

1 Nonparametric Tree-based Predictive Modeling of  
2 Storm Outages on an Electric Distribution Network

3 October 2015

4 **Abstract**

5 This paper compares two nonparametric tree-based models, quantile regression  
6 forests (QRF) and Bayesian additive regression trees (BART), for predicting storm  
7 outages on an electric distribution network in Connecticut, USA. We evaluated point  
8 estimates and prediction intervals of outage predictions for both models using high-  
9 resolution weather, infrastructure, and land use data for 89 storm events (including  
10 hurricanes, blizzards and thunderstorms). We found that QRF produced better results  
11 for high spatial resolutions, while BART predictions aggregated to coarser resolutions  
12 more effectively, which would allow for a utility to make better decisions about allo-  
13 cating pre-storm resources. We also found that the predictive accuracy was dependent  
14 on the season (e.g. tree-leaf condition, storm characteristics), and that the predictions  
15 were most accurate for winter storms. Given the comparable performance characteris-  
16 tics, we suggest that BART and QRF be implemented together to show the complete  
17 picture of a storm’s potential impact on the electric distribution network.

18 **KEY WORDS:** Weather hazards; electric distribution network; quantile regression  
19 forests; Bayesian additive regression trees; critical infrastructure outage modeling.

# 1 INTRODUCTION

Severe weather is among the major causes of damage to electric distribution networks and resultant power outages in the United States<sup>(1)</sup>. In addition to hurricanes, for which significant research has been focused, other more frequent weather systems (e.g. thunderstorms and frontal systems) have caused power outages in Connecticut lasting from several hours up to several days. Accurate prediction of the number of outages associated with storms would allow utility companies to restore power faster by allocating resources more efficiently. Furthermore, to effectively use these outage predictions in decision-making, models must exhibit acceptable accuracy in the spatial distribution of estimated outages and characterization of the prediction uncertainty.

A wide range of models have been employed in hurricane outage modeling, beginning with parametric statistical models. Generalized linear models (GLMs) were utilized by Liu *et al.*<sup>(2)</sup> using negative binomial regression with binary index variables representing storm similarity characteristics. Guikema *et al.*<sup>(3)</sup> explored the effects of tree trimming on hurricane outages with a GLM and a generalized linear mixed model (GLMM). Liu *et al.*<sup>(4)</sup> attempted the use of spatial GLMM for better inference on variables, but did not achieve improved prediction accuracies using random effects or spatial correlation modeling. Han *et al.*<sup>(5)</sup> suggested using more informative descriptive variables with GLM and performed principal components analysis (PCA) as a treatment to transform correlated variables and obtain stable parameter estimates.

Nonparametric models for hurricane outage prediction gained popularity shortly thereafter. Guikema *et al.*<sup>(6)</sup> compared multiple models including classification and regression trees (CART), generalized additive models (GAM), Bayesian additive regression trees (BART) and GLM, and discussed the advantages of nonparametric models over parametric models for outage prediction for hurricanes. Nateghi *et al.*<sup>(7)</sup> expanded the topic to outage duration modeling and concluded that multivariate adaptive regression splines (MARS) and BART had better results than traditional survival analysis models and CART, and

47 that BART produced the lowest prediction error. Guikema *et al.*<sup>(8)</sup> proposed a two-stage  
48 model using classification trees and logistic regression to deal with zero-inflation and GAM  
49 for over-dispersion, which helped balance the statistical assumption and prediction. Re-  
50 cently, Nateghi *et al.*<sup>(9)</sup> highlighted the modifiable areal unit problem (MAUP) and com-  
51 pared predictions from random forests (RF) and BART, concluding that RF benefited from  
52 its distribution-free setting and performed the best in outage duration prediction. Among  
53 all nonparametric models, tree-based models, and especially multiple trees or forests models,  
54 have been used widely and have been generally preferred in modeling hurricane outages.

55 Aiming to assist the largest utility company in Connecticut, Eversource Energy, in pre-  
56 storm decision making, we investigated two models and compared their predictions of spatial  
57 outage patterns and their ability to perform statistical inference. This study builds on our  
58 previous research that investigated the use of different model forcing data and methods for  
59 predicting power outages in Connecticut (Wanik *et al.*<sup>(10)</sup>). Prediction intervals of model  
60 estimates are as important for risk management as point estimations of storm outages; a point  
61 estimate only provides a single value at each location to describe the predicted storm outages,  
62 while prediction intervals provide a characterization of the uncertainty associated with the  
63 prediction. The lack of uncertainty characterization can affect the complex socioeconomic  
64 aspects of emergency response. In this paper, we compare two nonparametric tree-based  
65 models for the prediction of storm outages on the electric distribution network of Eversource  
66 Energy. In keeping with the most recent research in hurricane outage prediction, we selected  
67 quantile regression forests (QRF) and Bayesian additive regression trees (BART) as our  
68 candidate models because they were capable of both point estimation and prediction interval  
69 construction. BART was shown by Guikema *et al.*<sup>(6)</sup> to be the most accurate of the different  
70 hurricane outage prediction models evaluated in that study. QRF is derived from the random  
71 forest model, which has been demonstrated by Nateghi *et al.*<sup>(9)</sup> to provide robust spatial data  
72 aggregation and better power outage duration estimates than BART in terms of prediction  
73 error.

74 We seek to address the following three questions about these two models: 1) How accurate  
75 are these models in providing the point estimates (single predicted value per storm) of outages  
76 for storms of varying severity?; 2) How efficient are these models in evaluating the prediction  
77 uncertainty (i.e. the prediction interval)?; 3) Are these models able to predict outages for  
78 different spatial resolutions via aggregation?

