# BART-SIMP: a novel framework for flexible spatial covariate modeling and prediction using Bayesian additive regression trees

Alex Ziyu Jiang[*], Jon Wakefield[†]

**Abstract.** Prediction is a classic challenge in spatial statistics and the inclusion of spatial covariates can greatly improve predictive performance when incorporated into a model with latent spatial effects. It is desirable to develop flexible regression models that allow for nonlinearities and interactions in the covariate structure. Machine learning models have been suggested in the spatial context, allowing for spatial dependence in the residuals, but fail to provide reliable uncertainty estimates. In this paper, we investigate a novel combination of a Gaussian process spatial model and a Bayesian Additive Regression Tree (BART) model. The computational burden of the approach is reduced by combining Markov chain Monte Carlo (MCMC) with the Integrated Nested Laplace Approximation (INLA) technique. We study the performance of the method via simulations and use the model to predict anthropometric responses, collected via household cluster samples in Kenya.

**Keywords:** Bayesian additive regression trees, Spatial covariate modeling, Spatial prediction, Integrated nest Laplace approximation, Survey sampling.

## 1 Introduction

There have been growing interests in using covariates for spatial data modeling (Uwiringiyimana et al., 2022; Macharia et al., 2019), and previous studies show that the inclusion of influential spatial covariates can lead to improved prediction accuracy (Lindström et al., 2014; Zeraatpisheh et al., 2022). It is straightforward to include covariates within the linear predictor of spatial random effects models, but the inclusion of flexible covariate models is much more difficult. Flexible modeling is desirable to detect nonlinearities and interactions which can lead to improved predictive performance.

Currently, Gaussian random field (GRF) models are commonly used as a modeling framework for capturing spatial correlations, with wide application, for example, in population health modeling (Diggle and Giorgi, 2019). A common approach to computation in spatial modeling is Integrated nested Laplace approximation (INLA) (Rue et al., 2009), particularly in a low-and-middle-income countries (LMIC) context (Utazi et al., 2018; Burstein et al., 2019). INLA is a powerful tool for carrying out Bayesian inference, but requires the predictors to have a linear form, which fails to capture nonlinearities and may lead to reduced prediction performance.

---

[*]Department of Statistics, University of Washington, US, jiang14@uw.edu

[†]Department of Statistics and Biostatistics, University of Washington, US, jonno@uw.edu

On the other hand, there are many machine learning methods that allow for flexible covariate modeling, including the stacked generalization method (Davies and Van Der Laan, 2016) and random forests (Ren et al., 2018) and these have been applied to spatial modeling problems (Osgood-Zimmerman et al., 2018; Shi et al., 2021). However, these approaches suffer from drawbacks. A number of machine learning approaches in spatial modeling ignore the spatial dependence (Georganos et al., 2021) while others do not correctly propagate uncertainty as desired (Daw and Wikle, 2023).

Bayesian Additive Regression Trees (BART) (Chipman et al., 2010) provide reliable Bayesian inference by specifying a prior distribution on the 'sum-of-trees' structure that flexibly models the covariates. Previous applications of BART on spatial data modeling problems include Müller et al. (2007) who proposed a spatial BART model based on a conditional autoregressive (CAR) model, and Krueger et al. (2020) suggested using a matrix exponential spatial specification (MESS) (LeSage and Kelley Pace, 2007) as an extension of CAR. However, compared to the GRF that we will use, these models specify a simple variance-covariance structures and are designed for area-level data. The model considered by Spanbauer and Sparapani (2021) assumes a more generic covariance structure under a non-spatial setting, but is limited to models with few random effects and will not scale well to large spatial datasets. As a result, we note that the large amount of random effects in continous spatial modeling and the strong dependence in the posterior which lead to computational difficulties and pose serious challenges for incorporating BART into spatial models.

In this paper, we propose **BART-SIMP**, which is shorthand for '**BART** for **S**patial **I**NLA **M**odeling and **P**rediction', as a GRF spatial Bayesian modeling framework with a flexible covariate regression model. Our model leverages the flexibility in BART to capture the nonlinearity and interactions across covariates, while also recognizing the complex spatial correlation structure in the residuals which is modeled by the GRF. To ease computation, we use the INLA-within-MCMC method (Gómez-Rubio and Rue, 2018) to design a Metropolis-within-Gibbs sampler that integrates out the random effects using INLA and then performs MCMC updates on the remaining parameters. This model is the first to simultaneously model the covariate effects through BART and the spatial effects through a GRF, with valid fully Bayesian model uncertainty estimates. Our spatial BART model is also the first to use INLA as an approximate approach to reduce the computational burden.

We organize the paper as follows: we will describe our motivating example, which is spatial prediction for child anthropometric data in Section 2. The model formulation is in Section 3. We outline the Metropolis-within-Gibbs sampler which we use for model implementation in Section 4. We apply our model to simulated data experiments in Section 5 and return to the Kenya data in Section 6. Section 7 concludes the paper with a discussion.

## 2   Motivating Dataset

Our study is motivated by child undernutrition data collected in the 2014 Kenya Demographic and Health Survey (DHS). Specifically, we are interested in wasting for children

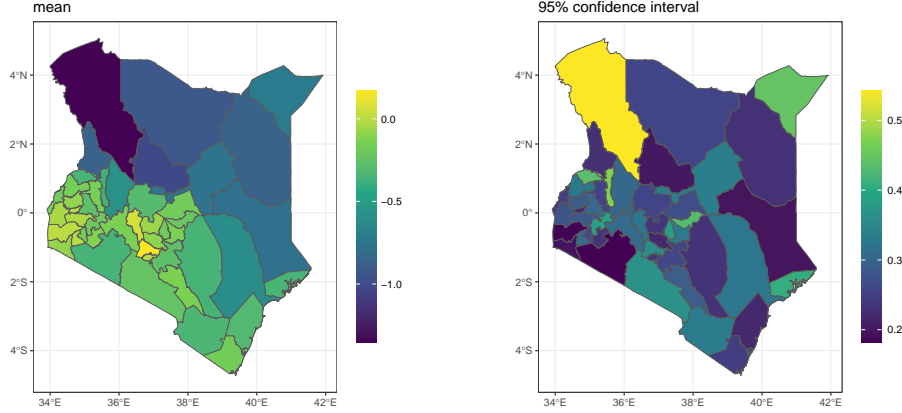| Variable | Years | Source |
|---|---|---|
| Population density | 2014 | WorldPop[1] |
| Night Time Light | 2013 | National Centers for Environmental Information[2] |
| Vegetation Index | 2013-2014 | NASA EOSDIS Land Processes DAAC[3] |
| Average Temperature | Mean value over period 1970-2000 | WorldClim[4] |
| Precipitation | Mean value over period 1970-2000 | WorldClim[5] |
| Access to Nearest City | 2000 | Joint Research Centre of the European Commission[6] |

Table 1: DHS Geospatial covariates used in our model.

under the age of five, which is measured using the weight-for-height Z-score (WHZ) metric. The Z-score can be interpreted as the number of times the weight for a child is higher or lower than the median value, compared to all children of the same height in the population (Mei and Grummer-Strawn, 2007). For example, wasting is defined by a WHZ score below $-2$ (Kassie and Workie, 2019).

The data is collected via stratified two-stage cluster sampling. In the first stage, 1584 enumeration areas (EAs) were sampled across the 92 strata. The 92 sampling strata are defined according to the 47 counties and the urban-rural status (Nairobi and Mombasa county have only urban areas). In the second stage, 40,300 households were sampled across the selected EAs, yielding 20,977 individual WHZ measurements. Additionally, we use geospatial covariates on the raster level, with a list of covariate descriptions shown in Table 1. In our example, for simplicity, we did not include the urban-rural status in our model, which may lead to bias due to oversampling of urban clusters but we believe that the inclusion of population density provides some protection.

