

1 Comparing Emulation Methods for a High-resolution Storm Surge Model

2 Grant Hutchings, Bruno Sansó, James Gattiker, Devin Francom,
3 and Donatella Pasqualini
4

5 Abstract

6 The availability of powerful computing resources has led scientists to increasingly utilize
7 simulation as a research tool. The statistical analysis of simulations, referred to as computer
8 experiments, has similarly grown. Gaussian Process (GP) models have proven themselves
9 exceptionally useful in this domain and have become a standard methodology for emulation
10 of simulator response. However, with moderately large training data, GP's require careful
11 implementation to scale appropriately. There are a number of reasonable emulation methods
12 available from ready to use software packages. In this paper we compare four such models:
13 BASS; BART; SEPIA; and RobustGaSP, by applying them to high-resolution hurricane in-
14 undation (flooding) data obtained from the Sea, Lake, and Overland Surges from Hurricanes
15 (SLOSH) simulator. Both SEPIA and RobustGaSP are based on Gaussian Process modeling,
16 while BASS implements a model based on adaptive splines, and BART is based on sums of
17 regression trees. We will describe the modeling strategies implemented in these four packages,
18 which run on R and Python, and then compare them in terms of computation time and a
19 variety of predictive metrics. The four models included in this comparison study were chosen
20 for their proven and distinct methodologies, their availability through easily accessible soft-
21 ware, and their ability to quantify prediction uncertainty in the context of our application.
22 The data in our case study form a large spatial grid with millions of response values. We
23 find that SEPIA and RobustGaSP provide exceptional predictive power, but cannot scale to
24 accommodate computer experiments as large as the one considered in this paper as effectively
25 as BASS and BART.

26 1. Introduction.

27 **1.1. Background.** The study of complex physical systems is often limited by the acqui-
28 sition of experimental data which can be expensive or even impossible to gather in many
29 fields. As a result, scientists turn to simulation to supplement experimental data, to gain un-
30 derstanding, and to make predictions (Sacks et al., 1989). Aided by advances in computing,
31 simulators based on mathematical models of physical processes have become a fundamental
32 tool to obtain scientifically motivated representations of a system of interest. Simulators de-
33 pend on a number of inputs (parameters) that control their behavior. We refer to the set

34 of all possible input values as the parameter space. To obtain a realistic description of the
35 system, and a good understanding of the simulator’s capabilities, we must analyze simulation
36 output at collections of points in the parameter space.

37 Computer simulations do not completely remedy the challenges associated with experimen-
38 tal data. The information that can be obtained from simulations is limited by the feasibility
39 of running the simulation at a given point in the parameter space. Depending on the sys-
40 tem of interest, simulations can take hours or even days to run for a given combination of
41 parameter values. This may make it impossible to do an exhaustive direct exploration of the
42 space, a problem that is compounded as the parameter space increases in size or dimension.
43 A variety of examples can be found in [Sacks et al. \(1989\)](#). Additionally, simulators are often
44 deterministic, which typically means that the primary source of uncertainty when running the
45 simulator is parameter uncertainty. Statistical models of these computer simulations, referred
46 to here as emulators, are designed to solve these problems.

47 A seminal work in the literature of computer experiments ([Sacks et al., 1989](#)) showed how
48 a Gaussian Process could be used to build a predictor with uncertainty quantification. GP’s
49 became and remain a common approach for emulating computer simulations with a vast lit-
50 erature encompassing a variety of approaches. The main purpose of an emulator is to provide
51 predictions at untried parameter settings with an estimate of the associated uncertainty, and
52 to do it much faster than running the actual computer simulation ([Salter and Williamson,
53 2016a](#)). GP’s are therefore very desirable for their predictive power and straightforward uncer-
54 tainty quantification. They are however not always practical; Gaussian Processes are limited
55 by the computational bottleneck of covariance matrix inversion which limits applicability to
56 large data. Many recent methods such as LaGP ([Gramacy and Apley, 2015](#)), TGP ([Gramacy
57 and Lee, 2008](#)), GPvecchia ([Katzfuss and Guinness, 2021](#)), and RobustGaSP ([Gu et al., 2017](#))
58 aim to tackle this scalability issue. More recently, competitive alternatives to GP’s have been
59 proposed such as BASS ([Francom and Sansó, 2020](#)) which implements an adaptive spline
60 model, and BART ([Sparapani et al., 2021](#)) which uses additive regression trees. These meth-
61 ods similarly provide accurate prediction with simple uncertainty quantification and often a
62 smaller computational footprint.

63 The analysis presented here will compare four emulation methods on simulated hurri-
64 cane induced flooding in the Delaware Bay. The simulator considered in this study allows
65 researchers to learn about hurricane flood risk to critical infrastructure on an accelerated
66 timeline, and explore different hurricane scenarios by changing the simulation parameters.
67 The comparison here is motivated by the need for emulation in further analysis based on this
68 model, as well as potential similar future models.

69 The goals of this study are to quantify the accuracy of predictions and understand the
70 computational requirements of each method for a range of training set sizes. In doing so, we
71 aim to understand how training set size effects predictions and run time. Additionally we

72 will compare the variable importance options given by each method. Investigating variable
73 importance for hurricane flooding models helps researchers understand which qualities of a
74 hurricane or a particular area are most influential in determining inland flooding. Some
75 emulators allow for spatially resolved variable importance and variance-based assessment of
76 importance (e.g., via the Sobol decomposition (Sobol, 2001)), which both benefit analyses
77 involving highly multivariate emulators.

78 The remainder of the paper is structured as follows: In Section 1.2 we give a brief overview
79 of the four methods included in our study and explain why they were chosen; Section 2 is
80 an overview of the simulations from SLOSH; Section 3 describes each of the four emulator
81 formulations; Section 4 presents our comparison study, highlighting a variety of predictive
82 metrics and scores; Section 5 gives an overview of the variable importance built into each
83 package; and we conclude with a discussion of our findings and recommendations to the
84 reader in Section 6.

85 **1.2. Emulation Methods.** The emulation methods we have chosen implement very dif-
86 ferent statistical models, all of which have proven themselves a reasonable choice for similarly
87 structured spatial data. We will consider two GP based models, SEPIA (Gattiker et al.,
88 2020b) and RobustGaSP; SEPIA fits a collection of independent GP models to coefficients
89 of an orthogonal basis representation of the simulation response data, while RobustGaSP im-
90 plements a Many Single approach, fitting an independent GP to each spatial location. We
91 also include the two non-GP based models mentioned above; BASS and BART. These four
92 models cover some diverse modeling strategies, but in no way cover the full spectrum of em-
93 ulation methodologies. While recognizing the limitations of only considering four models, we
94 would like to highlight the fact that this study customized implementation and computation
95 appropriately for each method for the application, an approach that represents a significant
96 investment in investigator and computational resources compared to a investigation based on
97 relatively limited customization and tailored test problems. Emulator comparisons have been
98 done in the past, often comparing on a host of test functions with relatively small amounts
99 of data, or focusing on parameter calibration rather than strictly emulation (e.g. Salter and
100 Williamson (2016b), Erickson et al. (2018)). The comparison here is motivated by the re-
101 quirements of this application which poses particular problems that are relevant to spatial
102 environmental modeling. What we present is a comparison which focuses only on a few mod-
103 els in greater detail, in an application driven big-data setting. This, to our knowledge, is not
104 prevalent in the literature.

105 The first of the four methods that we consider in this paper is the Simulation Enabled
106 Prediction Inference and Analysis (SEPIA) software that implements the Gaussian process
107 model described in Higdon et al. (2008). This model was originally implemented at Los Alamos
108 National Laboratory as the MATLAB code GPMSA (Gattiker et al., 2020a) and in 2020 was
109 refurbished and translated to python as SEPIA. SEPIA makes use of a basis representation,

110 typically empirical orthogonal functions (EOF) (also known as principle components analysis),
111 of the data to fit a Gaussian process to each of the basis coefficients. This is a tried and
112 true methodology for spatial modeling that has seen much success in the literature and in
113 applications.

114 Our implementation of Bayesian Adaptive Spline Surfaces (BASS) similarly makes use of
115 a basis representation, but takes a wholly different approach to modeling basis coefficients by
116 using adaptive splines. BASS has been recently applied to large spatial data from computer
117 experiments and has shown great results (see, for example, [Francom et al., 2019](#)).

118 The implementation considered in this work of Bayesian Additive Regression Trees (BART) ■
119 once again makes use of a basis representation where each basis coefficient is fit using an inde-
120 pendent BART model. The BART package does not inherently manage multivariate response
121 through basis representation (as in SEPIA and BASS), and so we extend the functionality
122 by explicitly supplying an EOF basis. The BART model fits the EOF weights and the pre-
123 dictions are expanded into the native space. This allows a more direct comparison to other
124 methods. We have explored this implementation in the past ([Francom et al., 2020](#)). Treed
125 models have seen success in the literature for their speed and flexibility, and BART has proven
126 to be effective in settings similar to the one considered in this paper, such as a recent analysis
127 of airborne particulate data over California ([Zhang et al., 2020](#)). Preliminary comparisons
128 of BASS and BART in [Francom et al. \(2019\)](#) showed that both approaches can be highly
129 accurate and efficient.

130 The fourth method considered in this work consists of Robust Gaussian Stochastic Process
131 Emulation (RobustGaSP) which handles multivariate response by fitting a GP to each point
132 in space, rather than reducing the modeling dimension through a linear projection as the other
133 methods in this comparison. This is made computationally feasible by both parallel compu-
134 tation, and the assumption of shared range parameters for all GP's. RobustGaSP does not
135 make use of Markov-chain Monte Carlo (MCMC) for model fitting like the other three models.
136 Instead parameters are fit using numerical optimization of marginal posterior distributions.
137 These major model differences make this an interesting inclusion to our comparison study.
138 RobustGaSP has also shown promising results on large scale computer model emulation of
139 large volcanic flow simulations ([Gu and Berger, 2016](#)).

140 Additionally, we include a simple linear model on the coefficients of an orthogonal basis
141 representation as a baseline to gauge the improvements provided by these complex models.

142 The models considered in this paper all show accurate predictions using quite different
143 methodologies. We will give a more detailed description of each model in Section 3.

144 **2. Simulator and Dataset.** The Sea, Lake, and Overland Surges from Hurricanes (SLOSH) ■
145 simulator ([Jelesnianski et al., 1992](#)) is a computer code developed by the National Weather
146 Service to estimate storm surge heights from hurricanes. Storm surge height is defined as the
147 maximum water height due to a hurricane at any single location. Our data consists of an

148 ensemble of 4,000 runs from the SLOSH simulator, corresponding to 4,000 simulated storms.
 149 Each storm in the ensemble is defined by a unique set of five input parameters:

- 150 • sea level rise in the year 2100 (lower: -20; upper: 350; units: cm)
- 151 • heading of the eye of the storm when it made landfall (lower: 204.0349; upper:
 152 384.0244; units: degrees, north is 0/360)
- 153 • velocity of the eye of the storm when it made landfall (lower: 0; upper: 40; units:
 154 knots)
- 155 • minimum air pressure of the storm when it made landfall (lower: 930; upper: 980;
 156 units: millibars)
- 157 • latitude of the eye of the storm when it made landfall (lower: 38.32527; upper:
 158 39.26811; units: degrees)

159 Input parameters for the ensemble use a space-filling Latin hypercube design over our five
 160 dimensional parameter space. Models are trained on subsets of this ensemble and tested on
 161 storms outside of the training sets.

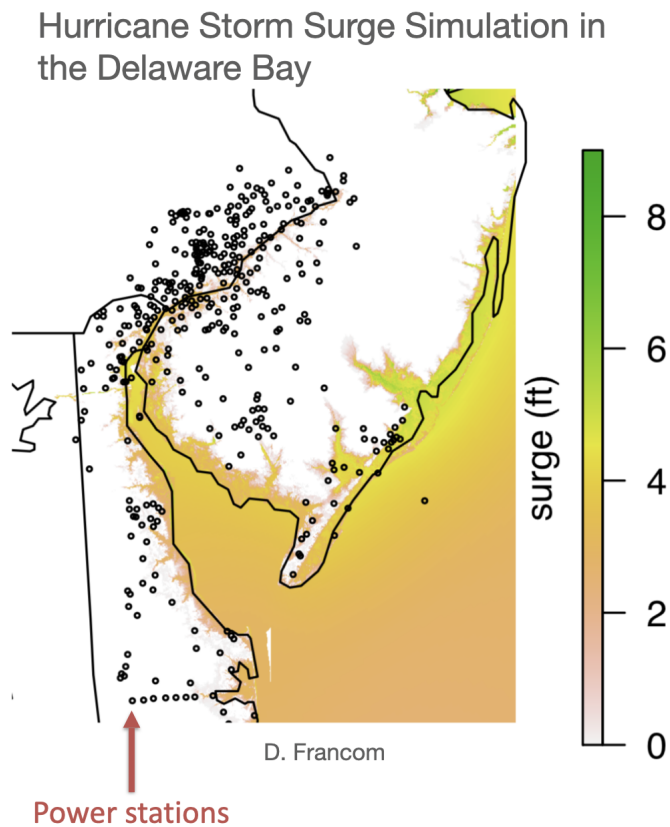


Figure 1: Surge output map from SLOSH

162 Our interest lies in prediction of hurricane-induced flooding in the Delaware Bay. Under
 163 our model setup, one output from SLOSH is a $4,520 \times 5,115$ grid of storm surge heights for
 164 each of the 23,119,800 locations. Figure 1 presents a spatial map of SLOSH output for a given
 165 combination of input parameters. This large number of spatial locations presents a formidable
 166 computational challenge which is fortunately eased by the fact that the majority of the points
 167 on the grid are far enough inland that there is no flooding for any of the 4,000 simulations.
 168 By modeling only cells which take non-zero values in at least one of the simulations we reduce
 169 the size of the field to 3,500,000 locations.

170 Accurate prediction of flooding is important for a variety of reasons including displacement
 171 of residents, and property/infrastructure damage. One area of specific interest for this project
 172 is possible damage to infrastructure, specifically power stations displayed as black dots in
 173 Figure 1. Power stations in this area are often fortified to handle four feet of flood water, any
 174 more can lead to catastrophic damage. We are therefore interested in the emulators' ability
 175 to accurately predict that a surge has reached four feet, as this information is very valuable
 176 for determining if an intervention (station shut down) is necessary due to an incoming storm.
 177 We will discuss predictions around this threshold of four feet in more detail in Section 4.

178 **3. Model Formulation.** The emulation problem considered in this paper presents the
 179 challenge of building emulators that are able to handle 4,000 runs from SLOSH, each with
 180 $n_y = 3.5 \times 10^6$ response values. One very common approach to reduce the dimension of
 181 a problem like this is to decompose the data into principal components (PCs; Ramsay and
 182 Silverman (1997)) using a singular value decomposition (SVD). The output vector $\mathbf{y}(\mathbf{x}) \in \mathbb{R}^{n_y}$
 183 from one SLOSH run, corresponding to inputs $\mathbf{x} \in \mathbb{R}^p$ can be represented on a set of orthogonal
 184 basis functions as $\sum_{j=1}^{\infty} w_j(\mathbf{x}) \mathbf{b}_j$ where $\mathbf{b}_j \in \mathbb{R}^{n_y}$ captures the spatial variation. By stacking
 185 the output obtained from each of the m storms in the training set, we obtain the matrix
 186 $\mathbf{Y} \in \mathbb{R}^{m \times n_y}$, which we center by subtracting the mean storm. \mathbf{Y}_{ik} then corresponds to the
 187 standardized output from storm i at location k . We compute $SVD(\mathbf{Y}) = \mathbf{U} \mathbf{D} \mathbf{V}^T$ where
 188 \mathbf{U}, \mathbf{V} are orthogonal matrices and \mathbf{D} is a diagonal matrix of singular values. \mathbf{V}^T and $\mathbf{U} \mathbf{D}$
 189 store the empirical $w_j(\mathbf{x})$ and \mathbf{b}_j respectively. We choose to truncate the sum at n_{pc} principal
 190 components, so that 99% of the variation in the data is captured by the basis representation.
 191 The number of principal components used varies by training set. The smallest set with only 50
 192 storms requires just $n_{pc} = 14$ principal components while the largest set with 3,636 requires
 193 $n_{pc} = 24$. The power of this decomposition comes from the fact that, rather than fitting
 194 an emulator to all n_y response values, we only need to fit n_{pc} scalar response models to
 195 the coefficients $w_j(\mathbf{x})$, which results in drastic computational savings. We utilize the identical
 196 matrix decomposition when fitting BASS, BART, SEPIA, and the linear model. RobustGaSP
 197 does not make use of this representation, as discussed. In Subsections 3.1-3.3 we will suppress
 198 the subscript j for simplicity and refer to an arbitrary $w_j(\mathbf{x})$ as $w(\mathbf{x})$.

199 **3.1. Simulation Enabled Prediction and Inference (SEPIA).** SEPIA is a python code
 200 developed by Jim Gattiker, Natalie Klein, Grant Hutchings and Earl Lawrence at Los Alamos
 201 National Laboratory (Gattiker et al., 2020b) and implements the model described in Higdon
 202 et al. (2008), with extensions. Here we use the emulator component only, without SEPIA’s
 203 full model calibration functionality. By utilizing the orthogonal basis representation described
 204 above, a Gaussian process is fit to each basis function coefficient $w(\mathbf{x})$.

$$205 \quad (3.1) \quad w(\mathbf{x}) \sim GP(0, \Sigma); \quad \Sigma = \sigma_n^2 \mathbf{I} + \sigma_p^2 \mathbf{C}$$

206 where $\mathbf{C}_{kl} = \exp\{-\frac{1}{2} \sum_{i=1}^p \beta_i (\mathbf{x}_{ki} - \mathbf{x}_{li})^2\}$ is the matrix obtained by applying the negative
 207 exponential squared (“Normal kernel”) correlation function to each pair of inputs, which is
 208 parameterized by length scale β . Σ incorporates process variance σ_p^2 and includes a noise
 209 process with variance σ_n^2 . This is a Bayesian model with priors on $\beta, \sigma_p^2, \sigma_n^2$. For a full
 210 model specification including discussion of priors, refer to Higdon et al. (2008). The resulting
 211 posterior distributions are explored via MCMC.

212 **3.2. Bayesian Adaptive Spline Surfaces (BASS).** BASS is an R package to fit Bayesian
 213 adaptive spline surfaces (Francom and Sansó, 2020). It implements a Bayesian version of
 214 multivariate adaptive regression splines (Friedman, 1991). Similar to the approach we took
 215 with SEPIA, we make use of a basis representation for the SLOSH output. BASS models each
 216 $w(\mathbf{x})$ as

$$217 \quad (3.2) \quad w(\mathbf{x}) = a_0 + \sum_{m=1}^M a_m Z_m(\mathbf{x}) + \epsilon(\mathbf{x}), \quad \epsilon(\mathbf{x}) \sim N(0, \sigma^2)$$

218 where a_0, a_1, \dots, a_M are constants and Z_1, \dots, Z_M are basis functions learned from the data.
 219 The basis functions have the form

$$220 \quad (3.3) \quad Z_m(\mathbf{x}) = \prod_{k=1}^{K_m} g_{km} [s_{km} \max(0, x_{v_{km}} - t_{km})]^\alpha$$

221 where $s_{km} \in \{-1, 1\}$ is the sign, $t_{km} \in [0, 1]$ is a knot, v_{km} selects a covariate, K_m is the degree
 222 of interaction and $g_{km} = [(s_{km} + 1)/2 - s_{km} t_{km}]^\alpha$ is a constant that makes the basis function
 223 have a maximum of one. The exponent α defines the degree of the polynomial splines. Note
 224 that variables can only be used once in each basis function.

225 To fit this model we need to estimate $\theta = \{\sigma^2, M, \mathbf{a}, \mathbf{K}, \mathbf{s}, \mathbf{t}, \mathbf{v}\}$. This is done via a
 226 reversible jump MCMC (RJMCMC) algorithm. For specifics on priors and the RJMCMC
 227 algorithm see Francom and Sansó (2020).

228 **3.3. Bayesian Additive Regression Trees (BART).** BART is a treed model with strong
 229 predictive power for non-linear responses. A recent example is the use of BART for spatial
 230 modeling of ambient fine particulate matter pollution (PM_{2.5}) over California (Zhang et al.,

231 2020). As detailed in Chipman et al. (2010), BART is a sum of trees model where scalar
 232 output $w(\mathbf{x})$ is approximated as

$$233 \quad (3.4) \quad w(\mathbf{x}) = \sum_{i=1}^I g(\mathbf{x}|T_i, M_i) + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

234 where each T_i is a regression tree that can incorporate one or more of the p inputs, corre-
 235 sponding to main and interaction effects. A tree T utilizing $\mathbf{x}_t \subseteq \mathbf{x}$ consists of a set of interior
 236 nodes with binary decision rules, and a set of leaf nodes containing parameter estimates. Let
 237 $M = \{\mu_1, \dots, \mu_b\}$ be the parameter estimates associated with the leaf nodes. The interior
 238 decision rules are binary splits of the predictor space, either $\mathbf{x}_t \in A$ or $\mathbf{x}_t \notin A$ where A is a
 239 subset of the range of \mathbf{x}_t . Then any fixed \mathbf{x}_t^* is assigned a μ^* by the function $g(\mathbf{x}|T, M)$ based
 240 on the sequence of decision rules leading to a leaf node.

241 This additive structure endows BART with a high degree of flexibility when the number
 242 of trees is large. This does however come at the price of complexity. BART needs to estimate
 243 $\{(T_1, M_1), \dots, (T_I, M_I), \sigma\}$ for I trees where T_i and M_i are not single parameters, but an entire
 244 tree structure fit with a set of decision rules, and a set of terminal nodes respectively. A
 245 backfitting MCMC algorithm is used for posterior sampling, which is designed to efficiently
 246 sample the many parameters in the additive tree structure. As a result, BART provides great
 247 flexibility with a relatively low computational cost. A key component of the model is a regu-
 248 larization prior which forces the effect from each tree to be small. This prevents individual tree
 249 effects from dominating the additive structure. Once posterior draws $(T_1^*, M_1^*), \dots, (T_I^*, M_I^*)$
 250 are available, predictions f^* can be obtained as

$$251 \quad (3.5) \quad f^*(\cdot) = \sum_{i=1}^I g(\cdot|T_i^*, M_i^*)$$

252 (Sparapani et al., 2021).

253 **3.4. Robust Gaussian Stochastic Process Emulation (RobustGaSP).** RobustGaSP (Gu
 254 et al., 2017) is a GP-based method that avoids the use of the basis function representation
 255 that we have used for SEPIA, BASS and BART. Also, unlike the other three models the
 256 estimation procedure relies on marginal likelihood optimization rather than MCMC. This
 257 has its drawbacks when it comes to uncertainty quantification as confidence bounds must be
 258 estimated using distributional assumptions. On the other hand it avoids the iterative sampling
 259 involved in MCMC, which incurs relatively large computational cost and memory footprint.

260 RobustGaSP implements a computationally feasibly alternative to the Many Single (MS)
 261 emulation approaches (Conti and O’Hagan, 2010; Lee et al., 2011, 2012). Individual emulators
 262 are fit to each coordinate of the output, which, in the context of our case study, consists of
 263 n_y independent Gaussian process emulators. Each emulator has its own mean function and

264 variance, but they all share the same correlation parameters $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$, which are
 265 estimated from the joint marginal likelihood of all emulators (Gu and Berger, 2016).

266 Let $i = 1, \dots, n_y$ index the locations so that $y_i(\mathbf{x})$ denotes the scalar response at location
 267 i with inputs \mathbf{x} . $y_i(\mathbf{x})$ is modeled with the Gaussian Process

$$268 \quad (3.6) \quad y_i(\mathbf{x}) \sim GP(\mu_i(\mathbf{x}), \sigma_i^2 c(\mathbf{x}, \mathbf{x}')), \quad i = 1, \dots, k$$

269 where $\mu_i(\mathbf{x})$ is the location specific mean function, σ_i^2 the location specific variance, and
 270 $c(\mathbf{x}, \mathbf{x}')$, by default, is the product of p Matèrn 5/2 correlation functions, each with its own
 271 range parameter $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$. Then for m runs of the simulator at inputs $\mathbf{x}_1, \dots, \mathbf{x}_m$ we
 272 have the multivariate likelihood

$$273 \quad (3.7) \quad (y_i(\mathbf{x}_1), \dots, y_i(\mathbf{x}_m) | \boldsymbol{\mu}_i, \sigma_i^2, \boldsymbol{\Sigma}) \sim \text{MVN}((\mu_{i\mathbf{x}_1}, \dots, \mu_{i\mathbf{x}_m}), \sigma_i^2 \boldsymbol{\Sigma})$$

274 where $\boldsymbol{\Sigma}$ is the correlation matrix obtained by applying $c(\mathbf{x}, \mathbf{x}')$ to each pair of input vectors.
 275 The mean function is modelled using a linear regression, $\mu_i(\mathbf{x}) = \sum_{l=1}^L h_l(\mathbf{x})\theta_l$, with basis func-
 276 tions $\mathbf{h}_i(\mathbf{x}) = (h_{i1}(\mathbf{x}), \dots, h_{iL}(\mathbf{x}))$ and unknown regression parameters $\boldsymbol{\theta}_i = (\theta_{i1}, \dots, \theta_{iL})$. An
 277 important aspect of this approach is the definition of the prior for the model parameters. This
 278 consists of the product of a standard objective prior is for the mean and variance parameters
 279 (Gu and Berger, 2016),

$$280 \quad (3.8) \quad \pi^R(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{n_y}, \sigma_1^2, \dots, \sigma_{n_y}^2) \propto \frac{1}{\prod_{i=1}^{n_y} \sigma_i^2}$$

281 and a jointly robust (JR) prior applied to the correlation parameters $\boldsymbol{\gamma}$. This prior was
 282 introduced in Gu (2018) and is called jointly robust because it cannot be written as the
 283 product of marginal priors and is robust in marginal posterior mode estimation.

284 First consider reparameterizing to the inverse range parameters $\beta_j = 1/\gamma_j, j = 1, \dots, p$.
 285 Then the JR prior is defined as

$$286 \quad (3.9) \quad \pi^{JR}(\beta_1, \dots, \beta_p) = C_0 \left(\sum_{l=1}^p C_l \beta_l \right)^\alpha \exp \left\{ -b \left(\sum_{l=1}^p C_l \beta_l \right) \right\},$$

287 where $C_0 = \frac{(p-1)! b^{a+p} \prod_{l=1}^p C_l}{\Gamma(a+p)}$, $a > -(p+1)$, $b > 0$ and $C_l > 0$ are parameters. We use the
 288 default values for these parameters; $a = 0.2$, $b = n^{-1/p}(a+p)$. The default values for C_l are
 289 not clearly given in the documentation. As we will discuss in Section 5, this prior facilitates
 290 the form of variable importance provided by the package.

291 The posterior distribution resulting from this model formulation is marginally optimized
 292 to obtain parameter estimates.

293 **3.5. Linear Model (LM).** For a baseline comparison, we include a simple linear model
294 on the EOF basis coefficients $w(\mathbf{x})$ with the form

$$295 \quad (3.10) \quad w(\mathbf{x}) = \sum_{i=1}^p \beta_i x_i + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

296 where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$ are unknown regression coefficients which we determine using the
297 function "lm()" from base R (R Core Team, 2020).

298 **4. Comparison Study.** This section presents assessment of the four different emulators
299 on the basis of out-of-sample predictive accuracy and computational feasibility. Predictive
300 accuracy is assessed using scores including root-mean-squared error (RMSE), energy score, and
301 coverage. RMSE assesses the mean prediction, whereas the energy score and coverage assess
302 the uncertainty associated with predictions. We will organize our comparison of predictive
303 accuracy into two Subsections, one for assessing the accuracy of the mean, and the other
304 considering estimates to be used in uncertainty quantification. Our results will show that in
305 these metrics Gaussian Process based emulators (SEPIA and RobustGaSP) produce better
306 mean predictions, however they appear less accurate in their predicted uncertainty.

307 We would like to be able to train our models with as few storms as needed for accuracy,
308 while minimizing computation time and leaving more examples in the model test set. To
309 examine the impact of training set size for each emulator we consider seven different training
310 sets; 50, 100, 500, 1000, 1750, 2500, and 3636 storms. 3636 was chosen as the largest training
311 set size because it is the largest number that permits a testing set size of 10% of the training
312 set (364 testing storms). The largest training set was sampled randomly from the full 4000
313 storm ensemble, and subsequent training sets sampled randomly from this set of 3,636. While
314 randomly subsampling a space-filling design is not optimal, the same selection is used for each
315 emulator, affording fair comparison.

316 Our comparison study involves training each of the four emulators on each of the seven
317 training sets, and computing all prediction metrics on the testing set. All models are tested
318 on the same 364 storms. This allows estimation of the impact of training set size, and com-
319 parison of performance both within and between these training set sizes. Computation time
320 is compared across training set sizes revealing the scaling properties of each algorithm. Our
321 results underline an important and well known fact that Gaussian Process, while providing
322 exceptional predictive power, becomes prohibitive with large data-sets. This is evident in that
323 we were only able to fit SEPIA and RobustGaSP with a maximum of 1000 and 500 training
324 storms respectively. We will discuss this further in Section 4.3.

325 BASS, BART, and SEPIA all make use of MCMC for parameter estimation. For each
326 model we have chosen to collect 10,000 MCMC samples, and discard the first 9000 to eliminate
327 transient state (so-called "burn-in"). Because of the size of the spatial field, we thin the
328 remaining samples down to 50, driven by memory constraints on our computational resources.

329 To fully appreciate the memory challenge, recall that our testing set is 364 storms. To generate
330 predictions using all 1000 posterior samples requires a double precision matrix of size $(364 \times$
331 $3,500,000 \times 1000)$, which requires 10 terabytes of storage. We are limited on our platforms to
332 a more modest 500 gigabyte matrix resulting from the use of 50 samples. This is one of the
333 many challenges involving an application dataset of this size. We appreciate that given the
334 relatively small number of initial samples (10,000) and even smaller number retained samples
335 (50), there may be questions regarding the convergence and mixing of our initial chains, and
336 of how well the 50 samples represent the posterior distributions. These software packages do
337 not provide methods to quantify convergence or mixing, and it is infeasible for us (and in
338 general practice) to tackle this problem for each combination of emulator, training set, and
339 EOFs. The results should be viewed with the understanding that poor convergence/mixing
340 and issues due to small sample set are potentially present in predictive metrics of accuracy
341 and coverage. For a practitioner interested in assessing MCMC convergence, they may want
342 to pursue a thorough analysis of chains which we do not consider here. For those who require
343 this analysis, we would have to recommend reducing the computation by further reduction of
344 the spatial data to make investigation tractable.

345 The following Subsection will present results for a variety of predictive metrics which can
346 be used to compare the models.

347 **4.1. Predictive Accuracy: Mean.** In this section we will assess the accuracy of mean
348 predictions from each emulator. For MCMC based models, this is the mean over our 50
349 posterior predictive samples and for RobustGaSP, the mean is returned to us by the package.
350 Our assessment will consider RMSE, mean absolute error (MAE), and our own loss function
351 designed specifically for flood risk analysis.

352 Figure 2a shows boxplots of RMSE for each emulation method and for each training set
353 size at which they were run. Samples in each correspond to the 364 test storms dataset.
354 As expected, RMSE is generally decreasing with training set size. The plots show diminish-
355 ing returns, with a reasonable conclusion that a training set size greater than 1000 runs is
356 unnecessary to achieve best performance in RMSE. Additionally, the figures show that that
357 methods utilizing Gaussian process, SEPIA and RobustGaSP, tend to have the lowest RMSE
358 at each training set size. Furthermore, they produce comparable RMSE to BASS trained on
359 the full 3,636 storms. BASS and BART produce fairly similar results, lagging behind SEPIA
360 and RobustGaSP, with BASS slightly outperforming BART. The results for MAE are very
361 similar to RMSE and are reported in the supplementary material.

362 In Section 2 we noted that a flooding threshold of four feet is of special interest. This
363 number has real implications in that many power stations are fortified to withstand this level
364 of flood water. ¹ Therefore, it is desirable for an emulator to correctly predict flood level above

¹Different flood impact thresholds can be found in the literature. The four foot threshold is driven by our

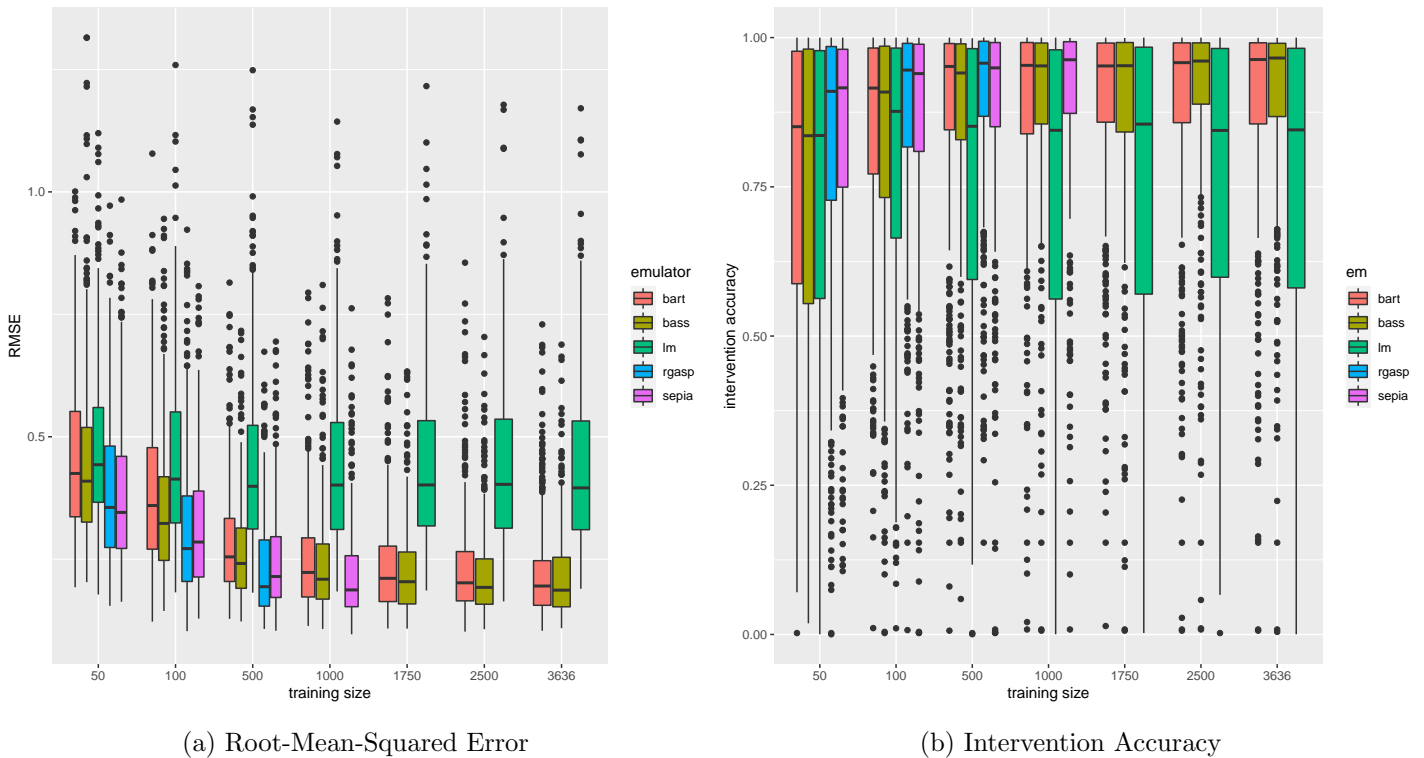


Figure 2: Predictive metrics for mean predictions by training set size.

365 four feet, a domain-relevant criterion for evaluation. We will evaluate this with standard em-
 366 ulation methods, rather than creating an emulator to satisfy the application-specific loss. To
 367 assess these emulators with respect to this feature, we consider the percentage of predictions
 368 that correctly indicate that an intervention is needed, which we call the intervention accuracy.
 369 To compute this metric for the mean prediction, we consider all cells in which the true SLOSH
 370 output is greater than four feet, and determine the percentage of cells in which the prediction
 371 is also greater than four feet. Figure 2b shows boxplots of our results where distributions are
 372 over the 364 testing storms. We see that SEPIA and RobustGaSP performance is better than
 373 BASS, BART, and the linear model at every training set size. Additionally, RobustGaSP with
 374 only 500 training storms is able to achieve indistinguishable performance to BASS with the
 375 full training set. SEPIA achieves a comparable performance with 1000 training storms. This
 376 is further evidence that GP-based methods are able to provide better mean predictions with
 377 less training data. This metric is especially interesting when viewed from a risk-management
 378 perspective; With SEPIA and RobustGaSP, we are less likely to miss an important interven-

application context of US infrastructure planning, and is indicative of threshold-based evaluation of emulators.

379 tion. There are a number of near zero values in Figure 2b which we found to be associated
 380 with storms for which an especially low number of locations reached the four foot threshold.
 381 One reason why this may result in low intervention accuracy is the smoothing associated with
 382 prediction. Fewer locations above the four foot threshold likely means those locations reside
 383 in smaller clusters which are more likely to be under-predicted due to smoothing effects.

384 **4.2. Predictive Accuracy: Uncertainty.** This section presents the results of predictive
 385 metrics which take uncertainty into consideration: *coverage probability*, *energy score*, and
 386 *interval score*.

387 **4.2.1. Coverage.** As all our emulators provide confidence intervals we are interested as-
 388 sessing in their level of coverage. In Figure 3a we present coverage probability distributions
 389 for a 95% interval over the 364 testing storms. Using the dashed red line at 0.95, we can see
 390 that the linear model, SEPIA, and RobustGaSP all consistently over-cover. BART tends to
 391 over-cover with small training sets and under-cover with larger training sets. BASS does the
 392 opposite, but seems to be consistently closest to the true 95% coverage. We will now extend
 393 our assessment of coverage by comparing the models using a score proposed in Gneiting and
 394 Raftery (2007), the interval score.

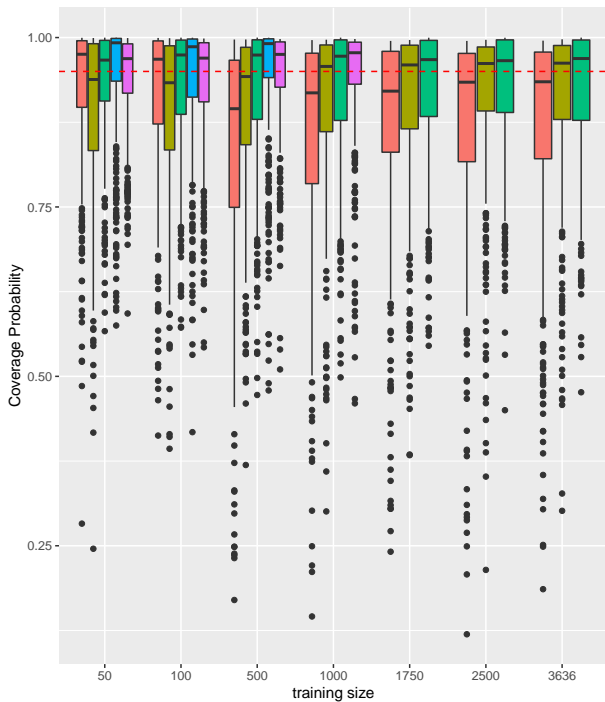
395 The interval score for confidence level α is defined as

$$396 \quad (4.1) \quad S_{\alpha}^{int}(l, u; x) = (u - l) + \frac{2}{\alpha}(l - x)\mathbb{1}\{x < l\} + \frac{2}{\alpha}(x - u)\mathbb{1}\{x > u\}.$$

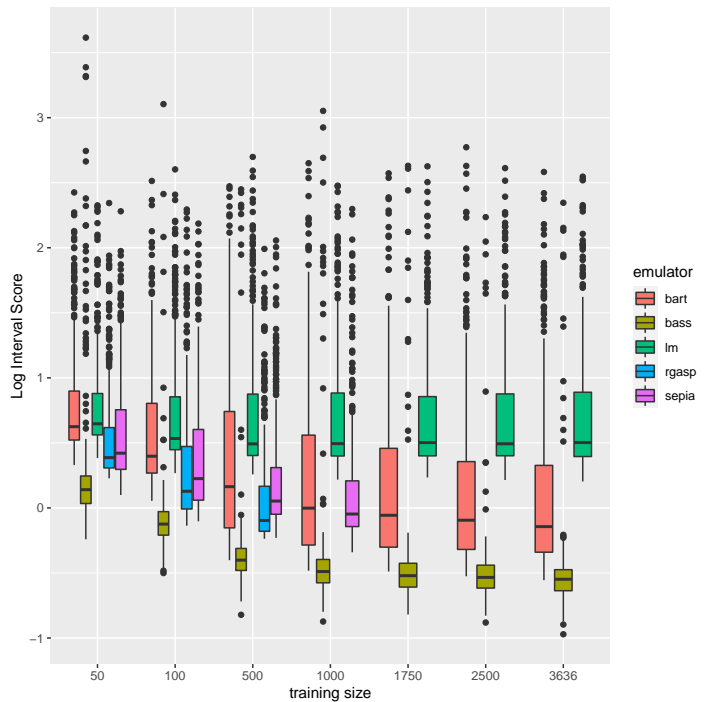
397 Where l, u are the lower and upper bounds of the $1 - \alpha$ confidence interval, and x is the true
 398 data value. This is a negative oriented score that is minimized at the width of the interval.
 399 The score then increases proportional to α if the true data value is outside the interval. This
 400 score provides more insight than coverage probability by consciously favoring models with the
 401 smallest possible intervals that still contain the data. In Figure 3b we present interval score
 402 distributions over the 364 testing storms where each storms score is an average over scores for
 403 each cell. BASS appears to be quite superior to the other emulators, while the linear model
 404 performs poorly in comparison. For small training sets, BART seems to do almost as poorly
 405 as the linear model, only catching up to SEPIA at 1000 training storms.

406 **4.2.2. Energy Score.** The energy score, a multivariate extension of the Continuous Rank
 407 Probability Score (CRPS) is proposed in Gneiting and Raftery (2007). This score takes
 408 into account not only the predictive accuracy of each sample from the posterior predictive
 409 distribution, but also the level of uncertainty in the distribution. For this reason, the CRPS
 410 and energy score have gained interest in recent literature as a model ranking mechanism
 411 (Heaton et al. (2018), Möller et al. (2013), Muniain and Ziel (2020)). With m draws from the
 412 posterior predictive distribution, $\tilde{Y} = \{\tilde{Y}_1, \tilde{Y}_1, \dots, \tilde{Y}_m\}$, we compute the energy score as

$$413 \quad (4.2) \quad es(Y, \tilde{Y}) = \frac{1}{m} \sum_{j=1}^m \|\tilde{Y}_j - Y\| - \frac{1}{2m^2} \sum_{j=1}^m \sum_{k=1}^m \|\tilde{Y}_j - \tilde{Y}_k\|,$$



(a) Coverage probability, 95% interval.



(b) Interval scores for 95% confidence level.

Figure 3: Coverage Metrics

414 where Y is the true response.

415 Results from Subsection 4.1 indicate that Gaussian Process models might be superior in
 416 terms of mean predictions. Interestingly, the energy score, which uses predictive samples
 417 rather than the mean tends to favor the tree and spline based models over the GP based
 418 models. So, while GP’s may provide very good mean predictions, results from this section
 419 indicate that they may not provide competitive uncertainty quantification to BASS.

420 **4.3. Computational Feasibility.** Computation time is an important aspect of any com-
 421 parison of emulators especially on large data sets where some methods are simply not feasible.
 422 All of our models were built on a Los Alamos compute cluster 1.5TB node with 96 cores, 2
 423 Xeon Platinum 8260 CPUs @ 2.40GHz, and 192GB of Dynamic RAM with the exception of
 424 RobustGaSP at 500 training storms and SEPIA at 1000 training storms, which were run on
 425 a similar but 3.4GHz node due to limits on clock run-time.

426 As expected, the baseline linear model is extremely fast and scales well but as shown above
 427 performs poorly. We can see that BASS remains relatively fast and scales well over the range
 428 of training set sizes, requiring a modest 1 minute of computation time to fit the full training

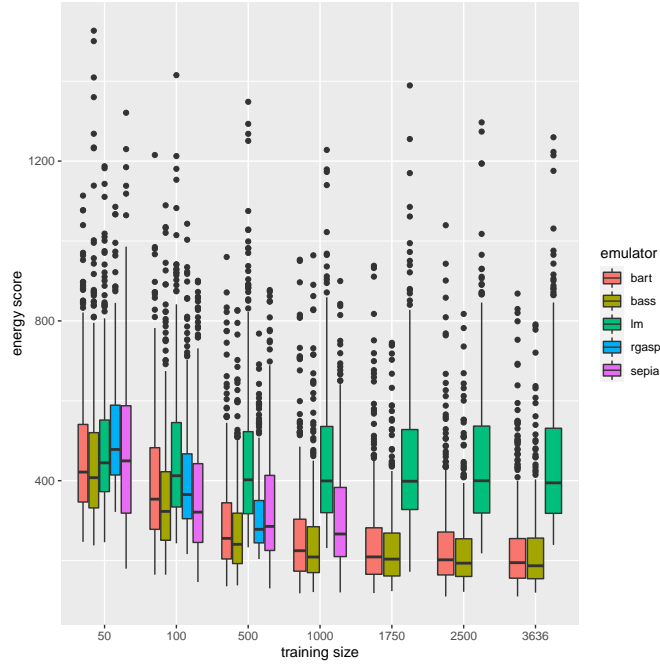


Figure 4: Energy score by training set size

429 set of 3,636 storms. BART scales similarly, requiring about 5.5 minutes for the full training
 430 set.

431 SEPIA and RobustGaSP scale relatively poorly. Both methods make use of Gaussian
 432 process which is inherently $O(n^3)$ scaling, so these methods quickly become infeasible. Ro-
 433 bustGaSP is the slowest of the four emulators, perhaps not surprising given the scope of the
 434 optimization problem it is addressing, on the native response space.

435 Parallel MCMC chain approaches may be able reduce execution time for SEPIA by a
 436 fixed factor admitting somewhat larger problems, but will not change the inherent scaling.
 437 Fortunately for RobustGaSP and SEPIA, in this application we showed that 3,636 training
 438 storms is not necessary to achieve near optimal predictive performance. We have seen that
 439 RMSE for surge height, flood area, and flood volume all reach best performance with around
 440 1000 training storms.

441 **4.4. Application Specific Metrics.** We also considered application-specific flooding and
 442 risk analysis related metrics which can be found in the supplementary materials. Specifically,
 443 we looked at predictions for the area and volume of catastrophic flooding, where area is defined
 444 as the number of land cells with greater than four feet of flood water, and volume is defined
 445 as the total water depth summed over all catastrophically flooded locations. We did not
 446 find our results to add a significant amount of information regarding the emulator methods

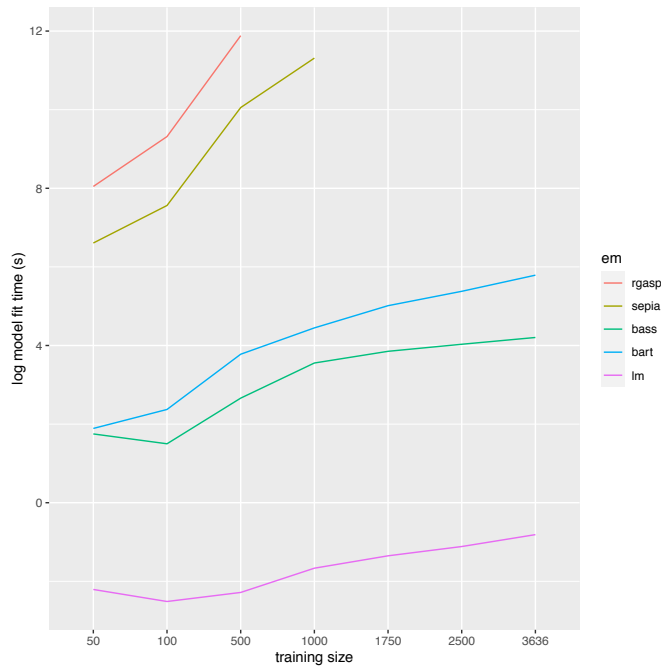


Figure 5: Model fit time

447 directly to our already rich comparison. To illustrate our point, we provide one example here
 448 in Figure 6 which shows the log RMSE for flood volume predictions. This shows little contrast
 449 to the information in Figure 2a. Additional results are available to the interested reader in
 450 the supplement.

451 There is also a description in the supplementary materials of an asymmetric loss function
 452 that we created to penalize emulators more heavily for under-prediction. This is of interest
 453 as a tunable metric that can express risk-aversion of decision-makers, especially surrounding
 454 the four foot threshold that results in power station damage. We applied this loss function to
 455 mean predictions and again found the results have no significant difference for the purposes
 456 of comparative evaluation, when compared to RMSE.

457 **5. Variable Importance.** Variable importance for computer models (often referred to as
 458 sensitivity analysis) consists of determining which inputs have the greatest (least) effect on
 459 the response. Validated emulators are useful for sensitivity analysis and variable importance
 460 calculations, as these operations typically require extensive evaluation of the response. Global
 461 sensitivity analysis consists of quantifying the percentage of the variability in the response
 462 due to each input, or combination of inputs, and is done through functional analysis of vari-
 463 ance (ANOVA) (Gu, 2018). More specifically, practitioners often use Sobol indices computed
 464 using draws from the emulator posterior predictive distribution (Sobol, 2001). An additional,

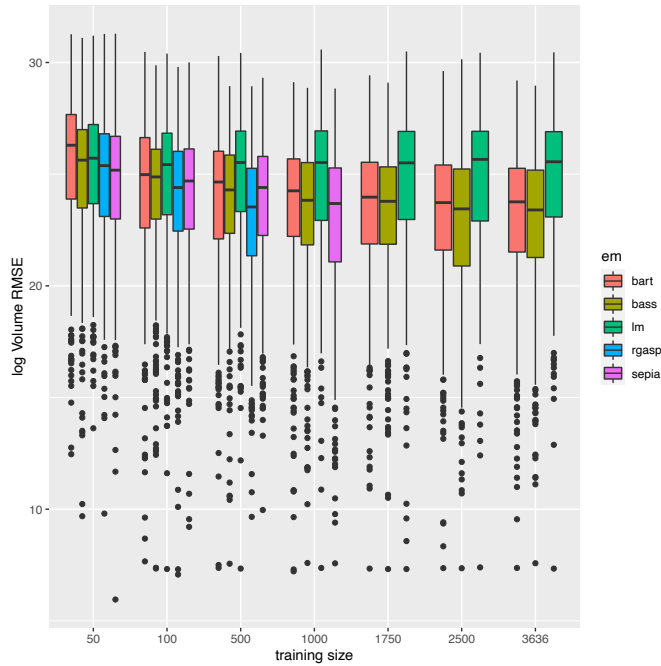


Figure 6: Flood volume log RMSE

465 very desirable property of Sobol indices is that different uncertainty distributions on the
 466 model inputs can be considered, and sensitivities can be compared across these distributional
 467 assumptions. This is very applicable to our case study as hurricane impacts are location
 468 specific, and there is no broad consensus on their spatial distributions (and the associated
 469 distributions in parameters).

470 SEPIA and BASS have built in functionality to compute Sobol indices, BART and Robust-
 471 GaSP do not. Methods for computing Sobol indices have been generalized in the R package
 472 “sensitivity” (Iooss and Pujol, 2021), so in principle sensitivity indices is available through
 473 extensions of the packages. However, the Sobol analysis requires many predictions from the
 474 emulator at various input settings, compounded by distribution samples, which would entail
 475 considerable computation. For this reason, we will instead compare the variable importance
 476 metrics that RobustGaSP and BART provide natively, rather than using those emulators to
 477 obtain Monte Carlo based Sobol indices.

478 The variable importance measures significantly differ in their implications and presenta-
 479 tion. This section is not a direct comparison of like quantities as above, but rather a presen-
 480 tation and qualitative comparison of the different information available from the methods to
 481 the user.

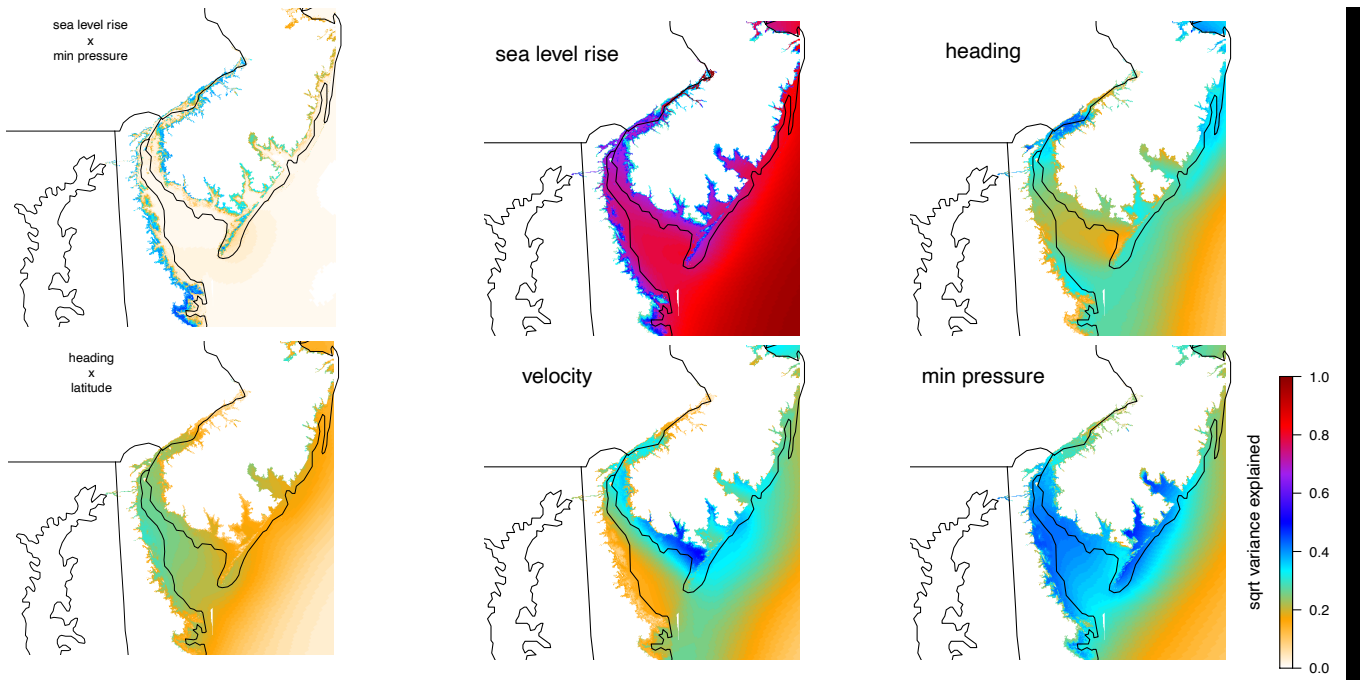


Figure 7: Bass Sobol indices, selected main and interaction effects.

482 **5.1. BASS.** BASS includes a closed-form technique for obtaining Sobol indices, facilitated
 483 by the underlying model form. The right four plots in Figure 7 show main effect Sobol indices
 484 colored by the square root of the explained variance. We can see that sea level rise explains
 485 most of the variation in the emulator and that velocity is most important at the northern
 486 opening to the bay, with a significant effect all along the northern coast.

487 We can also get sensitivity indices for interaction effects, shown in the left two plots
 488 of Figure 7. We see that interactions between sea level rise and minimum pressure play
 489 an important role in the furthest inland flooding. Our goal here is not to analyze these
 490 sensitivities, but rather to demonstrate the information provided by the Sobol decomposition.
 491 These results were generated using simple uniform priors over the input parameters.

492 **5.2. SEPIA.** SEPIA also has built in functionality for computing Sobol indices which
 493 provides sensitivities for the original response, not just the basis coefficients. Unfortunately,
 494 we found data of this size infeasible in the current implementation.

495 **5.3. BART.** BART offers a unique form of variable importance (and hence, sensitivity
 496 analysis) by keeping track of the number of times each input variable is included in the
 497 regression trees. For every posterior predictive sample, we calculate the percentage of trees
 498 containing each input variable. This gives a distribution of percentages over posterior draws.

499 The drawback is that information is only available for individual models corresponding to a
 500 single basis coefficient and we cannot simply aggregate over components to get sensitivity for
 501 the original response.

502 Figure 8a shows these distributions for the third PC and we notice that heading (theta),
 503 velocity (v), and latitude (lat) appear to be the most important inputs. This plot is more
 504 informative when combined with a visualization of the principal component as seen in Figure
 505 8b. Now we can see that these inputs explain variability mostly near the northern coast
 506 between 39 and 40 degrees latitude. Combining information from these figures gives us an
 507 idea of the locations in space where certain inputs are having an important effect. We show
 508 PC3 rather than another PC simply because it shows interesting structure and provides a
 509 good example of the results that are available from BART.

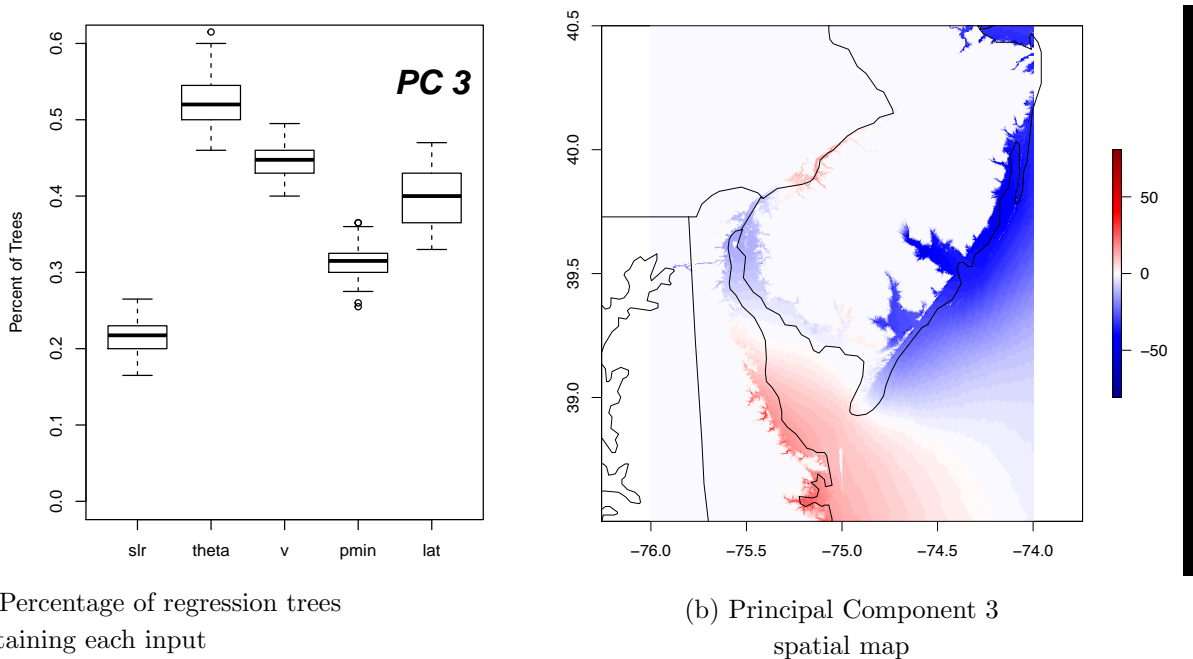


Figure 8: BART variable importance

510 **5.4. RobustGaSP.** RobustGaSP determines if an input is believed to be inert, or contrib-
 511 utes little to response variability. Inertness is decided through the estimated range parameters
 512 $\hat{\gamma}$. This is really more of a variable selection technique introduced in (Linkletter et al., 2006),
 513 but can be considered a form of variable importance or sensitivity analysis. If γ_l is inert,
 514 $\hat{\gamma}_l \rightarrow \infty$ and has little effect response variability (Gu, 2018). The JR prior we described in
 515 Section 3.4 is required for this to work. The key is that this prior, unlike the reference prior,

516 makes sure the marginal posterior for $\gamma > 0$ even if some $\hat{\gamma}_l \rightarrow \infty$. To decide whether a $\hat{\gamma}_l$ is
 517 sufficiently large to consider the associated input inert, we consider the normalized inverse

$$518 \quad (5.1) \quad \hat{P}_l = \frac{C_l \hat{\beta}_l}{\sum_{i=1}^{p_x} C_i \hat{\beta}_i}$$

519 where $\hat{\beta}_l = 1/\hat{\gamma}_l$ and C_l is a normalization constant to account for the different scales of the
 520 inputs (Gu, 2018). We can then set a threshold (default of 0.1) below which an input is
 521 determined to be inert. Table 1 shows the results for our RobustGaSP model trained on 500
 522 storms. We see that none of the inputs are found to be inert.

Table 1: Estimated normalized inverse range parameters

sea level rise	heading	velocity	min pressure	latitude
0.58	2.50	1.15	0.34	0.43

523 Albeit far less informative from a sensitivity analysis point of view than a Sobol decom-
 524 position, this is valuable information which comes for free as a byproduct of the model fit.

525 **6. Discussion.** Computer model emulation is most beneficial when applied to a simulator
 526 that is expensive to run. The SLOSH simulator is expensive enough to require emulation for
 527 analysis, but is not overly expensive; SLOSH’s relative speed is what allowed us to generate a
 528 generous ensemble of 4000 runs making a training set size study possible. The insight gained
 529 from this study can provide guidance for studies with more complex storm surge simulators
 530 like ADCIRC (Luettich and Westerink, 2015), which incorporates more physics, as well as
 531 modeling at greater resolution. As a final note about SLOSH, it was created by the National
 532 Weather Service and has thus proven to be the simulator of choice for analysis by government
 533 agencies. Given that SLOSH is so widely used, this comparison may be interesting to a wide
 534 audience of not only statisticians, but applied scientists exploring uncertainty quantification
 535 methods.

536 Figures 2a and 2b indicate that, for our case study, GP based models produced the most
 537 accurate mean predictions. This however comes at a significant computational cost as seen
 538 in Figure 5. Additionally we see evidence that our GP based models do not perform as well
 539 as BASS in terms of uncertainty quantification in Figures 3a, 3b. Therefore, we recommend
 540 SEPIA or RobustGaSP when the size of the ensemble is relatively small with correspondingly
 541 tractable computational time, and when uncertainty quantification is not the over-riding em-
 542 phasis. For most users, efficiency is likely to be very important and for these users we recom-
 543 mend BASS. BASS tends to outperform BART in our predictive metrics such as energy score,
 544 coverage probability, and interval score and it is relatively computationally tractable. Addi-
 545 tionally BASS supplies intuitive variable importance analysis through Sobol indices, relevant

546 for this application.

547 In future work we would like to confront some of the questions and limitations that arose
548 during this study. One of which is the outliers seen in all scores. It is clear that some storms are
549 performing very poorly for our predictive metrics, and although we examined some of these,
550 it is not clear whether or how these are systematic in the emulation application. An extension
551 of this work could examine whether these storms have particular features, for example a
552 particular region of the parameter space, and if outliers are consistent across methods. Another
553 limitation that comes with data of this size is the storing of large matrices, which led us
554 to use a relatively small number of posterior predictive samples. We would like to further
555 investigate whether each model has sufficiently converged. For SEPIA, BASS, and BART this
556 means analysis and diagnostics of the MCMC performance to ensure samples represent the
557 model posterior distributions, and for RobustGaSP running the optimization with a number
558 of different parameter initializations to ensure that we have not converged to a local mode. As
559 discussed, these analyses come with heavy computational burden and time that would likely
560 not be available in a typical applied analysis. Finally, in Section 4 we discussed the possibility
561 of reducing the area of particular interest to the application context of power grid impacts,
562 which would admit an effectively larger analysis within computational limitations.

563 **References.**

- 564 Chipman, H. A., George, E. I., and McCulloch, R. E. (2010). Bart: Bayesian additive regres-
565 sion trees. *The Annals of Applied Statistics*, 4(1):266–298.
- 566 Conti, S. and O’Hagan, A. (2010). Bayesian emulation of complex multi-output and dynamic
567 computer models. *Journal of Statistical Planning and Inference*, 140(3):640 – 651.
- 568 Erickson, C. B., Ankenman, B. E., and Sanchez, S. M. (2018). Comparison of Gaussian
569 process modeling software. *European Journal of Operational Research*, 266(1):179–192.
- 570 Francom, D. and Sansó, B. (2020). BASS: An R package for fitting and performing sensitivity
571 analysis of Bayesian adaptive spline surfaces. *Journal of Statistical Software*, 94(8):1–36.
- 572 Francom, D., Sansó, B., and Kupresanin, A. (2020). Landmark-based emulation for models
573 with misaligned functional response. Technical report, Technical Report UCSC-SOE-19-09,
574 University of California Santa Cruz.
- 575 Francom, D., Sansó, B., Bulaevskaya, V., Lucas, D., and Simpson, M. (2019). Inferring
576 atmospheric release characteristics in a large computer experiment using Bayesian adaptive
577 splines. *Journal of the American Statistical Association*, 114(528):1450–1465.
- 578 Friedman, J. (1991). Multivariate adaptive regression splines (with discussion). *Annals of*
579 *Statistics*, 19:1–67.
- 580 Gattiker, J., Higdon, D., and Williams, B. (2020a). lanl/gpmsa. <https://github.com/lanl/>
581 **GPMSA**.
- 582 Gattiker, J., Klein, N., Hutchings, G., and Lawrence, E. (2020b). lanl/sepia: v1.1. [https:](https://doi.org/10.5281/zenodo.4048801)
583 [//doi.org/10.5281/zenodo.4048801](https://doi.org/10.5281/zenodo.4048801).

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

586 Gramacy, R. B. and Apley, D. W. (2015). Local Gaussian process approximation for large
587 computer experiments. *Journal of Computational and Graphical Statistics*, 24(2):561–578.

588 Gramacy, R. B. and Lee, H. K. H. (2008). Bayesian treed Gaussian process models with
589 an application to computer modeling. *Journal of the American Statistical Association*,
590 103(483):1119–1130.

591 Gu, M. (2018). Jointly robust prior for Gaussian stochastic process in emulation, calibration
592 and variable selection.

593 Gu, M. and Berger, J. O. (2016). Parallel partial Gaussian process emulation for computer
594 models with massive output. *Ann. Appl. Stat.*, 10(3):1317–1347.

595 Gu, M., Wang, X., and Berger, J. O. (2017). Robust Gaussian stochastic process emulation.
596 Heaton, M., Datta, A., Finley, A., Furrer, R., Guinness, J., Guhaniyogi, R., Gerber, F., Gra-
597 macy, R., Hammerling, D., Katzfuss, M., Lindgren, F., Nychka, D., Sun, F., and Zammit-
598 Mangion, A. (2018). A case study competition among methods for analyzing large spatial
599 data. *Journal of Agricultural, Biological and Environmental Statistics*, 24.

600 Higdon, D., Gattiker, J., Williams, B., and Rightley, M. (2008). Computer model calibration
601 using high-dimensional output. *Journal of the American Statistical Association*, 103:570–
602 583.

603 Iooss, Sebastien Da Veiga, A. J. and Pujol, G. (2021). *sensitivity: Global Sensitivity Analysis*
604 *of Model Outputs*. R package version 1.24.0.

605 Jelesnianski, C., Chen, J., and Shaffer, W. (1992). Slosh: Sea, lake, and overland surges from
606 hurricanes.

607 Katzfuss, M. and Guinness, J. (2021). A general framework for vecchia approximations of
608 Gaussian processes. *Statistical Science*, 36(1).

609 Lee, L., Carslaw, K., Pringle, K., and Mann, G. (2012). Mapping the uncertainty in global
610 ccn using emulation. *ATMOSPHERIC CHEMISTRY AND PHYSICS*, 12:9739–9751.

611 Lee, L. A., Carslaw, K. S., Pringle, K. J., Mann, G. W., and Spracklen, D. V. (2011). Em-
612 ulation of a complex global aerosol model to quantify sensitivity to uncertain parameters.
613 *Atmospheric Chemistry and Physics*, 11(23):12253–12273.

614 Linkletter, C., Bingham, D., Hengartner, N., Higdon, D., and Ye, K. Q. (2006). Variable
615 selection for Gaussian process models in computer experiments. *Technometrics*, 48:478 –
616 490.

617 Luettich, R. and Westerink, J. (2015). ADCIRC. <http://www.adcirc.org>. [Accessed: May
618 2021].

619 Muniain, P. and Ziel, F. (2020). Probabilistic forecasting in day-ahead electricity markets:
620 Simulating peak and off-peak prices. *International Journal of Forecasting*, 36(4):1193–1210.

621 Möller, A., Lenkoski, A., and Thorarinsdottir, T. L. (2013). Multivariate probabilistic fore-

622 casting using ensemble Bayesian model averaging and copulas. *Quarterly Journal of the*
623 *Royal Meteorological Society*, 139(673):982–991.

624 R Core Team (2020). *R: A Language and Environment for Statistical Computing*. R Founda-
625 tion for Statistical Computing, Vienna, Austria.

626 Ramsay, J. and Silverman, B. W. (1997). *Functional Data Analysis*. Springer.

627 Sacks, J., Welch, W. J., Mitchell, T. J., and Wynn, H. P. (1989). Design and Analysis of
628 Computer Experiments. *Statistical Science*, 4(4):409 – 423.

629 Salter, J. M. and Williamson, D. (2016a). A comparison of statistical emulation methodologies
630 for multi-wave calibration of environmental models. *Environmetrics*, 27(8):507–523.

631 Salter, J. M. and Williamson, D. (2016b). A comparison of statistical emulation methodologies
632 for multi-wave calibration of environmental models. *Environmetrics*, 27(8):507–523.

633 Sobol, I. (2001). Global sensitivity indices for nonlinear mathematical models and their monte
634 carlo estimates. *Mathematics and Computers in Simulation*, 55(1):271–280. The Second
635 IMACS Seminar on Monte Carlo Methods.

636 Sparapani, R., Spanbauer, C., and McCulloch, R. (2021). Nonparametric machine learning
637 and efficient computation with Bayesian additive regression trees: The BART R package.
638 *Journal of Statistical Software*, 97(1):1–66.

639 Zhang, T., Geng, G., Liu, Y., and Chang, H. H. (2020). Application of Bayesian additive
640 regression trees for estimating daily concentrations of pm2.5 components. *Atmosphere*,
641 11(11).