## 79 **2 STUDY AREA AND DATA DESCRIPTION**

80 The analysis was performed on a dataset of 89 storms of multiple temporal and spatial  
81 scales (i.e. deep and shallow convective events, hurricanes, blizzards and thunderstorms)  
82 that occurred during a ten-year period (2005-2014). We selected the explanatory variables  
83 based on their potential contribution to outages on the overhead lines when interactions of  
84 overhead lines and vegetation occur. All data including distribution system infrastructure,  
85 and land cover information (Table 1) were processed on a high-resolution gridded domain  
86 (grid spacing: 2x2 km<sup>2</sup>) to represent the average conditions in the corresponding Numerical  
87 Weather Prediction (NWP) model grid spacing. Further, a seasonal categorization variable  
88 was created for each of the 89 storms (Table I) to represent the actual tree-leaf conditions  
89 (e.g. leaf-on, leaf-off or transition) at the time of each storm. The study area was the  
90 Connecticut service territory of Eversource Energy (Figure 1), which spans 149 towns in  
91 Connecticut and is organized into four divisions (central, west, east and south). The outage  
92 predictions were analyzed over the corresponding NWP model grid cells, and subsequently  
93 spatially aggregated into coarser resolutions (town, division and territory) to investigate the  
94 effects of multiple scales.

### 95 **2.1 Storm Outage Data**

96 An outage is defined as a location where a two-man restoration crew needs to be sent for  
97 manual intervention to restore power. Storm outage records were acquired from the utility's

98 outage management system (OMS) and to improve data quality, duplicate outage records  
99 and records with cause codes irrelevant to weather (e.g. vandalism or vehicular outage)  
100 were deleted from the data. Outages were recorded at the location of the nearest upstream  
101 isolating device (i.e. fuses, reclosers, switches, transformers) from where the damage on the  
102 overhead line occurred, which may be different from where the actual outage occurred. We  
103 made no differentiation of the different outage types to the overhead lines (i.e. a tree leaning  
104 on a conductor, a malfunctioning isolating device, or a snapped pole, etc.) because such  
105 data were not available to us.

## 106 **2.2 Weather Simulation and Verification**

107 The Weather Research and Forecasting Model (WRF; Skamarock *et al.*<sup>(11)</sup>) was devised  
108 to simulate the 89 storm events used in our study. The WRF model simulations were  
109 initialized and constrained at the model boundaries using NCEP Global Forecast System  
110 (GFS) analysis fields<sup>(12)</sup>. The NWP model is set up with three nested domains with 18, 6  
111 and 2-km of increasing grid spacing (Figure 2). The simulated meteorological variables were  
112 summarized into maximum and mean values (Table I). The wind-related variables in the  
113 NWP model included wind at 10m, gust winds, and wind stress. The precipitation-related  
114 variables comprised of total accumulated precipitation, the precipitation rate, and snow  
115 water equivalent (SWE). The mean values of the selected meteorological variables represent  
116 the lasting impact of a storm. The means were calculated over the 4-hour window defined by  
117 the simulated wind speed, to which hereafter we refer to as the sustained period of the storm.  
118 This period is defined by the highest averaged value in the 4-hour running window across  
119 the NWP simulation length. The wind-based sustained period was then used to calculate  
120 the mean of the other meteorological variables. The maximum values of the meteorological  
121 variables represent the peak severity that occurred during the storm; they correspond to  
122 the nominal variable value at the time of highest simulated wind speed. Complementing  
123 the mean and maximum variables, the duration of winds and gusts above defined thresholds

124 (9m/s for wind, 18m/s and 27m/s for gust) were used to relate the duration of damaging  
125 winds to outages (Table I).

126 The NWP model simulations were verified by comparing the sustained wind speed (pre-  
127 viously defined) for three major events (Hurricane Irene in August 2011; Hurricane Sandy  
128 in October 2012; Blizzard Nemo in February 2013) against METAR observations (airport  
129 meteorological station data) provided by the National Centers for Environmental Prediction  
130 (NCEP) ADP Global Upper Air and Surface Weather Observations <sup>(13)</sup>. Though not shown  
131 here, the NWP model simulations showed acceptable agreement with the airport station  
132 data (e.g. low mean bias, and high correlation between actual and simulated sustained wind  
133 speed). The reader may refer to Wanik *et al.* <sup>(10)</sup> for more details on this verification exercise.

### 134 **2.3 Seasonal Categorization**

135 Storms affecting the distribution network can have a wide range of weather attributes (e.g.  
136 heavy snow or rain) that interact with overlying vegetation, and can have differing impact  
137 on the outage magnitude and frequency depending on the tree-leaf condition. For exam-  
138 ple, high winds usually have a greater impact on trees with leaves due to increased wind  
139 loading <sup>(14), (15)</sup>. To capture this dynamic, we grouped our data by season (Table I), which  
140 resulted in separate fits for each of the three different seasonal categories. Of the 89 storm  
141 events in our dataset, there were 38 storms and 1 hurricane (Irene) during the summer (leaves  
142 on) months (June to September); there were 24 storms and 1 hurricane (Sandy) during the  
143 spring and fall (transition) months (October, November, April and May); and there were 25  
144 storms during the winter (leaves off) months (December to March).

### 145 **2.4 Infrastructure and Land Use**

146 The same infrastructure and land use data from Wanik *et al.* <sup>(10)</sup> study was used in this  
147 paper. The sum of isolating devices (e.g. sum of fuses, reclosers, switches, transformers)  
148 in each 2-km grid cell was an important predictor in our models, which we attribute to

149 the outage recording methodology (Section 2.1); if only one outage can be recorded at an  
150 isolating device, a grid cell with more isolating devices has more chances to record more  
151 outages than a grid cell with less isolating devices. Given that outages were recorded at the  
152 nearest isolating device and not the actual outage location, the different types of isolating  
153 devices were summed up into a single variable (“sumAssets”) instead of modeling outages  
154 by isolating device type. This variable sets the upper limit on the number of outages that  
155 could occur in a 2-km grid cell.

156 Accurate tree-specific data (i.e. height, species, and health) around overhead lines are  
157 difficult to acquire, so we used land cover data aggregated around the overhead lines as a  
158 surrogate for the actual tree data. This aggregation differs with previous research (Quir-  
159 ing *et al.* <sup>(16)</sup>) that used the percentage of all land cover types in a grid cell, regardless of  
160 whether or not certain land cover types in that grid cell interacted with the overhead lines  
161 (i.e. a waterbody that is in the grid cell but is not close enough to the overhead lines to  
162 cause influence). Land cover data (30m resolution) were obtained from the University of  
163 Connecticut Center for Land Use Education and Research (CLEAR) <sup>(17)</sup> and were used to  
164 generate percentages of land use per grid cell. Details about the calculation of land use are  
165 available in Wanik *et al.* <sup>(10)</sup>.

### 166 **3 METHODOLOGY**

167 As two systematically different examples of nonparametric tree-based models, QRF and  
168 BART utilize different assumptions and techniques for their application. We briefly introduce  
169 the benefits and known issues of these two models, followed by measurements and methods  
170 for analysis and comparison of the models.

### 171 3.1 Quantile Regression Forests

172 Based on the well-known random forests algorithm by Breiman <sup>(18)</sup>, Meinshausen <sup>(19)</sup> created  
173 the quantile regression forests (QRF) with the idea of quantile regression from econometrics.  
174 Similar to the weighted average of all the trees for predicted expected value of response,  
175 QRF utilizes the same weights to calculate the empirical distribution function:

$$\hat{F}(y|X = x) = \sum_{i=1}^n w_i(x) 1_{\{Y_i \leq y\}}. \quad (1)$$

176 The algorithm of QRF can be summarized as:

- 177 1. Grow  $k$  trees  $T_t$ ,  $t = 1, \dots, k$ , as in random forests. However, for every terminal node of  
178 every tree, take note of all observations, not just their average.
- 179 2. For a given  $X = x$ , drop  $x$  down all trees. Compute the weight  $w_i(x, T_t)$  of observation  
180  $i = 1, \dots, n$  for every tree. Compute weight  $w_i(x)$  for every observation  $i = 1, \dots, n$  as an  
181 average over  $w_i(x, T_t)$ ,  $t = 1, \dots, k$ . These weight calculations are the same as in random  
182 forests.
- 183 3. Compute the estimate of the distribution function as in (Equation 1), using the weights  
184 from Step 2.

185 Given the estimated empirical distribution, quantiles and percentiles are readily available.  
186 In this case, mean from the estimated distribution is a natural point estimate, which is  
187 exactly the same as random forests. Meinshausen <sup>(20)</sup> has provided an R package called  
188 “quantregForest”, with dependency on the “randomForest” package by Liaw *et al.* <sup>(21)</sup>. Our  
189 analysis was based on a slightly modified version of “quantregForest” package providing RF  
190 prediction as well as desired quantiles.

191 QRF has already been used in several aspects of natural phenomena, but have not been  
192 implemented in storm outage prediction directly. Juban *et al.* <sup>(22)</sup> tested QRF in approximat-  
193 ing the kernel density of short-term wind power, indicating comparatively wide prediction

194 intervals by QRF. Francke *et al.* <sup>(23)</sup> addressed the better performance of QRF over GLMs  
 195 in flood-based analysis of high-magnitude sediment transport. Zimmermann *et al.* <sup>(24)</sup> also  
 196 took advantage of QRF to study erosion in rainforests.

197 Suppose we need to generate a prediction interval for a storm event by aggregation.  
 198 Accurate empirical distributions are highly preferred, but demand a large number of obser-  
 199 vations. For example, to generate percentiles (actually 101 quantiles including maximum and  
 200 minimum), we prefer more than 100 observations in each terminal node. However, enforcing  
 201 too many observations in terminal nodes could introduce bias if the sample size is not large  
 202 enough. In this study, we compared QRF to BART in order to get a deeper understanding  
 203 about the importance of prediction intervals in characterizing model performance.

## 204 3.2 Bayesian Additive Regression Trees

205 Bayesian additive regression trees model (BART), introduced by Chipman, George and  
 206 McCulloch <sup>(25, 26)</sup>, is a high performance derivation of Bayesian classification and regression  
 207 trees model (CART). It takes advantage of a backfitting MCMC algorithm <sup>(27)</sup> in generating  
 208 the posterior sample of CART. Instead of a single regression tree (the mode of posterior  
 209 tree sample), a sum of regression trees is utilized to estimate the response under normal  
 210 assumption:

$$Y = \underbrace{\sum_{j=1}^m g(x; T_j, M_j)}_{\text{Mean Model}} + \underbrace{\epsilon, \epsilon \sim N(0, \sigma^2)}_{\text{Variance Model}}. \quad (2)$$

211 Here  $T_j$  stands for the  $j^{\text{th}}$  regression tree;  $M_j$  stands for the  $j^{\text{th}}$  set of terminal nodes;  $m$   
 212 stands for total number of trees. The prior for probability of splitting node  $\eta$  (depth= $d_\eta$ ),  
 213 which is also presented by Chipman, George and McCulloch <sup>(28)</sup>:

$$p_{\text{split}} = \alpha(1 + d_\eta)^{-\beta}, 0 \leq \alpha \leq 1 \text{ and } \beta \geq 0. \quad (3)$$

214 Similar to Friedman's gradient boosting <sup>(29)</sup>, each terminal nodes  $\mu_{ij}$  is determined by

215  $N(0, \sigma_\mu^2)$ , where  $\sigma_\mu = 0.5/k\sqrt{m}$ . An inverse chi-square distribution is set as the prior of  $\sigma^2$ ,  
216 parameterized with  $\nu$  and  $P(\sigma < \hat{\sigma}) = q$ . All of these hyper parameters can be optimized  
217 via cross-validation.

218 Chipman *et al.*<sup>(30)</sup> provided an R package “BayesTree” based on C and Fortran, with  
219 their original work. Pratola *et al.*<sup>(31)</sup> offered a standalone C++ implementation with fast  
220 parallel computation. Kapelner and Bleich *et al.*<sup>(32)</sup> made the R package “bartMachine”  
221 based on rJava, including features like parallel cross-validation and interaction detection,  
222 which we used in this paper.

223 BART has been widely used in risk analysis and the prediction of natural hazards.  
224 Guikema *et al.*<sup>(6)</sup> conducted a comparison of multiple models for estimating the number of  
225 damaged poles during storms, and concluded that BART and an ensemble model with BART  
226 outperformed other parametric regression methods. Nateghi *et al.*<sup>(7)</sup> compared BART with  
227 traditional survival models in predicting power outage durations in Hurricane Ivan, 2004,  
228 and concluded that BART had better performance over parametric survival models. Blat-  
229 tenberger *et al.*<sup>(33)</sup> implemented BART in predicting binary response of avalanches crossing  
230 streets. They compared BART classification with linear and logistic regressions by altering  
231 the cutoff probabilities and concluded that BART excelled in predicting binary response.

232 As a well-defined Bayesian statistical model, BART naturally offers prediction intervals  
233 under its model assumptions, but the error term can be misspecified with respect to storm  
234 outage modeling. Both modeling the number of outages<sup>(6)</sup> and outage durations<sup>(9)</sup> involve  
235 errors that do not necessarily follow a normal distribution. In our study, the response  
236 variable (the number of outages) seemed to follow a Poisson distribution in grid cells and  
237 towns, while a zero-truncated normal distribution seemed to fit better in divisions and the  
238 territory (Figure 3; hurricanes are excluded for extreme values.). That is, these errors could  
239 approximately satisfy normality and homogeneity of variance in some situations, while the  
240 distribution of data aggregated with different spatial resolution can vary greatly. This issue  
241 was first discovered by Gehlke and Biehl<sup>(34)</sup> and later discussed in details by Openshaw<sup>(35)</sup>.

242 To understand the impact in our study, we would like to study the prediction intervals given  
243 by BART in more detail.

### 244 **3.3 Metrics of Model Performance Evaluation**

245 We will compare the two models using the following metrics:

246 **Mean absolute error** (MAE) is an absolute measurement of the point estimate error  
247 of  $n$  predictions, which is calculated by

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|. \quad (4)$$

248 **Root mean square error** (RMSE) measures the magnitude of error as well as its  
249 variability, which is defined as

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}. \quad (5)$$

250 MAE is less than or equal to RMSE. Closer difference between MAE and RMSE indicates  
251 smaller variability in the point estimate of error. A combination of MAE and RMSE is a  
252 common tool in model comparison with respect to point estimates.

253 **Relative error** (RE) is also known as the relative percentage error, which is computed  
254 by the following normalized average:

$$RE = \frac{\hat{y}_i - y_i}{y_i}. \quad (6)$$

255 RE is useful for diagnostics of over-prediction or under-prediction, typically offering an  
256 indication of bias.

257 **Nash-Sutcliffe efficiency** (NSE) is the generalized version of R-squared from paramet-  
258 ric regression, and is widely used in hydrology. NSE was introduced by Nash and Sutcliffe <sup>(36)</sup>

259 and summarized by Moriasi *et al.* <sup>(37)</sup>. It is calculated by the following:

$$NSE = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \stackrel{\text{unbiased}}{=} 1 - \frac{\hat{Var}(error)}{\hat{Var}(y)}. \quad (7)$$

260  $R^2$  is intuitively known as “percent variance explained”, and pseudo  $R^2$  is a built-in  
 261 statistic of both “randomForest” package and “bartMachine” package in R. However, since  
 262  $R^2$  is always measured with in-sample data, we use the name of NSE to highlight its capability  
 263 in validation. Without bias and overfitting, NSE is a powerful tool measuring predictions of  
 264 spatial variability for nonparametric models; NSE values range from negative infinity to 1.  
 265 For example, NSE for a single storm validation could be negative, while the average NSE for  
 266 the validation of all the storms remains positive, indicating that predictions for this specific  
 267 storm is worse than a mean-only model in terms of spatial variability. When calculating  
 268 NSE, we corrected the bias of total predicted value of each storm by scaling in order to focus  
 269 on the spatial variability of each storm event.

270 When  $R^2$  is positive, it is weakly increasing as  $p/n$  increases, where  $p$  stands for the  
 271 number of predictors and  $n$  stands for number of predictions. In our case, we aggregate  
 272 predicted values from high resolution to low resolutions, which actually decreases  $n$  with  
 273 predictors fixed. When NSE is positive ( $<1$ ) for grid cells, we can expect NSE to be positive  
 274 and even closer to 1 for towns and divisions, since the pseudo  $p/n$  increases. Similarly,  
 275 negative NSE in high resolution may result in even smaller NSE in low resolutions. Thus we  
 276 may observe a “polarization” effect after aggregation.

277 **Uncertainty ratio** (UR) is a benchmark statistic of prediction uncertainty. Denote the  
 278 prediction interval as  $(Q_{lower_i}, Q_{upper_i})$ . Similar to the UR used in Özkaynak *et al.* <sup>(38)</sup>, our  
 279 definition of UR is a generalized version for asymmetric intervals:

$$UR = \frac{\sum_{i=1}^n (Q_{upper_i} - Q_{lower_i})}{\sum_{i=1}^n y_i}. \quad (8)$$

280 While UR is computed by summing up all the ranges of prediction intervals, the formal

281 calculation of a prediction interval for each storm is done by summing up the simulated  
 282 sample of prediction for each grid cell. Larger UR indicates relatively wider ranges of the  
 283 prediction intervals (which may be less useful than narrower intervals, as they provide less  
 284 detailed information).

285 **Exceedance probability** (EP) is a measurement of the probability that actual value  
 286 will exceed the prediction interval. In this paper, we calculate the EP for each storm by

$$\hat{P}\{exceed.\} = 1 - \frac{1}{n} \sum_{i=1}^n 1_{\{Q_{lower_i} < y_i < Q_{upper_i}\}}. \quad (9)$$

287 Similar to UR, large EP is unfavorable, and it implies a large chance to have actual values  
 288 outside the prediction intervals.

289 **Coverage percentage** is the opposite of exceedance probability and is defined as

$$Coverage \% = \frac{1}{n} \sum_{i=1}^n 1_{\{Q_{lower_i} < y_i < Q_{upper_i}\}} * 100\%. \quad (10)$$

290 In contrast to exceedance probability, we pursue high coverage rate of prediction intervals  
 291 on actual values.

292 **Rank histogram:** A rank histogram provides a diagnosis of bias and dispersion in en-  
 293 semble predictions (detailed interpretation of this plot can be found in Hamill <sup>(39)</sup>). Suppose  
 294 we have  $m$  predictions for each observation  $y_i$ , denoted as  $\{\hat{y}_{ij}\}_{j=1,2,\dots,m}$ , and then the “en-  
 295 semble” prediction is  $\hat{y}_i = \overline{\hat{y}_{ij}}$ . A good model implies that the response is a realization of the  
 296 prediction distribution, namely:

$$E\{P[y_{i,j-1} < y_i < y_{i,j}]\} = \frac{1}{m+1}. \quad (11)$$

297 A rank histogram could be generated by collecting all the ranks of actual value in their  
 298 prediction samples, denoted by  $\mathbf{R} = (r_1, r_2, \dots, r_{m+1})$ , where  $r_j$  is the following average over

299 all the  $i$ 's:

$$r_j = \overline{\hat{P}\{y_{i,j-1} < y_i < y_{i,j}\}}. \quad (12)$$

300 An ideal rank histogram should display a uniform distribution. A rank histogram with  
301 a convex function indicates under-dispersive predictions; a concave function implies over-  
302 dispersion. A rank histogram with larger values on the right than on the left addresses  
303 negative bias of predictions, while positively biased predictions yield large value on the left.  
304 In short, rank histogram above average means the distribution of actual values is “denser”  
305 than the distribution of predictions and vice versa.

### 306 **3.4 Methods of Model Performance Comparison**

307 Our analysis is based on dataset containing 89 storms that occurred in the Eversource  
308 (Connecticut) service territory, which resulted in 253,739 observations (89 storms with 2,851  
309 grid cells per storm) for model training and validation. Within each storm, we randomly  
310 selected two thirds of the observations ( $n = 1,989$ ) per storm for model training, and training  
311 data were grouped by season. The rest of the data ( $n = 952$  grid cells per storm) were used for  
312 model validation. Only model validation results will be presented in our model verification  
313 statistics.

314 For the QRF model, we specified a random forest of 1,000 trees (default setting for  
315 the “quantregForest” package<sup>(20)</sup>) and 200 minimal instances for each terminal node to  
316 generate percentiles. Quantile regression was introduced to obtain 101 percentiles (including  
317 minimum and maximum) and predicted empirical distribution for each grid cell. The mean of  
318 the predicted empirical distribution (the same as random forest predictions), are recorded as  
319 point estimates at the grid cell level. We calculated 80% prediction intervals for each grid cell  
320 with the 10% and 90% prediction quantiles from QRF. We then sampled from the predicted  
321 empirical distribution 10,000 times per grid cell and aggregated these prediction samples to  
322 get empirical distributions in town, division and territory resolution. The weights obtained

323 from step 2 in Section 3.1, were normalized as probabilities to draw prediction sample from  
324 training data responses. The mean and quantiles are consistent with predicted distribution,  
325 proved by Bickel *et al.* <sup>(40)</sup>. After that, sample means and 80% intervals are calculated for  
326 these 3 granularities. In the end, we combined these point estimations (means) and intervals  
327 for different seasons and generated plots of statistics, such as NSE, UR and rank histogram.

328 For the BART model, a 5-fold cross-validation indicated the following settings for the  
329 parameters: 50 trees,  $k = 2$ ,  $q = 0.99$ ,  $\nu = 3$ , while other parameters remained default. In  
330 order to reach the convergence of MCMC, we performed 10,000 burn-in iterations, which  
331 produced a momentum sample (discarded). After that, we ran another 10,000 iterations to  
332 get the prediction sample of the model, which was used for prediction and validation. Point  
333 estimations and prediction samples for each observation in the testing dataset were computed  
334 in the Bayesian way. That is, sampling from the mean model (Equation 2) posterior sample  
335 and variance model (Equation 2) posterior sample under model assumptions. In contrast  
336 to QRF, BART generated prediction samples from a well-defined distribution instead of  
337 empirical distribution. Prediction intervals were calculated for BART in a way similar to  
338 QRF: the 10% and 90% quantiles from prediction samples for grid cells, or aggregated  
339 prediction samples for resolution lower than the grid cell.

340 We plotted statistics grouped by season and varying spatial resolutions (grid cell, town,  
341 division and territory) for evaluation. By inspecting the plots, we intuitively summarized  
342 the different behaviors of QRF and BART, followed by a discussion of the observed patterns  
343 and their causes.

## 344 4 RESULTS AND DISCUSSION

345 First, we will discuss the QRF and BART performance for predicting Hurricane Irene (2011)  
346 and Hurricane Sandy (2012) outages. Then, we will evaluate the consistency and prediction  
347 intervals of QRF and BART for different types of storms in our database. This evaluation

348 is based on the statistical metrics described above, evaluated at different spatial resolutions  
349 ranging from the 2-km grid cell to town and regional averages.

## 350 **4.1 Hurricane Outage Modeling**

### 351 **4.1.1 Point Estimate Results**

352 With respect to hurricanes, Table II shows that both QRF and BART performed well in  
353 terms of point estimates, compared with a mean model (assuming uniform outages across grid  
354 cells). The mean model performed well in predicting total number of outages for Hurricane  
355 Sandy, because the randomly selected validation partition happened to capture two thirds  
356 of total outages in this case. However, the town level MAE and RMSE reveal that the  
357 mean model did not spatially predict the actual outages for Irene and Sandy. Both QRF  
358 and BART exhibited small MAE values (5.50 - 8.86 outages per town) with RMSE values  
359 close to MAE values, which indicates moderate variance of the point estimates and a lack of  
360 large residuals. For these tropical storm cases, BART showed less error than QRF in town  
361 resolution but did not exhibit an overwhelming advantage. QRF predictions were slightly  
362 more spread-out than BART as shown in the scatter plots of Figure 4, which were consistent  
363 to the error metrics in Table II. Figure 5 illustrates the similar capability of QRF and BART  
364 in explaining the spatial variability of the predictions using the town-aggregated estimates;  
365 both QRF and BART predicted that the majority of outages from Irene and Sandy would  
366 be in central and southwestern Connecticut.

### 367 **4.1.2 Prediction Interval Results**

368 Both models exhibited different characteristics in terms of their prediction intervals. Figure  
369 4 shows that QRF produced more conservative town-resolution predictions by offering longer  
370 prediction intervals and a higher coverage rate than BART; coverage values for Irene and  
371 Sandy were 84% (69%) and 87% (77%) for QRF (BART) using the 80% confidence intervals.  
372 However, BART had narrower prediction intervals and was able to the cover actual number

373 of outages for both Irene and Sandy in Table II, while QRF failed to cover Irene’s actual  
374 number of outages in its interval. We also noticed that BART produced symmetric intervals  
375 from the normal distribution, while QRF generated asymmetric intervals from the empirical  
376 distribution. Although we observed comparatively good results for BART, we were unable to  
377 conclude which model was superior in terms of predicting storm outages from analyzing these  
378 two hurricanes alone. Further investigation of the model performance for both hurricanes  
379 and the remaining less severe weather events follows next.

## 380 **4.2 Comparison of Model Performance for All Storm Events**

### 381 **4.2.1 Point Estimate Results**

382 In this section, we will examine how QRF and BART explained the magnitude and spatial  
383 variation of outages (i.e. point estimates), and also examine how both models’ prediction  
384 intervals explained the variability of predicted outages. Our analysis will highlight depen-  
385 dencies of our analysis on storm severity, season and leaf condition.

386 Figure 6 summarizes the overall fit of QRF and BART aggregated by storm events. Both  
387 models show over-prediction for low impact events ( $< 100$  outages), while QRF also shows  
388 under-prediction for medium-high impact events (between 100 and 1000 outages). QRF  
389 (BART) exhibits coverage rates of prediction interval around 28% (36% to 60%), which were  
390 below our expectation (further analysis will reveal the cause of this phenomenon). We can see  
391 variations in performance across different leaf conditions and storm severities; both models  
392 did especially well at predicting hurricanes and for minor events with leaves off (winter).

393 First we will investigate the point estimate predictions of both models. Figure 7 (a), (b)  
394 and (c) illustrate how NSE varied for different spatial resolutions (grid, town, division) as a  
395 function of magnitude (outages per storm), while Figure 8 (a), (b) and (c) illustrate the same  
396 subject vs. different leaf conditions in boxplots. We see that the majority of storm events  
397 enjoy positive NSE values in Figure 7 and the 25% quantiles of NSE were close to or above 0  
398 in Figure 8, indicating that the most of predictions were informative in predicting the spatial

399 variability of outages. In addition, both models show that the NSE for the three resolutions  
400 generally exhibited a positive correlation between the accuracy of spatial variability modeling  
401 and the magnitude of the event outages (Figures 7). As expected (Section 3.3), we observed  
402 in both Figures 7 and 8 that for those events with positive NSE in grid cell resolution, the  
403 NSE's increased as the scale resolution became coarser (hence, we have many more events  
404 with NSE close to 1 at division resolution than at the grid resolution). Conversely, NSE  
405 tended to worsen at coarser scale aggregations for events where the models exhibited negative  
406 NSE at grid resolution. This “polarization” effect was so significant in division resolution  
407 that we should be cautious in using the aggregated results in this resolution, because the  
408 predicted outages may not reflect the true spatial distribution of outages for some storms.  
409 For all three resolutions, BART yielded better (more positive) NSE than QRF.

410 Figures 7 (d) and 8 (d) show how the relative error (RE) of aggregated storm total  
411 predictions varied in event magnitude and leaf condition. In the territory resolution (on the  
412 storm event scale), there is no NSE defined and instead we are more interested in how the  
413 point estimate performed in predicting the actual magnitude (outages per storm). The RE  
414 exhibited a negative correlation with magnitude in Figure 7, which suggests that both models  
415 were accurate for the most severe events. In Figure 8, BART frequently yielded RE's with  
416 smaller variance than QRF, which is consistent to its less spread-out predictions in Figure 6.  
417 Moreover, BART had the lowest (closest to zero) RE as well as highest NSE for the leaves off  
418 season (Figure 8). We attribute BART's improved winter (leaves off) season error metrics to  
419 the similarity of the data in the grouping (e.g. only minor events (no hurricanes), the ground  
420 is more likely to be frozen). In short, we conclude that BART yielded better point estimates  
421 than QRF because it had higher NSE and lower (more close to zero) RE than QRF.

#### 422 **4.2.2 Prediction Interval Results**

423 We also examined the prediction intervals provided by BART and QRF with uncertainty  
424 ratio (UR), exceedance probability (EP) and rank histogram. Prediction intervals that are

425 very wide offer no value to decision makers, because they suggest any amount of storm outage  
426 may occur; conversely, a too-narrow prediction interval may not give a useful estimate of the  
427 extent of possible outages. These widths of intervals are captured by UR in Figure 9 and a  
428 negative correlation between UR and magnitude was observed. Although narrow prediction  
429 intervals for high-magnitude events were favorable for their high certainty, the coverage rate  
430 becomes an important issue. In contrast to the UR, the actual value of outages exceeded  
431 the interval more frequently in severe events than in moderate events according to Figure 10  
432 (a) and (b). We see a reverse trend of UR and EP vs. response magnitude. For BART, this  
433 is due to the homogeneity of variance assumption that made BART to offer intervals with  
434 similar absolute widths; for QRF, this is due to the fixed minimal number of instances in  
435 each terminal node that treated the nonzero or extreme responses the same as zero responses  
436 (more than 80% of responses as zeros at grid cell level in Figure 3 (a)). In practice, we look  
437 for prediction intervals that have acceptable small UR and stable EP which associated with  
438 the confidence level (e.g. a stable probability of 0.2 given 80% intervals in our case). This  
439 suggests that more flexible assumptions and settings are needed for BART and QRF to  
440 capture the variation in response magnitude.

441 During aggregation, UR reduced step by step from grid cell resolution to territory reso-  
442 lution (Figure 9), which is favorable. In contrast, the EP generally increased step by step  
443 via aggregation for both models (Figure 10), which led to the low interval coverage rates in  
444 Figure 6. In fact, QRF offered both low UR and low EP, implying superior performance  
445 in the highest resolution (i.e. grid cell). However, QRF's UR and EP became similar to or  
446 larger than BART's after aggregation, suggesting weakness in QRF's spatial aggregation.  
447 Since there are only four different divisions, EP can only be 0, 0.25, 0.5, 0.75 or 1 at division  
448 level in Figure 10 (c); Similarly, EP can only be 0 or 1 at storm event level in Figure 10 (d).  
449 Specifically, Figure 10 (c) and (d) address that QRF suffered more 1's of EP than BART.

450 To further elaborate the nature and issues of both models, we introduce rank histograms  
451 in Figures 11 and 12. In practice, a uniformly distributed rank histogram means the pred-

452 icated distribution (including its quantiles and intervals generated by quantiles) reflects  
453 the variability of the actual response. Overall QRF (Figure 11 (a)) did well in grid cell  
454 resolution as evidenced by the near-uniform distributed rank histogram with a moderate  
455 under-prediction issue. In comparison, BART (Figure 12 (a)) produced biased predicted  
456 distribution by assuming normal distribution on Poisson-distributed actual outages (Figure  
457 3 (a)) for grid cells. However, spatial aggregation appears to undermine the QRF prediction  
458 by accumulating biasedness (Figure 11). It is interesting to see BART (Figure 12) benefited  
459 a little from spatial aggregation. In fact, BART yielded better predictions for storm event  
460 totals, where the normal distribution becomes a better approximation of combined Poisson  
461 distributions (Figure 3). This explains why BART ended up with better interval coverage  
462 rates in Figure 6, even though QRF started from more accurate predicted distribution for  
463 grid cells. Note that biasedness also differs from location to location and aggregating lo-  
464 cations with under-estimates and locations with over-estimates could result in complicated  
465 bias which is hard to predict. In conclusion, BART produced better prediction intervals for  
466 divisions and whole territory, while QRF did better for grid cells and towns.

### 467 4.3 Discussion

468 Similar to previous works in parametric modeling like Liu *et al.* <sup>(4)</sup>, the models we utilized  
469 offer prediction intervals as well as point estimates of outages. Instead of simply identifying  
470 the potentials in quantifying prediction uncertainty, we took one more step in this study to  
471 evaluate QRF and BART with real-world data for their uncertainty measures. For hurri-  
472 canes, BART model exceeded QRF in both predicting the outage magnitude (e.g. effective  
473 prediction intervals) and spatial variation of hurricane outages (Table II). BART under-  
474 predicted Irene by 1.9% and over-predicted Sandy by 2.4%, while the best ensemble decision  
475 tree model in our previous work (Wanik *et al.* <sup>(10)</sup>) over-predicted Irene by 11% and under-  
476 predicted Sandy by 4.4%. However, caution must be exercised when directly applying BART  
477 to storm outage modeling. Nateghi *et al.* <sup>(9)</sup> illustrated the weakness of BART when com-

478 pared to random forest in survival analysis of hurricane outages where the response variable  
479 did not follow normal distribution. Most storms cause much less outages than hurricanes and  
480 even result in zero outage in some grid cells (Figure 3 (a)). Alternatively, QRF is promis-  
481 ing in generating predictions and intervals without normality. Our analysis suggests that  
482 QRF suffered minor bias (Figure 11 (a)) in dealing with zero-inflated number of outages,  
483 while BART suffered significant bias (Figure 12 (a)). Compared to previous research (e.g.  
484 Guikema *et al.* <sup>(8)</sup>), we used real zero-inflated response variable based on storm events instead  
485 of simulated data or hurricanes, to suggest proper treatments to zeros. However, we found  
486 that BART was at least as good as QRF with respect to aggregated point estimates (Figure  
487 7 and 8) and was better at generating aggregated prediction intervals (Figure 9 and 10). In  
488 short, different models could be utilized based on different interests or scales of application.

489 There are limitations for our study and results. Unlike Hurricane Ivan, studied by Nateghi  
490 *et al.* <sup>(9)</sup>, Hurricanes Irene and Sandy did not make landfall in the service territory of our  
491 study. Since these storms were not at their peak when they impacted Connecticut, our  
492 research does not necessarily reflect the “worst case scenario” for outages. We did not  
493 include ice storms in this research due to their fundamentally different characteristics with  
494 other events in our database. The categorization of leaf conditions according to seasonal  
495 periods and spatial aggregation strategy according to geographic boundaries was based upon  
496 utility’s demand to integrate the models with their emergency planning efforts.

## 497 **5 CONCLUSIONS AND FUTURE WORK**

498 This article has developed and validated outage prediction models for an electric distribu-  
499 tion network. We incorporated high-resolution weather simulations, distribution infrastruc-  
500 ture and land use for modeling storm outages using quantile regression forests (QRF) and  
501 Bayesian additive regression trees (BART). For hurricanes, BART model exceeded QRF in  
502 both predicting the outage magnitude and spatial variation of hurricanes. In our study, we

503 found that outages caused by storms were not normally distributed and followed different  
504 distributions in different spatial resolutions. Hence, QRF was better at characterizing storm  
505 outages in high resolution, but did not aggregate well (from grid to division resolution).  
506 In contrast, BART did well at aggregating predictions for the storm outage total, but did  
507 not fully characterize the distribution of storm outages at higher resolution (e.g. grid res-  
508 olution). In an operational context, utility companies might like to use maps of pre-storm  
509 outage predictions at the town resolution while also viewing a broader summary of the spa-  
510 tial variability, point estimates, and prediction intervals for the whole service territory. We  
511 suggest presenting the results from BART at coarser resolutions (e.g. division and service  
512 territory) and results from QRF for higher resolutions (e.g. grid and town) to best present  
513 the potential storm outages. Doing so will ensure that decision-makers get a complete idea  
514 of the overall severity of the event at a coarser resolution while also providing the detailed  
515 information supporting a pre-storm response at a higher resolution.

516 There are many opportunities for improvement in storm outage modeling on electric dis-  
517 tribution networks. From a methodological point of view, both models could be modified  
518 to deal with the Poisson-distributed sparse (i.e. more than 80% as zeros) response variable.  
519 For QRF, empirical distributions for extreme observations are very different from the ma-  
520 jority (mostly zeros). Varied treatments for the majority vs. real signals (i.e. outages) could  
521 help. By experimenting with the minimum number of observations required for the terminal  
522 node (according to overall magnitude of observations in the node), the accuracy of point  
523 estimates and prediction intervals may be improved. For BART, an application in general-  
524 ized linear models (GLMs) with more flexible assumptions and link functions becomes very  
525 important, because Poisson-distributed data or zero-inflated data appear frequently in high-  
526 resolution analysis. For example, Poisson, negative binomial and zero-truncated normal with  
527 heterogeneous variance could perform as priors to assist BART in storm outage modeling. In  
528 addition to the two-staged model (similar to the classification-GAM model used by Guikema  
529 *et al.* <sup>(8)</sup>), a zero-inflated BART model optimized simultaneously for both zero-inflated clas-

530 sification and zero-truncated signals could be implemented. For spatial aggregation, the  
531 ideal unbiased prediction may not be available for every location, thus getting accurate pre-  
532 dictions for another resolution based on biased results is challenging and important. There  
533 are already some advanced techniques (e.g. Reilly *et al.*<sup>(41)</sup>) to aggregate point estimates  
534 into multiple scales while eliminating bias and error by utilizing spatial patterns. Similar  
535 techniques to aggregate predicted distributions for each location could be investigated in a  
536 future study.

537 From a modeling point of view, the inconsistent performance of both models for varying  
538 season categories (tree-leaf condition) implies difficulties in predicting storm outages with  
539 leaves on trees. Future research may consider including additional effective explanatory  
540 variables that represent the localized tree conditions such as Leaf Area Index (LAI)<sup>(42)</sup>,  
541 vegetation management (e.g. tree trimming) data or detailed tree density, location, height  
542 and species data to better capture this phenomenon.

## REFERENCES

1. North American Electric Reliability Corporation (NERC). Events Analysis: System Disturbance Reports, 1992 - Continuing [Internet]. Atlanta (GA): North American Electric Reliability Corporation; 2010 [cited 2015 Jul 1]. Available from: <http://www.nerc.com/pa/rrm/ea/System%20Disturbance%20Reports%20DL/Forms/AllItems.aspx>.
2. Liu H, Davidson RA, Rosowsky DV, Stedinger JR. Negative binomial regression of electric power outages in hurricanes. *Journal of Infrastructure Systems*, 2005; 11(4):258-267.
3. Guikema SD, Davidson RA, Liu H. Statistical models of the effects of tree trimming on power system outages. *Power Systems, IEEE Transactions on*, 2006; 11(3):1549-1557.
4. Liu H, