1908) (Legation (21) in Gneiting and Raftery, 2007) and the
root measures of prediction error (RMSPE) (Yeniay and Goktas, 2002)
when $q_w = 8$. The coverage of the 95% credible intervals of
the prediction was close to the nominal value for all q_w ; however, the width
of a rapidly decreases as q_w approaches the true number of spatial
factors.

Both the fitted values for the spatial signals and the out-of-sample pre-

Table 1: Out-of-sample prediction comparison across models with different numbers of spatial factors.

q_w	CRSP	RMSPE	95% Coverage	95% CI Width
3	0.85	1.61	95.82	6.14
5	0.67	1.28	95.43	4.79
8	0.45	0.83	94.78	3.10
10	0.45	0.83	94.84	3.10

dictions with $q_w = 8$ and $q_w = 10$ are practically n each other. Furthermore, the model with $q_w = 8$ accurately recovers all he true factor loadings (Figure 3). Interestingly, with $q_w = 10$ pection of the estimates for columns 1 through 6 in indicates that this model accurately estimates the corresponding true pa lues (see Figure 4 in the online supplement). However, i lel, the estimated parameter \mathbf{sam} values in columns 7 and 8 epartures from their true values. for all unconstrained elements in the 9th Furthermore, the 95% cre and 10th column ero (see Figure 3 in the online supplement). guid ce on the number of factors q_w to use. Because These results the model with $q_w = 10$ over that with $q_w = 8$ there ve accuracy or parameter fit, the results favor the more odel of the two. parsimon

5. Modeling LiDAR Signals and Forest Structure

Our focus in the subsequent analysis is to assess and interpret the utility of SF-NNGP spatial factors to explain the variability in the three forest outcomes defined in Section 2, measured on the BCEF. Following the two-stage model developed in Section 3.2, we fit (3.5) using $q_z \in \{1, 2, 3, 4, 6, 7, 8\}$ spatial factors and m = 10 neighbors to the second LiDAR data constant $n_z = 50,197$ signals, each of length $n_z = 57$. The model mean and admit an intercept. The specification for the priors follows second Lib the support for elements in ϕ_w adjusted to match the BCEF spatial ent.

The n_y =197 locations with h_y =3 forest outcomes were and the secondstage model (3.6). To more clearly inte the spatial factors' ability to explain the variability in forest outcomes, v to avoid potential issues with spatial confounding Hanks 201set $\mathbf{v}(\mathbf{s})$ to zero. In practice, however, if ot main object lize predictive performance, then should likely be included in the model. this residual spatial vand In addition to t the second-stage model was informed by \blacksquare cap predictor variables defined in Section 2, the three La tercept, were included in $X_y(s)$. Importantly, these which. are available across the entire BCEF; hence, given the predicted va. of the spatial factors at unobserved locations, we can create ge forest outcome maps.

Posterior inference for all candidate models was based on three chains of 50,000 post-burn-in MCMC samples. Chains converged by 20,000 MCMC

iterations. Using the same computer configuration detailed in Section 4, the total runtime for the most demanding model, $q_w = 8$, was ~ 36 hours.

The eight candidate models, specified by q_w , were assessed based on their ability to inform the forest outcome predictions. This was done by fitting each of the first-stage models, then fitting their corresponding second cage models using data from 99 of the 197 available trons in \mathcal{T}_y . There exists outcomes were then predicted for the emaining 98 of the cations. The scoring rules and other summandation distributions for the 98 out-of-sample locations are presented in Σ , the 2.

Increasing the number of spatial factors improves the RMSPE for each forest outcome shown in Table 2 Exploratory analysis showed that the gains in predictive performance were $(e_w)^{-1}$ beyond $q_w = 4$ for AGB and $q_w = 5$ for TD and BA. Given that the $q_w = 5$ model generally yielded the "best" predictions, it is said to be exposition below.

Table 2: Cross-validation prediction summary for forest outcomes given increasing number of spatial factors q_w . Bold values identify lowest CRPS and RMSPE.

	q_w	CRSP	RMSPE	95% Coverage	95% Width
AGB	1	26.21	51.37	91.88	161.24
	2	26.36	52.02	92.39	162.14
	3	23.64	46.95	9	155.71
	4	23.53	46.93	9	15a
	5	24	47.54	90 15	157.
	6	24.47	47.8	94.	Q.
	7	24.75	47.84	95.43	11-11-
	8	24.76	48.02	96.45	182.12
	1	1017.7	1980.62	92.39	
TD	2	1006.02	1957.54	93.4	ə944.81
TD	3	1007.72	1954.87	93.4	6068.29
	4	$997. \ \ 2$	1955.2		6040.06
	5	$989.\overline{31}$	1930.76	14	6182.2
	6	998.3	1944.7		6223.73
	7	1005.2	196 8	510	6450.5
	8	1004.36		96.95	6503.17
ВА	1	5.53	29	91.88	36.34
	2	4	0.	94.42	36.85
DA			54	93.91	35.16
	4		52	93.4	36.21
	1	16	9.58	93.4	36.51
		5	9.59	96.45	38.62
	7		9.73	95.43	38.34
		5.27	9.72	94.42	37.93

Table 3: Elements of Λ_y median and 95% credible intervals for the $q_w = 5$ model. Bold entries indicate where the 95% credible interval excludes zero.

$\lambda_{AGB,1}^{(y)}$ -6.65 (-8.89, -4.23) $\lambda_{AGB,2}^{(y)}$ 27.20 (-14.11, 65.1) $\lambda_{AGB,3}^{(y)}$ -278.29 (-324) 2.28) $\lambda_{AGB,4}^{(y)}$ -46.15 (-16, 6, 75.91) $\lambda_{AGB,5}^{(y)}$ -308.81 (-524, 2, 90.45) $\lambda_{TD,1}^{(y)}$ -1.77 (-21.35, 17.60) $\lambda_{TD,2}^{(y)}$ -357.49 (-718.82, -7.86) $\lambda_{TD,3}^{(y)}$ 269.03 (-137.51, 667.0	1
$\lambda_{AGB,2}^{(y)}$ 27.20 (-14.11, 65.1) $\lambda_{AGB,3}^{(y)}$ -278.29 (-324) 2.28) $\lambda_{AGB,4}^{(y)}$ -46.15 (-16, 6, 75.91) $\lambda_{AGB,5}^{(y)}$ -308.81 (-524) 90.45) $\lambda_{TD,1}^{(y)}$ -1.77 (-21.35, 17.60) $\lambda_{TD,2}^{(y)}$ -357.49 (-718.82, -7.86) $\lambda_{TD,2}^{(y)}$ 269.03 (137.51, 667)	
$\lambda_{AGB,3}^{(y)}$ -278.29 (-324 2.28) $\lambda_{AGB,4}^{(y)}$ -46.15 (-16 6, 75.91) $\lambda_{AGB,5}^{(y)}$ -308.81 (-524 2.290.45) $\lambda_{TD,1}^{(y)}$ -1.77 (-21.35, 17.60) $\lambda_{TD,2}^{(y)}$ -357.49 (-718.82, -7.86) $\lambda_{TD,2}^{(y)}$ 269.03 (137.51, 667)	
$egin{array}{cccccccccccccccccccccccccccccccccccc$)
$egin{array}{cccccccccccccccccccccccccccccccccccc$	
$\lambda_{TD,1}^{(y)}$ -357.49 (-718.82, -7.86)	
$\lambda^{(y)}$ 260.03 (137.51.667)	
$\lambda_{\text{TD}}^{(y)}$ 269.03 (-137.51, 667.	
A _{TD,3} 209.00 (-101.01, 001.0	
$\lambda_{ extstyle extstyle au, 4}^{(y)}$ -1777.21 (-2696.67, -708.08	3)
$\lambda_{TD,5}^{(y)}$ 2457.52 (18, 4337.97))
$\lambda_{BA,1}^{(y)}$ -2.93 (- 975)	
$\lambda_{BA,2}^{(y)}$	
$\lambda_{BA,3}^{(y)}$ 9.79, -76.24)	
$\lambda_{BA,4}^{(y)}$ 72. 20.60, -23.00)	
$\lambda_{\rm B}^{\rm BA,4}$ 0.55 (-177.44, 20.51)	

Table 3 covide at increase for the second-stage model's spatial factor regress $\mathbf{v}_{\mathbf{v}}$ is in that is, the elements in $\mathbf{\Lambda}_{y}$. These results show that specifically specifically a substantial portion of the variability in the forest extremes. It is, however, difficult to interpret the different sense of what characteristic of $\mathbf{z}(\mathbf{s})$ the spatial factors are capturing. When considered with the estimates in Table 3, Figure 5 provides a biological interpretation of the spatial factors. Specifically, each panel in

Figure 5 represents a spatial factor. The 50 lines in each panel are observed LiDAR signals, with the lines corresponding to the 25 largest (lighter colored lines) and 25 smallest (darker lines) estimated spatial factor values.

There are some general biological relationships between to forest canopy structure and AGB, TD, and BA. A very low maximum canopy hei indicative of a young regenerating forest (e.g. wth after a fire and low BA would be characterized by low AGB, high T B, low of trees in a forest have a high canopy height, TD, and high BA (i.e., a few large-diameter mature trees dominated) e area). When the forest is characterized by trees of many diff i.e., tree crowns in several vertical strata), then we wight expect moderate/high AGB, moderate TD, and mod rate/high BA. Son se expected relationships are observed when comparing Tal 3 and ure 5. For example, the top left panel in Sigure 5 diffe regenerating forests and all other forest structures, that is, ter s show a spike of energy returned at or near ground evel line which show the majority of the energy of several meters. Hence, we have negative is returned at or e a $_{B,1}$ and $\lambda_{BA,1}^{(y)}$ in Table 3. The LiDAR signals shown egre ion coe el in Figure 5 differentiate between young and old singlein the to e., all trees were regenerated around the same time and there variation in canopy height); hence, we have negative $\lambda_{AGB,3}^{(y)}$ is little ver Table 3. The top middle and bottom left panels in Figure 5 generally separate the signal for mature 20+ and ~ 20 meter canopy heights (lighter lines), respectively, from the lower stature ~ 10 meter canopy height

forest (darker lines). Consistent with the biological expectation, the negative $\lambda_{TD,2}^{(y)}$ and $\lambda_{TD,4}^{(y)}$ suggest that forests associated with red LiDAR signals have higher tree density relative to the older taller forests.

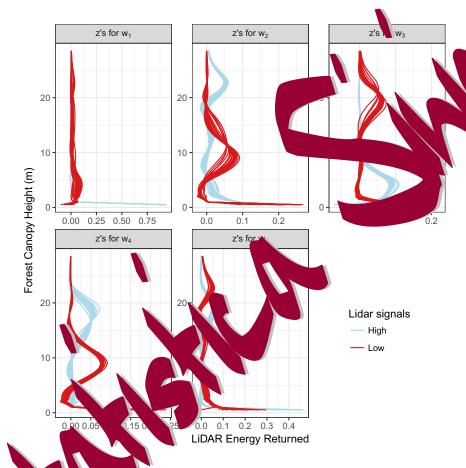


Figure 5. AR signals with the 25 largest (*High* in the legend) w in the legend) values of $\mathbf{w}(\mathbf{s})$ from the $q_w = 5$ model.

As det in Section 1, complete-coverage maps of the forest outcomes with associated uncertainty estimates are important data products that can be delivered by the proposed two-stage model. Following Section 3.4 and using the full data set depicted in Figure 1, we predicted the forest outcomes

on a 30×30 m grid over the BCEF. Figure 6 provides the median and 95% credible interval width maps for each outcome. Nonforested areas are omitted (white regions on the maps). The posterior predictive point estimates match well with the distribution of the forest outcomes a ross the BCEF and are clearly informed by the LiDAR factors, which are capturing kenforprediction unc est structure characteristics. Most importar maps, displayed in the right column of Figu 6, accurately of information for prediction units that are where LiDAR data are available, that is, we achieve more precise posor predictive distributions along and adjacent to locations data are available. Far from the LiDAR flight limes, prediction is only informed by the Landsat 8 tassel cal predictor variable. in this study explained very little variability in the forest

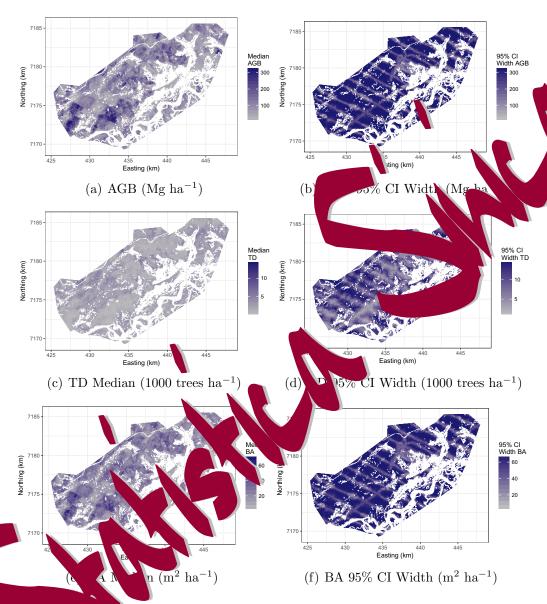


Figure 6. $q_w = 5$ posterior predictive distribution median and 95% CI width for B, TD, and BA forest variables over Bonanza Creek Experiment.

6. Concluding Remarks

We formulated an approach to model high-dimensional spatial data over a large set of locations, and developed an efficient implementation in C++. The SF-NNGP enables the analysis of multivariate spatially referenced data sets that, due to their magnitude, could not be rigo pusly explored by ore. It does so by combining the ability of SFN compress the signal high-dimensional structures into a few dimensions with the comparable scalability of NNGPs.

The algorithm was used to exploit the information from the mensional LiDAR signals to jointly model and generate LiDAR-bases aps of multiple forest variables. Importantly, the propo two-stage model provides a viable approach to producing spatially contil s from sparsely sampled LiDAR and forest measurement the model delivers spatially rthe explicit uncer ainty quantil ptures the irregular distribution nail. f interest. Such frameworks will become of information across the increasingly imp ng LiDAR systems, such as GEDI, come on-Sutu The approaches can also be extended to help guide LiDA equisition to minimize prediction uncertainty.

number of its results q_w to use in the model; there are different strategies to all the left decreases the consider out-of-sample evaluation metrics for different choices of q_w and select the one where the curves flatten out. This is a pragmatic solution, similar in spirit to cross-validation approaches commonly

used to tune hyper-parameters in richly parametrized models. Like any other cross-validation approach, this leads to additional computation, but parallel computing opens the possibility of conducting simultaneous MCMC runs for different values of q_w . As shown, both in the simulation experiment and in the BCEF data analysis, this heuristic provides sufficiently good roults. Other automated rank selection schemes are the literature of the as those proposed in Lopes and West (2004) of in Ren and Equipment (2004), however, these drastically increase the computationally costly problem.

In future research, we would like to explore an analysis spatiotemporal data. For this type of data, it is necessary to posit a strategy to select the neighbors in the spatio-temporal properties of following the discussion presented in Datta et al. (2016a).

Although our method results spin antial improvement in terms of scalability over existing approves, where efforts are required to scale multivariate spatial methods will be a large massive data sets. For instance, the ultimate goal for forest with an apping assisted by sampled LiDAR in interior Alash is a complete contrage map of the entire domain (e.g., 46 million ha), which contrage models capable of assimilating LiDAR signals in more contractions.

S tary Materials

The supplementary materials include (1) background information on NNGPs and spatial factor models, (2) the sampling algorithm for the SF-NNGP, and

(3) additional simulation results.

Acknowledgments

The research presented in this study was partially supported by NASA's Arctic-Boreal Vulnerability Experiment (ABoVE) and Carbon Monitoring System (CMS) programs. Additional support was a ovided by the Ucted States Forest Service Pacific Northwest Research Station. Fixley apported by National Science Foundation (NS) DMS-1513481, 137509, and EF-1241874, and Finley and Taylor-Rodriguez were supported by EF-1253225. Banerjee was supported by NSF DMS-1513654, NSF 562303, and NIH/NIEHS 1R01ES027027-01.

References

Aguilar, O. and West, M. (201) Baye in Oynan Faccor Models and Portfolio Allocation.

Journal of B siness & Econom. attis :338–357.

Andersen, H.-E., Strunk, J. and Green, (2011). Using airborne light detection and ranging as a sampling to the static crest biomass resources in the upper Tanana Valley of interior Alam Wester Journal of Applied Forestry, 26(4):157–164.

Anderson, and Statistics, Hoboken, NJ.

Asner, G., Hughen, C., Varga, T., Knapp, D., and Kennedy-Bowdoin, T. (2009). Environmental antrols over aboveground biomass throughout a tropical rain forest. *Ecosystems*, 12(2):261–278.

Babcock, C., Finley, A. O., Andersen, H.-E., Pattison, R., Cook, B. D., Morton, D. C., Alonzo,

- M., Nelson, R., Gregoire, T., Ene, L., Gobakken, T., and Næsset, E. (2017). Geostatistical estimation of forest biomass in interior alaska combining landsat-derived tree cover, sampled airborne lidar and field observations. *ArXiv e-prints*. https://arxiv.org/pdf/1705.03534.pdf.
- Babcock, C., Finley, A. O., Bradford, J. B., Kolka, R., Birdsey, R., and Ryan, M. G. 115).

 Lidar based prediction of forest biomass using his condels with spatially coefficients. Remote Sensing of Environment, 169:11 -127.
- Babcock, C., Matney, J., Finley, A., Weiskittel, A., and regression models for predicting individual tree structure variables using lide ta. *IEEE Journal of Selected Topics in Applied Earth Observations and Research* (1, SI):6–14.
- Baig, M. H. A., Zhang, L., Shuai, T., and Tong, Q. 14). Derivation of a tasselled cap transformation based on lands 8 at-satellite reflects by Sensing Letters, 5(5):423–431.
- Banerjee, S. (2017). High-dimensional by each geost Bayesian Anal., 12(2):583-614.
- Bechtold, W. A. Patro, R. (2005). The Enhanced Forest Inventory and Analysis Program (2005). Design and Estimation Procedures. US Department of Agriculture e. L. L. Lern Research Station Asheville, North Carolina.
- Blackford, L. S. mmel, J., Dongarra, J., Duff, I., Hammarling, S., Henry, G., Heroux, M., Lumsdaine, A., Petitet, A., Pozo, R., Remington, K., and Whaley, R. C. (2001). An updated set of basic linear algebra subprograms (blas). *ACM Transactions on Mathematical Software*, 28:135–151.

- Bolton, D. K., Coops, N. C., and Wulder, M. A. (2013). Measuring forest structure along productivity gradients in the Canadian boreal with small-footprint lidar. *Environmental monitoring* and assessment, 185(8):6617–6634.
- Bonanza Creek LTER (2016). Bonanza Creek Experimental Forest. http://www.lter.uaf.edu/research/study-sites-bcef. Accessed: 12-16-2017.
- Christensen, W. F. and Amemiya, Y. (2002). Latent variables of multivariate spaces of multivariate spaces of multivariate spaces. *Journal of the American Statistical Association*, 97 (7):302–317.
- Dagum, L. and Menon, R. (1998). Openmp: an interpretable standard api for shared-memory programming. Computational Science & Engineeric 1, 1998):46–55.
- Datta, A., Banerjee, S., Finley, A., Han et, A., and Schap, M. (2016a). Non-Separable Dynamic Near t-Neighbor Gaus are related as for Large Spatio-Temporal Data With an Application to Particle Mate Analysis. *Annals of Applied Statistics Statistics*, 44(2):629–659.
- Datta, A., Bane C., S., P., A. and Gelfand, A. E. (2016b). Hierarchical Nearest-Neighbor Gat. Viscos St. Large Geostatistical Datasets. *Journal of the American Statistical* 1(514):800–812.
- Datta, A., Baner S., Finley, A. O., and Gelfand, A. E. (2016c). On nearest-neighbor Gaussian for massive spatial data. Wiley Interdisciplinary Reviews: Computational Statistics, 8(5):162–171.
- Ene, L. T., Gobakken, T., Andersen, H.-E., Nsset, E., Cook, B. D., Morton, D. C., Babcock,

- C., and Nelson, R. (2018). Large-area hybrid estimation of aboveground biomass in interior alaska using airborne laser scanning data. *Remote Sensing of Environment*, 204(Supplement C):741 755.
- Finley, A. O., Banerjee, S., and Cook, B. D. (2014a). Bayesian hierarchical models for spatially misaligned data in R. *Methods in Ecology and Evolution*, 5(6) 514–523.
- Finley, A. O., Banerjee, S., Ek, A. R., and McRobert (1998). Bayesian multiprocess modeling for prediction of forest attributes. *Journal of Agricultu* (1998). Environmental Statistics, 13(1):60.
- Finley, A. O., Banerjee, S., Weiskittel, A. R., Babcock, C., and Cook, B. D. (2012) Dynamic spatial regression models for space-varying forest stand tables.
- Finley, A. O., Banerjee, S., Zhou, Y., Cook, B. D., and S., C. (2017). Joint hierarchical models for sparsely sampled high-divergent and Lie and forest variables. *Remote Sensing of Environment*, 190:149–161.
- G-LiHT (2016). Goddard's lid a pers trai and thermal (G-LiHT) imager. http://www.gliht.gsfc.na. 11-2017.
- GEDI (2014). what stem ynamics investigation lidar. http://science.nasa.gov/mis d: 8-11-2017.
- r, W. (2015). Cross-Covariance Functions for Multivariate Geostatistics.

 Statistical 2 vce, 30(2):147–163.
- G. (1996). Measuring the pricing error of the arbitrage pricing theory.
- Gneiting, T. and Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation.

 *Journal of the American Statistical Association, 102(477):359–378.

- Hanks, E. M., Schliep, E. M., Hooten, M. B., and Hoeting, J. A. (2015). Restricted spatial regression in practice: geostatistical models, confounding, and robustness under model misspecification. *Environmetrics*, 26(4):243–254.
- Heaton, M. J., Datta, A., Finley, A. O., Furrer, R., Guhaniyogi, R., Gerber, 7., Gramacy, R. B., Hammerling, D., Katzfuss, M., Lindgren, F., Nychka, D. W., Sun, F., and Zammit-Marion, A. (2017). Methods for analyzing large spatial data was an and comparison.
 prints. https://arxiv.org/abs/1710.05013.
- Hogan, J. W. and Tchernis, R. (2004). Bayesian factor analysta, with application to summarizing area-level material deprivation from census data.

 **nal of the American Statistical Association, 99(466):314–324.
- ICESat-2 (2015). Ice, cloud, and land elevation Wite-2. http://icesat.gsfc.nasa.gov/icesat2. Accessed: 8-11-117.
- Jakubowski, M. K., Guo, Q., and Kelly M. 2013). If the street meast ement accuracy. It is a few first meast ement accuracy. It is a few few from the few from the
- Junttila, V. and La variables under small field sample size. Remote Sensing of En. pp. 1-19 ment C):45 57.
- 199 Fraphical Models. Clarendon Press, Oxford, United Kingdom.
- Lopes, H. F. and St., M. (2004). Bayesian model assessment in factor analysis. Statistica Sinica,
- Murphy, K. (2012). Machine Learning: A probabilistic perspective. The MIT Press, Cambridge, MA.

- Næsset, E. (2011). Estimating above-ground biomass in young forests with airborne laser scanning. *International Journal of Remote Sensing*, 32(2):473–501.
- Nelson, R., Gobakken, T., Næsset, E., Gregoire, T., Ståhl, G., Holm, S., and Flewelling, J. (2012).

 Lidar sampling using an airborne profiler to estimate forest biomass i Hedmark County,

 Norway. Remote Sensing of Environment, 123:563–578.
- Nelson, R., Margolis, H., Montesano, P., Sun, G., Cook L., Andersen, H.-E., B., Pellat, F. P., Fickel, T., Kauffman, J., and Price, S. (2017). Lidar-like aboveground biomass in the continental us and mexicons as a satellite observations. Remote Sensing of Environment, 188(Supplement C):127 140
- Ren, Q. and Banerjee, S. (2013). Hierarchical Factor Models for Large angled Data:

 A Low-Rank Predictive Process Approach. Biotechnology (1913).
- Schmidt, A. M. and Gelfand, A. E. (2003). A bayesian of station approach for multivariate pollutant data. *Journal of Geophys. al.* search search spheres, 108(D24).
- Stein, M. L. (2011). Limitations of an analysis and mations for covariance matrices of spatial data. Spatial Statistics, 8:1
- White, J. C., Wulder, V. Cook, B. D., Pitt, D., and Well M. (2). A lest practices guide for generating forest inventory attributes from Victory land and ling data using an area-based approach. The Forestry Chronicle,
- Woodall, C. W. Julston, J. W., Domke, G. M., Walters, B. F., Wear, D. N., Smith, J. E., Clough, B. J., Cohen, W. B., Griffith, D. M., et al. (2015). The US forest carbon accounting framework: Stocks and stock change, 1990-2016.
- Yeniay, O. and Goktas, A. (2002). A comparison of partial least squares regression with other

prediction methods. Hacettepe Journal of Mathematics and Statistics, 31(99):99-101.

Zhang, H. (2007). Maximum-likelihood estimation for multivariate spatial linear coregionalization models. *Environmetrics*, 18(2):125–139.

Zhang, X. (2016). An optimized blas library based on gotoblas2. https://libraryi/OpenBLAS/. Accessed 2015-06-01.

Department of Mathematics & Statistics, Portland State Iniversity, Portland

E-mail: dantayrod@pdx.edu

Department of Forestry, Michigan State University, East Lansing, MI

E-mail: finleya@msu.edu

Department of Biostatistics, Johns Hopkins Univer Baltimore, MA

E-mail: abhidatta@jhu.edu

School of Environmental and Forest Scil com Iniversal Vashington, Seattle, WA

E-mail: babcoc7 Juw.edu

USDA Forest Service Pacific No. Search Station, Seattle, WA

E-mail: handersen@vs

Biospheric Scient Labory, No. A Goddard Space Flight Center, Greenbelt, MD

E-mail: vc d 2

ratory, NASA Goddard Space Flight Center, Greenbelt, MD

E-mail: douglas. ton@nasa.gov

Latistics, University of California Los Angeles, Los Angeles, CA

E-mail: sudipto@ucla.edu