The aim of the analysis is to provide predictions of WHZ at various administrative levels. Admin 0 denotes the national level, Admin 1 one below (for example, states in the United States) and Admin 2 one below that (counties in the United States). In Kenya, the Admin 1 level contains 47 counties, while Admin 2 contains 290 constituencies. This is a problem in small area estimation (SAE). If there are sufficient data in each area, a weighted (so-called direct) estimator can be used. But often there are insufficient area-based data and models must be introduced In an area-based model (Fay and Herriot, 1979) the direct estimator is modelled and random effects are introduced. In a unit level model point level data are modeled (Battese et al., 1988), and this is the path we follow. Rao and Molina (2015) provides a comprehensive overview of SAE. For these data, direct estimators tend to produce unreliable uncertainty estimates, as shown for WHZ Admin 1 predictions in Figure 1. There are higher levels in the northwest but it is hard to make definitive statements given the large uncertainty in this part of the map. Hence, we aim to reduce uncertainty using spatial smoothing and covariate modeling.

Figure 1: Maps of Admin 1 level WHZ mean estimates (left) and 90% confidence interval widths.



## 3   Model Description

To describe the data framework and introduce notation, we will describe a conventional traditional spatial model with linear predictors in Section 3.1. We will then introduce BART in Section 3.2 and propose our spatial version in Section 3.3. We describe the priors in Section 3.4.

### 3.1   Model formulation and preliminaries

We consider the spatial dataset $(\mathbf{Y}, \boldsymbol{x})$, where $\mathbf{Y}$ is the collection of observed outcomes and $\boldsymbol{x}$ is the covariate matrix. Let $n$ be the number of spatial locations observed (in our example this corresponds to clusters), noting that we allow multiple outcomes (these are children in our case) to be observed at the same location so that $\boldsymbol{y} = (\boldsymbol{y}_1, ..., \boldsymbol{y}_n)^\mathsf{T}$ is a jagged array of vectors, with $n_i$ being the number of observations in cluster $i$ and $\boldsymbol{y}_i = (y_{i1}, \ldots, y_{in_i})^\mathsf{T}$ being the vector of all observations at the same location. Specifically, we let $y_{ik} := y_k(\boldsymbol{s}_i)$ be the $k$-th response observed at location $\boldsymbol{s}_i$, $i = 1, \ldots, n$. We let $P$ be the number of covariates (including the intercept), let $\boldsymbol{x}_i = (x_{i1}, \ldots, x_{iP})$ be the vector of covariates observed at location $\boldsymbol{s}_i$ where $x_{ip}$ is the value of covariate $p$ observed at cluster $i$ for $p = 1, \ldots, P$ (for simplicity of notation, we assume that observations at the same spatial unit shares the same covariate values). We assume Gaussian measurement error with variance $\sigma_e^2$. Also, we let the latent random spatial component at location $\boldsymbol{s}_i$ be $z(\boldsymbol{s}_i)$, following a Gaussian random field (GRF) with Matérn covariance function

---

[1]http://www.worldpop.org.uk/data/get_data/.

[2]https://ngdc.noaa.gov/eog/dmsp/downloadV4composites.html.

[3]https://lpdaac.usgs.gov/dataset_discovery/modis/modis_products_table/mod13a3_v006.

[4]http://worldclim.org/version2.

[5]http://worldclim.org/version2.

[6]http://forobs.jrc.ec.europa.eu/products/gam/download.php.

over the region. The overall linear mixed effects model with a spatial random field can be written as:

$$y_{ik} \mid \boldsymbol{x}_i, \boldsymbol{\beta}, \boldsymbol{z}, \sigma_e^2 \overset{i.i.d}{\sim} N\left(\boldsymbol{x}_i^\top \boldsymbol{\beta} + z_i, \sigma_e^2\right), \tag{3.1}$$

$$\boldsymbol{z} \mid \boldsymbol{\psi} \sim GF(0, \boldsymbol{\Sigma}), \tag{3.2}$$

where the $n \times n$ matrix $\boldsymbol{\Sigma}$ is a Matérn covariance matrix so that

$$\Sigma_{ij} = \sigma_m^2 \, \mathrm{Corr}_M\left(z\left(\boldsymbol{s}_i\right), z\left(\boldsymbol{s}_j\right)\right) = \frac{2^{1-\nu}\sigma_m^2}{\Gamma(\nu)}\left(\kappa \left\|\boldsymbol{s}_i - \boldsymbol{s}_j\right\|\right)^\nu K_\nu\left(\kappa \left\|\boldsymbol{s}_i - \boldsymbol{s}_j\right\|\right), \tag{3.3}$$

with $\sigma_m^2$ and $\kappa$ being respectively the marginal variance and the scale parameter of the GRF. Note that the smoothness parameter $\nu$ is usually fixed *a priori*, and there is a one-to-one relationship between the scale parameter $\kappa$ and the range parameter (which is the more commonly interpreted parameter when it comes to spatial modeling) is given by $\rho = \frac{\sqrt{8\nu}}{\kappa}$. We let $\boldsymbol{\theta} = (\sigma_e^2, \boldsymbol{\psi})$ with $\boldsymbol{\psi} = (\sigma_m^2, \rho)$ the set of Matérn covariance parameters.

For latent Gaussian models (LGMs), a closed-form expression for $\pi(\boldsymbol{z}, \boldsymbol{\beta}, \boldsymbol{\theta} \mid \boldsymbol{y})$ is usually computationally expensive, largely due to the computation of the inverse and determinant of the dense $n \times n$ covariance matrix. Efficient MCMC sampling algorithms are difficult for spatial models because of the strong dependence between parameters (Knorr-Held and Rue, 2002). Rue et al. (2009) showed that the INLA method, an approximation Bayesian inference procedure based on the Laplacian approximation and numerical integration rules for LGMs is available for models where the latent Gaussian random field is distributed as a Gaussian Markov random field (GMRF). The INLA approach provides estimates for the marginal distributions of the latent Gaussian random field and the hyperparameters, along with a reliable estimate for the marginal distribution $\pi(\boldsymbol{y})$ (Hubin and Storvik, 2016).

For models with continuously indexed Gaussian fields, Lindgren et al. (2011) showed that a discrete GMRF can be used to approximate the continuously indexed data with GRF generated through a Matérn covariance function via a stochastic partial differential equation (SPDE) approach. This not only reduces the computation complexity from $\mathcal{O}(n^3)$ to $\mathcal{O}(n^{3/2})$, but also enables model inference to be carried out through INLA.

## 3.2 The BART model

The Bayesian Additive Regression Trees (BART) model (Chipman et al., 2010) can be viewed as the summation of a fixed number of highly flexible nonparametric regression functions, where each function $g(\cdot; T_l, \boldsymbol{\mu}_l): \mathbb{R}^p \to \mathbb{R}$ is a random function that maps the $p$-dimensional covariate space into the real space, with a total number of $l = 1, \ldots, m$ functions. The component functions are defined by two sets of parameters, the binary tree structure $\mathbf{T} = (T_1, \ldots, T_m)$ and the terminal node values $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_m)$. Each binary tree can be viewed as a set of decision rules that splits the covariate space into a finite number of regions: at any internal nodes for a given tree, a 'splitting variable' is chosen from the $p$ covariates and a threshold value is chosen that splits the current

region into two sub-regions, and a fitted value is assigned to the corresponding region at each end node, according to the terminal node values $\boldsymbol{\mu}$. By splitting according to different covariate variables and threshold values, BART is capable of modeling complex nonlinear relationships and interactions among covariates. Furthermore, a regularization prior that penalizes tree complexity in order to avoid model overfitting is assigned to the tree structure parameters so that the trees functions are 'weak learners' (Chipman et al., 2010).

## 3.3   Sampling Model and Latent Field

The sampling model for BART-SIMP is a combination of the linear mixed effects model described in Section 3.1 and the BART terms in Section 3.2. We substitute the linear covariate effects for the sum-of-tree model to obtain:

$$
y_{ik} \mid z_i, x_{ip}, \mathbf{T}, \boldsymbol{\mu}, \sigma_e^2 \overset{i.i.d}{\sim} N \left( \sum_{l=1}^{m} g\left( \boldsymbol{x}_{ip}; T_l, \boldsymbol{\mu}_l \right) + z_i, \sigma_e^2 \right)
$$

with the distribution of $\boldsymbol{z}$ and Matérn covariance as previously defined in Equations (3.2) and (3.3).

## 3.4   The Priors

We assume a priori independence between the tree parameters $(\mathbf{T}, \boldsymbol{\mu})$, the residual variance parameter $\sigma_e^2$ and the spatial hyperparameters $\boldsymbol{\psi}$:

$$
p\left( \mathbf{T}, \boldsymbol{\mu}, \sigma_e^2, \boldsymbol{\psi} \right) = \left[ \prod_{l=1}^{m} p\left( \boldsymbol{\mu}_l \mid T_l \right) p\left( T_l \right) \right] p\left( \sigma_e^2 \right) p\left( \boldsymbol{\psi} \right)
$$

For the tree parameters and the residual variance parameter, we follow the prior specifications in Chipman et al. (2010) and decompose the tree structure prior into three hierarchical parts: the probability of a node being non-terminal, the probability of choosing one of the covariates as the splitting variable for a non-terminal node and the probability of choosing a splitting value given the chosen splitting variable. For a tree node at depth $d$, the probability of it being a non-terminal node is $\alpha(1+d)^{-\beta}$, with $\alpha \in (0, 1), \beta \in [0, \infty)$. In our examples in the paper, we defer to the default settings in Chipman et al. (2010) and choose $\alpha = 0.95$ and $\beta = 2$. Further, we assume that the splitting variable is uniformly chosen among all covariates for each internal node, and the splitting value is uniformly chosen from all distinct values from the selected splitting variable. We assume an independent and identically distributed normal prior for the terminal node values $\boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_m$ given their corresponding tree structure $T_l$. The mean and variance for the normal prior are chosen such that each tree will only function as a 'weak learner' that contributes a small part in the model; we refer to Chipman et al. (2010) for more details.

For the residual variance parameter $\sigma_e^2$, we take the default choice in Chipman et al. (2010) and specify a scaled inverse-Gamma distribution $\sigma_e^2 \sim \nu\lambda/\chi_\nu^2$, where $\nu$ is

the degrees-of-freedom. The values of the hyperparameters $(\nu, \lambda)$ are chosen through a exploratory data analysis procedure such that $\hat{\sigma}_r^2$, an empirical guess of $\sigma_e^2$ through a tentative working model (for example, an SPDE model with all covariates as linear predictors), matches the $q$-th quantile of the scaled inverse Gamma prior. Following Chipman et al. (2010), we let $(q, \nu) = (0.9, 3)$. Hence, the prior depends on the data in a weak fashion. Finally, we place a penalized complexity (PC) prior on the Matérn hyperparameters $\psi$, following Fuglstad et al. (2019). The joint PC prior for $\psi$ is

$$p(\psi) = \frac{d}{2} \tilde{\lambda}_1 \tilde{\lambda}_2 \rho^{-d/2-1} \exp\left(-\tilde{\lambda}_1 \rho^{-d/2} - \tilde{\lambda}_2 \sigma_m\right), \sigma_m > 0, \rho > 0,$$

where $d = 2$, $\tilde{\lambda}_1 = -\log(\alpha_1)\rho_0^{d/2}$ and $\tilde{\lambda}_2 = -\log(\alpha_2)/\sigma_0$. The hyperparameters $(\alpha_1, \alpha_2, \rho_0, \sigma_0)$ guarantee that $P(\rho < \rho_0) = \alpha_1$ and $P(\sigma_m > \sigma_0) = \alpha_2$. We let $\alpha_1 = \alpha_2 = 0.5$ and set $(\rho_0, \sigma_m)$ to the crude estimates obtained from fitting a model with only intercept and spatial random field.

## 4 Algorithm

Under our model framework, a standard MCMC algorithm requires the sampling and storage of $z$, which is computationally undesirable. Hence, our strategy is to integrate out $z$ and operate on the corresponding marginal likelihood. However, the computation of the marginal likelihood involves computing the determinant of large and dense variance-covariance matrices, which becomes prohibitive as the number of clusters becomes large. As a result, we consider the 'INLA-within-MCMC' technique proposed by Gómez-Rubio and Rue (2018). Specifically, the algorithm calls INLA to approximate the aforementioned marginal likelihood, which is used to compute the Metropolis-Hastings acceptance ratio in the MCMC routine for the remaining model parameters.

We let $T_{-j}, \mu_{-j}$ respectively be the set of trees and terminal node values, except for $T_j$ and $\mu_j$. We give an outline of the MCMC routine in Algorithm 1.

---

**Algorithm 1:** An overview of the BART-SIMP algorithm.

**Input:** Initialized values for $\mathbf{T}, \mu, \sigma_e^2, \psi$ for a number of MCMC iterations $B$.
**for** $b \leftarrow 1$ **to** $B$ **do**
    **for** $l \leftarrow 1$ **to** $m$ **do**
        **Update** $T_l \mid y, T_{-j}, \mu_{-j}, \sigma_e^2, \psi$.
        **Update** $\mu_l \mid y, T_j, T_{-j}, \mu_{-j}, \sigma_e^2, \psi$.
        **Update** $\sigma_e^2 \mid y, \mathbf{T}, \mu, \psi$.
        **Update** $\psi \mid y, \mathbf{T}, \mu, \sigma_e^2$.
    **end**
**end**

---

We now describe each of the updates.

## 4.1   Update $T_l \mid y, T_{-j}, \mu_{-j}, \sigma_e^2, \psi$

To update the sum-of-trees model, we follow the Bayesian backfitting MCMC algorithm, outlined in Section 3.1 of Chipman et al. (1998). Specifically, we do a sequential update of the set of tree structure parameters and terminal node values $\{(T_1, \mu_1), \ldots, (T_m, \mu_m)\}$, updating each tree, one at a time. To update each of the tree structure parameters $T_j$, we furthermore integrate out $\mu_j$ and do a Metropolis-within-Gibbs update of $T_j$, while keeping the other parameters, $(T_{-j}, \mu_{-j}, \sigma_e^2, \psi)$, fixed. The acceptance ratio for $T_j$, is

$$\alpha_{T_j} = \min\left\{1, \frac{\pi\left(y \mid T_j^*, T_{-j}, \mu_{-j}, \sigma_e^2, \psi\right) \pi\left(T_j^*\right) q\left(T_j \mid T_j^*\right)}{\pi\left(y \mid T_j, T_{-j}, \mu_{-j}, \sigma_e^2, \psi\right) \pi\left(T_j\right) q\left(T_j^* \mid T_j\right)}\right\}.$$

Here $T_j^*$ is the proposed structure for tree $j$ according to Chipman et al. (2010), through which we can compute the transition kernel ratio $\frac{q\left(T_j \mid T_j^*\right)}{q\left(T_j^* \mid T_j\right)}$. To compute the likelihood ratio, we employ the backfitting algorithm as follows. For $\pi\left(y \mid T_j, T_{-j}, \mu_{-j}, \sigma_e^2, \psi\right)$, we see that given the tree parameters of all trees except for $T_j$, we can compute the residual values with respect to the remaining $m - 1$ trees and denote the vectors of residuals as $r_1^{(j)}, \ldots, r_n^{(j)}$, where

$$r_{ik}^{(j)} = y_{ik} - \sum_{l \neq j} g\left(x_i; T_l, \mu_l\right), l = 1, \ldots, m, \quad i = 1, \ldots, n, \quad k = 1, \ldots, n_i.$$

We see that the original model is now equivalent to a single-treed model (along with the spatial field and the Gaussian noise) with response $(r_1^{(j)}, \ldots, r_n^{(j)})$, where the superscript $(j)$ is with respect to the particular tree being 'left out'. Since the tree structure for $T_j$ is conditioned upon, the tree part is equivalent to a linear model $\mathbf{C}^j \mu_j$ where $\mu_j = (\mu_{j1}, \ldots, \mu_{jb_j})$ is the set of terminal node parameters in $\mu_j$, and $\mathbf{C}^j$ is the $n \times b_j$ covariate matrix such that

$$C_{it}^j = 1 \text{ if } x_{[i]} \text{ belongs to terminal node } t \text{ in tree } T_j, \text{ otherwise } C_{it}^j = 0,$$

where $C_{it}^j$ is the element in the $i$-th row and $t$-th column of $\mathbf{C}^j$. Thus, given $T_j, T_{-j}, \mu_{-j}$, the whole model is equivalent to the following model

$$\begin{aligned} r^{(j)} \mid x, \sigma_e^2 &\overset{i.i.d}{\sim} N\left(\mathbf{C}^j \mu_j + z, \sigma_e^2\right) \\ z \mid \psi &\sim GF(\mathbf{0}, \mathbf{\Sigma}) \end{aligned} \tag{4.1}$$

where $\mathbf{\Sigma}$ follows the expressions in (3.3). This corresponds to a linear mixed effect model with Gaussian latent spatial effects. Thus we can approximate the Gaussian field with a GMRF using SPDE, with the spatial parameters fixed at $\psi$, the residual variance fixed at $\sigma_e^2$ and approximate the marginal likelihood of the model $\pi(r^{(j)} \mid \sigma_e^2, \psi)$, which is exactly $\pi(y \mid T_j, T_{-j}, \mu_{-j}, \sigma_e^2, \psi)$. We update $\mu_j$ given $(T_j, T_{-j}, \mu_{-j}, \sigma_e^2, \psi, y)$ using a similar Metropolis-within-Gibbs approach. Let $\mu_j^*$ be the proposed values during an update for $\mu_j$, the acceptance ratio is defined as follows:

$$\alpha_{\mu_j} = \min\left\{1, \frac{\pi\left(y \mid \mathbf{T}, \mu_{-j}, \mu_j^*, \psi, \sigma_e^2\right) \pi\left(\mu_j^*\right) q\left(\mu_j \mid \mu_j^*\right)}{\pi\left(y \mid \mathbf{T}, \mu_{-j}, \mu_j, \psi, \sigma_e^2\right) \pi\left(\mu_j\right) q\left(\mu_j^* \mid \mu_j\right)}\right\}.$$

Similarly, we use the backfitting algorithm to compute the residuals $\boldsymbol{r}$:

$$r_{ik} = y_{ik} - \sum_{l=1}^{m} g\left(\boldsymbol{x}_i; \mathbf{T}_l, \boldsymbol{\mu}_l\right), l = 1, \ldots, m, \quad i = 1, \ldots, n, \quad k = 1, \ldots, n_i.$$

This is equivalent to a model with the following sampling model:

$$\boldsymbol{r} \mid \boldsymbol{z}, \sigma_e^2 \overset{i.i.d}{\sim} N\left(\boldsymbol{z}, \sigma_e^2\right)$$
$$\boldsymbol{z} \mid \boldsymbol{\psi} \sim GF(\mathbf{0}, \boldsymbol{\Sigma}).$$

We can then approximate $\pi(\boldsymbol{r} \mid \sigma_e^2, \boldsymbol{\psi})$ using SPDE, which is equivalent to the conditional $\pi\left(\boldsymbol{y} \mid T_j, \mathbf{T}_{-j}, \boldsymbol{\mu}_{-j}, \boldsymbol{\mu}_j, \boldsymbol{\psi}, \sigma_e^2\right)$. We can also carry out the same procedure for $\pi\left(\boldsymbol{y} \mid T_j, \mathbf{T}_{-j}, \boldsymbol{\mu}_{-j}, \boldsymbol{\mu}_j^*, \boldsymbol{\psi}, \sigma_e^2\right)$.

Similar to the tree structure parameters $(T_1, \ldots, T_m)$, we use the Metropolis-within-Gibbs method to update the residual variance parameter $\sigma_e^2$ and use the INLA-within-MCMC technique to approximate the Metropolis-Hastings acceptance ratio, $\alpha_{\sigma_e^2}$, defined as follows:

$$\alpha_{\sigma_e^2} = \min\left\{1, \frac{\pi\left(\boldsymbol{y} \mid \mathbf{T}, \boldsymbol{\mu}, \psi, \sigma_e^{2*}\right) \pi\left(\sigma_e^{2*}\right) q\left(\sigma_e^2 \mid \sigma_e^{2*}\right)}{\pi\left(\boldsymbol{y} \mid \mathbf{T}, \boldsymbol{\mu}, \psi, \sigma_e^2\right) \pi\left(\sigma_e^2\right) q\left(\sigma_e^{2*} \mid \sigma_e^2\right)}\right\}.$$

To ensure that the variance parameter is positive, we use a Gaussian proposal for $\log \sigma_e^2$. We again use the backfitting technique to compute the overall residual $\boldsymbol{r}$, defined as all $m$ tree terms subtracted from the response value:

$$r_{ik} = y_{ik} - \sum_{l=1}^{m} g\left(\boldsymbol{x}_i; \mathbf{T}_l, \boldsymbol{\mu}_l\right), \quad i = 1, \ldots, n.$$

This is equivalent to a model with the following sampling model:

$$\boldsymbol{r} \mid \boldsymbol{z}, \sigma_e^2 \overset{i.i.d}{\sim} N\left(\boldsymbol{z}, \sigma_e^2\right)$$
$$\boldsymbol{z} \mid \boldsymbol{\psi} \sim GF(\mathbf{0}, \boldsymbol{\Sigma}).$$

By fitting a linear mixed effect model with residual variance fixed at $\sigma_e^2$ and spatial hyperparameters fixed at $\boldsymbol{\psi}$, we can compute the approximated marginal likelihood of the model $\pi(\boldsymbol{r} \mid \sigma_e^2, \boldsymbol{\psi})$, which is equivalent to $\pi\left(\boldsymbol{y} \mid \mathbf{T}, \boldsymbol{\mu}, \boldsymbol{\psi}, \sigma_e^2\right)$. We can update the spatial hyperparameters in a similar way.

## 4.2   R and C++ Integration

We implemented the backbone of the MCMC algorithm of BART-SIMP based on the C++ code in the `BART` package (Spanbauer and Sparapani, 2021). In order to use the functions in the `R-INLA` package (Martins et al., 2013) to carry out the INLA computation during the Metropolis-Hastings adjustment step, we used Rcpp (Eddelbuettel and François, 2011; Eddelbuettel, 2013; Eddelbuettel and Balamuta, 2018) as an interface between C++ and R. To provide an open-source computing software for BART-SIMP, we developed the R package `BARTSIMP`, with source code available at https://github.com/AlexJiang1125/BARTSIMP.

# 5   Simulation Experiments

## 5.1   Simulation setting

In this section, we study several scenarios in which the spatial and covariate signals have different strengths. We consider a $50 \times 50$ grid surface over a study region $[0,1] \times [0,1]$ and denote the set of grid cells as $G$. For the covariates associated with each grid, we independently generate five covariates from a uniform distribution on $[0,1]$. The actual deterministic field evaluated at cell $g$, denoted as $f_g(\boldsymbol{x})$, is defined as follows:

$$f_g(\boldsymbol{x}) = (1-\omega)\boldsymbol{z}_g^* + \omega f_{0g}(\boldsymbol{x}), \quad \forall g \in G,$$

where the 'raw' spatial field $\boldsymbol{z}^*$ is generated from a GRF with Matérn parameters $\kappa = 2.5$, $\sigma_m^2 = 0.5$. Likewise, we let $f_0(\boldsymbol{x})$ be the 'raw' covariate surface generated based on the Friedman function (Friedman, 2001) as follows:

$$f_{0g}(\boldsymbol{x}) = \sin\left(\pi x_{g1} x_{g2}\right) + 2\left(x_{g3} - 0.5\right)^2 + 2x_{g4} + x_{g5}, \quad \forall g \in G.$$

where $x_{gp}, g \in G, p = 1, ..., 5$ represents the $p$-th covariate for grid cell $g$. Finally, the scalar $\omega$ is fixed at different values and can be interpreted as the proportion of 'covariate signal' among the overall signal. Here we consider four different scenarios for $\omega$: 1 (covariate signal only), 0.8 (strong covariate signal), 0.5 (medium covariate signal) and 0.2 (weak covariate signal), which we denote as scenarios 1, 2, 3 and 4. We then randomly select 250 cells from the grid and randomly sample one location uniformly within each cell, ending up with 250 spatial points over the study region, mimicking the spatial locations for the clusters. The number of observations for each spatial location is sampled from a uniform distribution over $\{5, 6, 7, 8, 9, 10\}$. For each observation at a given spatial location, its value is defined as the sum of the deterministic field value from the grid cell it belongs to, and the Gaussian random noise with $\sigma_e^2 = 1$. We simulate 10 datasets for each scenario. For the BART-SIMP and BART model, we used ensembles of five trees, and collected 2,000 posterior samples after a burn-in period of 2,000.

## 5.2   Performance Criteria

We compare our model against the following models: a standard BART model without spatial correlations (**BART**), a GMRF spatial model with Matérn covariance function fitted by the SPDE approach, including all covariates in a linear model (**SPDE**) and a spatial model similar to SPDE, with intercept only (**SPDE0**). We focus on accuracy of interval estimates since this is clearly key for prediction and machine learning methods are often poor with respect to this aspect.

We consider the following model performance measures over the gridded surface:

- **Average coverage rate (ACR)**: Let $(L_{g,\alpha}^j, U_{g,\alpha}^j)$ denote the $100 \times (1-\alpha)\%$ prediction interval for $f(g)$ given by method $j$, the ACR is defined as the proportion of correct coverages of $f(g)$:

$$\mathrm{ACR}^j = \frac{1}{|G|} \sum_{g \in G} \mathrm{I}\left(f(g) \in L_{g,\alpha}^j, U_{g,\alpha}^j\right).$$

In our analysis, we let $\alpha = 0.10$.

- **Average Interval score (AIS, Gneiting and Raftery, 2007)**: the AIS is an integrated metric for prediction interval accuracy defined as follows:

$$\text{AIS}^j = \frac{1}{|G|} \sum_{g \in G} \left[ (U_{g,\alpha}^j - L_{g,\alpha}^j) + \frac{2}{\alpha} \left( L_{g,\alpha}^j - f(g) \right) \cdot \mathbf{1}(f(g) < L_{g,\alpha}^j) + \right.$$
$$\left. \frac{2}{\alpha} \left( f(g) - U_{g,\alpha}^j \right) \cdot \mathbf{1}(f(g) > U_{g,\alpha}^j) \right].$$

Note that AIS can be broken down into three parts, where the first part penalizes wider intervals and the second and third parts penalize low coverage rates.

## 5.3   Results

Tables 2 and 3 show comparisons of ACR and AIS across all four methods under the four scenarios. With respect to coverage, Table 2 shows that BART-SIMP has the best performance with only slight undercoverage. BART has less than half the nominal coverage in all 4 scenarios, because the spatial dependence in the residuals is ignored, the intervals are overly optimistic. SPDE always overestimates and has coverage rates very close to 100% (across all scenarios and in all simulated datasets). SPDE0 performs better than BART but still has undercoverage which improves slightly as the spatial signal increases. As a combination of both coverage and width metrics, we see that BART-SIMP performs best in scenarios 1 and 2 (with greater covariate signals), while SPDE0 achieves the best scores in the more spatial scenarios.

Table 2: Average coverage rates (ACR) for four different methods (nominal coverage is 90%) under four scenarios (covariate only, strong, medium and weak covariate signals). The averaged values over 10 replications are given in the cells, with standard deviations in brackets.

|  | Scenario 1 ($\omega = 1$) | Scenario 2 ($\omega = 0.8$) | Scenario 3 ($\omega = 0.5$) | Scenario 4 ($\omega = 0.2$) |
|---|---|---|---|---|
| BART-SIMP | **82**% (0.016) | 79% (0.023) | 59% (0.005) | 78% (0.007) |
| BART | 32% (0.069) | 33% (0.064) | 37% (0.096) | 35% (0.061) |
| SPDE | 98% (0.000) | **99**% (0.000) | **100**% (0.000) | **100**% (0.000) |
| SPDE0 | 51% (0.020) | 54% (0.014) | 57% (0.027) | 63% (0.037) |

Table 3: Average interval score (AIS) rates for four different methods under four scenarios (covariate only, strong, medium and weak covariate signals). The averaged values over 10 replications are given in the cells, with standard deviations in brackets. (Low scores are preferred.)

|            | scenario 1 $(\omega = 1)$ | scenario 2 $(\omega = 0.8)$ | scenario 3 $(\omega = 0.5)$ | scenario 4 $(\omega = 0.2)$ |
|------------|---------------------------|------------------------------|------------------------------|------------------------------|
| BART-SIMP  | **0.286** (0.014)         | **0.290** (0.019)            | 0.860 (0.011)                | 0.433 (0.007)                |
| BART       | 0.503 (0.108)             | 0.380 (0.077)                | 0.224 (0.050)                | 0.224 (0.052)                |
| SPDE       | 0.372 (0.012)             | 0.356 (0.005)                | 0.348 (0.006)                | 0.344 (0.006)                |
| SPDE0      | 0.462 (0.069)             | 0.314 (0.038)                | **0.152** (0.008)            | **0.053** (0.003)            |

# 6  Application

In this section, we investigate the prediction performance of BART-SIMP, compared with other commonly used methods, evaluated on the WHZ measurements from the 2014 Kenya DHS. To compare the prediction performances across different models, we apply cross-validation and split the 1584 clusters in 2014 DHS into training and test datasets of cluster sizes 1267 and 317 respectively. To preserve the stratification structure in both sets, we consider a stratified sampling of the clusters such that the strata proportions in the training set roughly matches the test set. Finally, we repeated the procedure 10 times to reduce sampling variation caused by using a single data split. We choose ACR and AIS as performance criteria of all four methods, similar to Section 5.2 and used the same settings for the BART and BART-SIMP models. The algorithm took three days to run for BART-SIMP on Ubuntu 18.04, with 50GB of memory.

Table 4 shows ACR (nominal coverage is 90%) and AIS for all four models. The results are reported as the average over 10 test datasets, with standard deviations in brackets. We see BART-SIMP dominates the other three methods in terms of coverage being the closest to the nominal coverage, while the SPDE model has the lowest AIS.

|   | model     | ACR             | AIS               |
|---|-----------|-----------------|-------------------|
| 1 | BART-SIMP | **86**% (0.029) | 0.101 (0.008)     |
| 2 | BART      | 77% (0.035)     | 0.106 (0.007)     |
| 3 | SPDE      | 81% (0.019)     | **0.095** (0.007) |
| 4 | SPDE0     | 81% (0.019)     | 0.096 (0.007)     |

Table 4: Prediction performance measures over all four competing methods. The nominal coverage is 90% and small values of AIS are preferred. The averages over 10 test datasets are shown with standard deviations shown in brackets.

## 6.1  Surface Prediction

In this section, we conduct spatial predictions on a grid surface over the study region and generate aggregated estimates at the Admin 1 and Admin 2 levels. At location $\boldsymbol{s}$, we let $\mathrm{WHZ}(\boldsymbol{s})$ be the spatial surface of the height-for-weight Z-scores and $d_5(\boldsymbol{s})$ be the under-five population density. The areal level WHZ is a weighted average over the under-five population density $d_5(\boldsymbol{s})$, as the Z-scores were evaluated for children under age five. The defined $d_5(\boldsymbol{s})$ values are obtained from WorldPop. The WHZ for an administrative region is defined as

$$\mathrm{WHZ}_{R_i} = \frac{\int_{R_i} \mathrm{WHZ}(\boldsymbol{s}) d_5(\boldsymbol{s})}{\int_{R_i} d_5(\boldsymbol{s})}, i = 1, 2, \ldots, |\mathcal{R}|. \tag{6.1}$$

where $\mathcal{R} = \{R_1, R_2, \ldots\}$ is the set of administrative regions.

We approximate the integrals in (6.1) by a weighted sum over observations on grid cells over the regions. Let $\mathrm{WHZ}(g)$ be the height-for-weight Z-score and $d_5(g)$ be the under-five population density evaluated at grid cell $g$, the regional WHZ approximated on the grid level is

$$\mathrm{WHZ}_{R_i} \approx \frac{\sum_{g \in R_i} \mathrm{WHZ}(g) d_5(g)}{\sum_{g \in R_i} d_5(g)}. \tag{6.2}$$

Using formula (6.2), we can calculate the posterior mean and 95% credible interval quantiles for regional WHZ at all Admin 1 and Admin 2 areas, based on the four methods mentioned. As a comparison, we also consider the direct weighted areal-level estimates and 95% confidence intervals. Figure 2 shows the posterior median and 95% credible/confidence intervals, derived from all five methods, and Figure 3 shows the predicted areal-level WHZ for all 48 Admin 1 areas in Kenya. The predicted posterior mean for WHZ given by BART-SIMP ranges from -1.10 to 0.05., showing that there is large within country variation in WHZ in Kenya. Among the 48 Admin 1 regions, Kiambu and Nairobi have the highest WHZ predictions – these are the most populated counties in Kenya. Low WHZ scores occured in other areas and these could be targeted for interventions.

We see from Figure 2 and 3 that while all five methods give quantitatively similar patterns of the areal estimate across the country region, BART-SIMP and the spatial methods provide similar point estimates, while BART-SIMP gives relatively wider credible interval lengths. The latter is consistent with the simulations, and suggest that these are more appropriate. Finally, BART does not yield reliable point estimates, and gives interval estimates that are too narrow.

We also provided areal-level predictions on Admin 2 levels, with results shown in the supplemental materials. Examination of these results show that the direct estimates for Admin 2 regions have more unreliable point estimates and much wider confidence intervals, due to insufficient samples observed in each region.
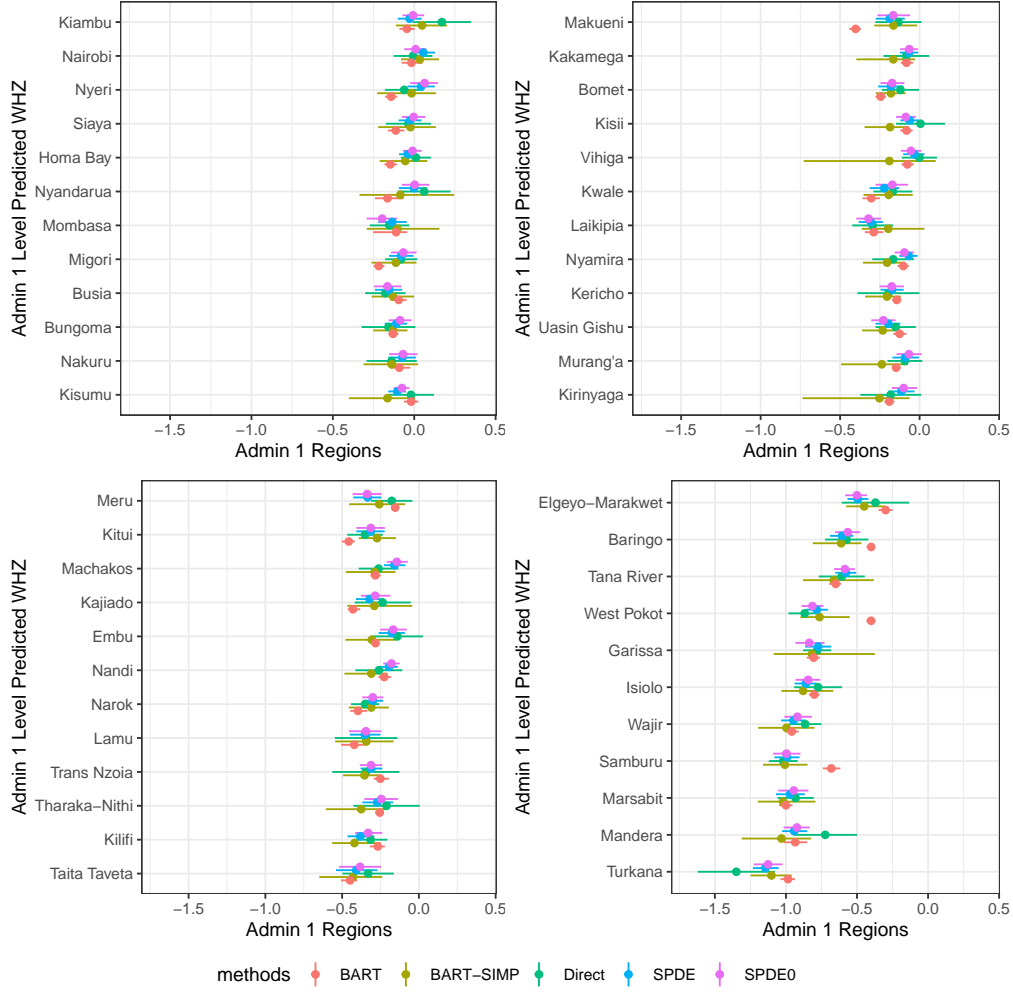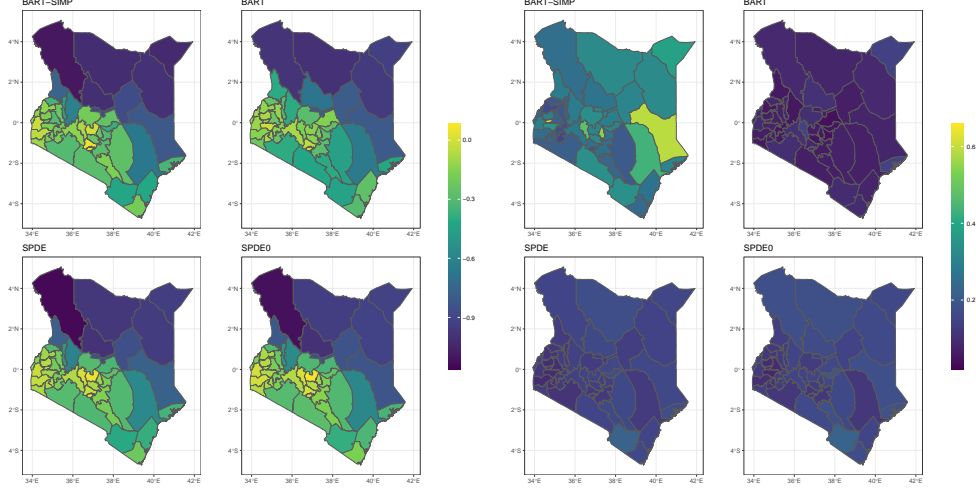
Figure 2: The Admin 1 level posterior median and 95% credible/confidence intervals of WHZ for BART, BART-SIMP, SPDE0, SPDE and direct estimates. The Admin 1 areas are arranged according to the predicted posterior mean given by BART-SIMP and grouped into four graphs.

## 6.2   Partial dependence

In our dataset, it is also of particular interest to study the marginal influence of a certain variable. Various measures on partial dependence have been proposed in the machine learning literature (Breiman, 2001; Goldstein et al., 2015). For models based on BART, the partial dependence function (Friedman, 2001) is a commonly used measure. Let $f(\cdot) : \mathbb{R}^p \to \mathbb{R}$ be a multivariate function defined on $p$ variables $(x_1, \ldots, x_p) = \boldsymbol{x}$, let $\boldsymbol{x}_t$

Figure 3: Maps of Admin 1 level WHZ posterior median (left) and 90% credible interval lengths (right) for BART-SIMP, BART, SPDE, SPDE0.



be the set of variables we want to study out of the $p$ variables, and let $\boldsymbol{x}_c = \boldsymbol{x}/x_t$ be its compliment, suppose we have $n$ observations of such multivariate variable, the partial dependence function for $f(\cdot)$ with respect to $x_t$ can be defined as

$$f^{pd}(x_t) = \frac{1}{n} \sum_{i=1}^{n} f(x_t, \boldsymbol{x}_{ic}),$$

where $\boldsymbol{x}_{i.}, i = 1, \ldots, n$ represents the $i$-th observation. For BART-SIMP and BART, we can calculate the partial dependence function of the sum-of-trees function for all six variables in our dataset. Figure 4 shows the partial dependence function and its 95% credible interval based on 1000 MCMC samples for the BART-SIMP and BART methods. As a comparison, we also included the locally-weighted smoother function estimated from the scatterplot of centered WHZ versus each variable. We observe that for all six variables, BART-SIMP and BART have similar partial dependence patterns that roughly resemble the pattern of the raw data. Substantively, holding all other variables constant, WHZ increases with increasing population density, vegetation index and precipitation, and decreases with increasing average temperature. The assocation with access is nonlinear but there is great uncertainty at low level of access. Note that BART has much narrower credible intervals, which support our previous finding that BART tends to underestimate model uncertainty.

# 7   Discussion

In this work, we have proposed BART-SIMP as a novel framework for flexible covariate modeling and prediction for spatial datasets. We incorporated the nonparametric nature
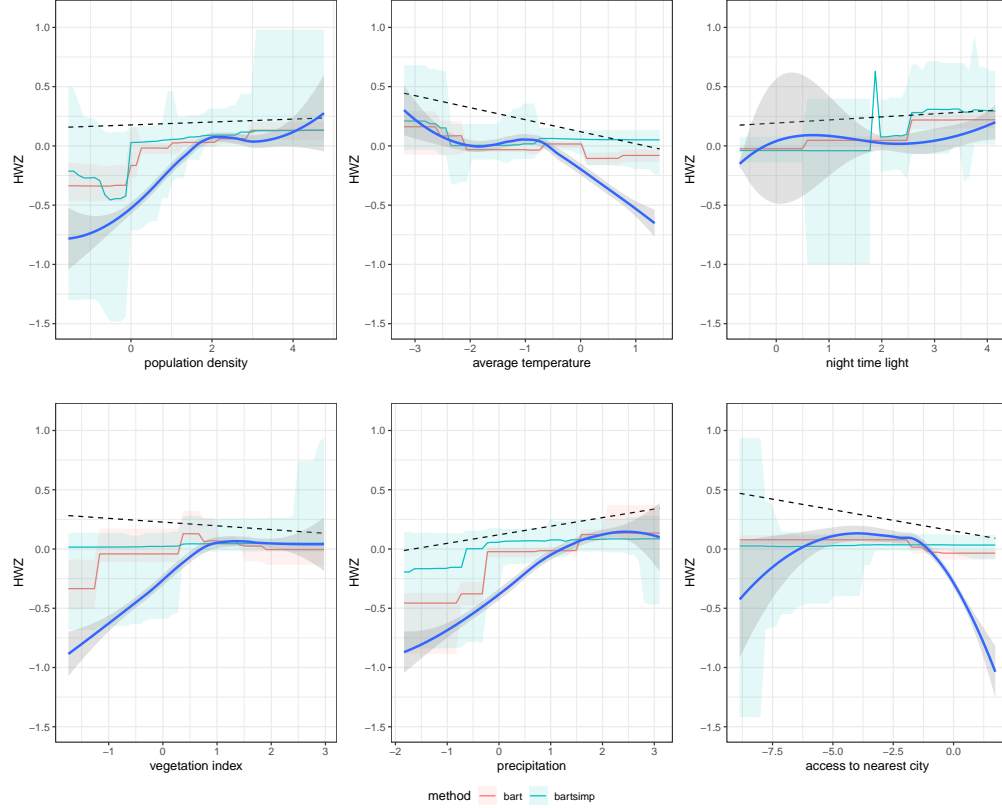
Figure 4: Partial dependence function for all six covariates (populational density, average temperature, night time light, vegetation index, precipitation and access to nearest city) and its 95% pointwise credible interval provided by the BART-SIMP (blue) and BART (red) methods. The blue lines and grey bands represent the locally-weighted smoother function and its 95% pointwise confidence interval.

of BART into continuous spatial models and leveraged the flexibility of BART to model and predict the complex structure to allow nonlinear covariates and interactions. We have developed a sampling-based inference algorithm for BART-SIMP based on the INLA-within-MCMC technique.

BART-SIMP has a number of limitations which require more investigation. A cross-validation analysis showed that the BART-SIMP method yields average coverage rates closer to the nominal coverage compared to other methods while having poorer estimation performance than other methods when the covariate signal is weak. This suggests that when the covariate signal in the dataset is not strong compared to the spatial signal then it is not worth attempting any flexible covariate modeling. This is supported by our simulation studies. We hope our study provides inspiration for future spatial

modeling studies with complex covariate patterns. Another potential extension is to consider non-Gaussian likelihood models where the outcome variable is discrete (i.e. counts and proportions) using the Pólya–Gamma data augmentation technique (Polson et al., 2013).

The key challenge when combining spatial modeling with machine learning technique is estimating uncertainty appropriately. In general, uncertainty estimation is difficult with machine learning techniques, since the bootstrap does not work in many instances. (For example, with sparse estimators this occurs because the limiting distribution is complex and may not be continuous, Dezeure et al., 2015). When combining ML techniques with spatial models, this aspect becomes even more challenging.

# References

Battese, G. E., Harter, R. M., and Fuller, W. A. (1988). An error-components model for prediction of county crop areas using survey and satellite data. *Journal of the American Statistical Association*, 83(401):28–36. 3

Breiman, L. (2001). Random forests. *Machine Learning*, 45(1):5–32. 14

Burstein, R., Henry, N. J., Collison, M. L., Marczak, L. B., Sligar, A., Watson, S., Marquez, N., Abbasalizad-Farhangi, M., Abbasi, M., Abd-Allah, F., et al. (2019). Mapping 123 million neonatal, infant and child deaths between 2000 and 2017. *Nature*, 574(7778):353–358. 1

Chipman, H. A., George, E. I., and McCulloch, R. E. (1998). Bayesian CART model search. *Journal of the American Statistical Association*, 93(443):935–948. 8

Chipman, H. A., George, E. I., and McCulloch, R. E. (2010). BART: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4(1):266–298. 2, 5, 6, 7, 8

Davies, M. M. and Van Der Laan, M. J. (2016). Optimal spatial prediction using ensemble machine learning. *The International Journal of Biostatistics*, 12(1):179–201. 2

Daw, R. and Wikle, C. K. (2023). REDS: Random ensemble deep spatial prediction. *Environmetrics*, 34(1):e2780. 2

Dezeure, R., Bühlmann, P., Meier, L., and Meinshausen, N. (2015). High-dimensional inference: confidence intervals, p-values and r-software hdi. *Statistical Science*, pages 533–558. 17

Diggle, P. J. and Giorgi, E. (2019). *Model-Based Geostatistics for Global Public Health: Methods and Applications*. CRC Press. 1

Eddelbuettel, D. (2013). *Seamless R and C++ Integration with Rcpp*. Springer. 9

Eddelbuettel, D. and Balamuta, J. J. (2018). Extending R with C++: a brief introduction to Rcpp. *The American Statistician*, 72(1):28–36. 9

Eddelbuettel, D. and François, R. (2011). Rcpp: Seamless R and C++ Integration. *Journal of Statistical Software*, 40:1–18. 9

Fay, R. E. and Herriot, R. A. (1979). Estimates of income for small places: an application of James-Stein procedures to census data. *Journal of the American Statistical Association*, 74(366a):269–277. 3

Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, pages 1189–1232. 10, 14

Fuglstad, G.-A., Simpson, D., Lindgren, F., and Rue, H. (2019). Constructing priors that penalize the complexity of Gaussian random fields. *Journal of the American Statistical Association*, 114(525):445–452. 7

Georganos, S., Grippa, T., Niang Gadiaga, A., Linard, C., Lennert, M., Vanhuysse, S., Mboga, N., Wolff, E., and Kalogirou, S. (2021). Geographical random forests: a spatial extension of the random forest algorithm to address spatial heterogeneity in remote sensing and population modelling. *Geocarto International*, 36(2):121–136. 2

Gneiting, T. and Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378. 11

Goldstein, A., Kapelner, A., Bleich, J., and Pitkin, E. (2015). Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics*, 24(1):44–65. 14

Gómez-Rubio, V. and Rue, H. (2018). Markov chain monte carlo with the integrated nested laplace approximation. *Statistics and Computing*, 28(5):1033–1051. 2, 7

Hubin, A. and Storvik, G. (2016). Estimating the marginal likelihood with Integrated nested Laplace approximation (INLA). *arXiv preprint arXiv:1611.01450*. 5

Kassie, G. W. and Workie, D. L. (2019). Exploring the association of anthropometric indicators for under-five children in Ethiopia. *BMC Public Health*, 19(1):1–6. 3

Knorr-Held, L. and Rue, H. (2002). On block updating in Markov random field models for disease mapping. *Scandinavian Journal of Statistics*, 29(4):597–614. 5

Krueger, R., Bansal, P., and Buddhavarapu, P. (2020). A new spatial count data model with Bayesian additive regression trees for accident hot spot identification. *Accident Analysis and Prevention*, 144:105623. 2

LeSage, J. P. and Kelley Pace, R. (2007). A matrix exponential spatial specification. *Journal of Econometrics*, 140(1):190–214. Analysis of spatially dependent data. 2

Lindgren, F., Rue, H., and Lindström, J. (2011). An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(4):423–498. 5

Lindström, J., Szpiro, A. A., Sampson, P. D., Oron, A. P., Richards, M., Larson, T. V., and Sheppard, L. (2014). A flexible spatio-temporal model for air pollution with spatial and spatio-temporal covariates. *Environmental and Ecological Statistics*, 21:411–433. 1

Macharia, P. M., Giorgi, E., Thuranira, P. N., Joseph, N. K., Sartorius, B., Snow,

R. W., and Okiro, E. A. (2019). Sub national variation and inequalities in under-five mortality in Kenya since 1965. *BMC Public Health*, 19(1):1–12. 1

Martins, T. G., Simpson, D., Lindgren, F., and Rue, H. (2013). Bayesian computing with inla: new features. *Computational Statistics and Data Analysis*, 67:68–83. 9

Mei, Z. and Grummer-Strawn, L. M. (2007). Standard deviation of anthropometric Z-scores as a data quality assessment tool using the 2006 WHO growth standards: a cross country analysis. *Bulletin of the World Health Organization*, 85:441–448. 3

Müller, P., Shih, Y.-C. T., and Zhang, S. (2007). A spatially-adjusted Bayesian additive regression tree model to merge two datasets. *Bayesian Analysis*, 2(3):611–633. 2

Osgood-Zimmerman, A., Millear, A. I., Stubbs, R. W., Shields, C., Pickering, B. V., Earl, L., Graetz, N., Kinyoki, D. K., Ray, S. E., Bhatt, S., et al. (2018). Mapping child growth failure in Africa between 2000 and 2015. *Nature*, 555(7694):41–47. 2

Polson, N. G., Scott, J. G., and Windle, J. (2013). Bayesian inference for logistic models using pólya–gamma latent variables. *Journal of the American statistical Association*, 108(504):1339–1349. 17

Rao, J. N. and Molina, I. (2015). *Small Area Estimation*. John Wiley & Sons. 3

Ren, Z., Zhu, J., Gao, Y., Yin, Q., Hu, M., Dai, L., Deng, C., Yi, L., Deng, K., Wang, Y., Li, X., and Wang, J. (2018). Maternal exposure to ambient PM10 during pregnancy increases the risk of congenital heart defects: Evidence from machine learning models. *Science of The Total Environment*, 630:1–10. 2

Rue, H., Martino, S., and Chopin, N. (2009). Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(2):319–392. 1, 5

Shi, T., Hu, X., Guo, L., Su, F., Tu, W., Hu, Z., Liu, H., Yang, C., Wang, J., Zhang, J., et al. (2021). Digital mapping of zinc in urban topsoil using multisource geospatial data and random forest. *Science of the Total Environment*, 792:148455. 2

Spanbauer, C. and Sparapani, R. (2021). Nonparametric machine learning for precision medicine with longitudinal clinical trials and Bayesian additive regression trees with mixed models. *Statistics in Medicine*, 40(11):2665–2691. 2, 9

Utazi, C. E., Thorley, J., Alegana, V. A., Ferrari, M. J., Takahashi, S., Metcalf, C. J. E., Lessler, J., and Tatem, A. J. (2018). High resolution age-structured mapping of childhood vaccination coverage in low and middle income countries. *Vaccine*, 36(12):1583–1591. 1

Uwiringiyimana, V., Osei, F., Amer, S., and Veldkamp, A. (2022). Bayesian geostatistical modelling of stunting in Rwanda: risk factors and spatially explicit residual stunting burden. *BMC Public Health*, 22(1):1–14. 1

Zeraatpisheh, M., Garosi, Y., Reza Owliaie, H., Ayoubi, S., Taghizadeh-Mehrjardi, R., Scholten, T., and Xu, M. (2022). Improving the spatial prediction of soil organic carbon using environmental covariates selection: A comparison of a group of environmental covariates. *CATENA*, 208:105723. 1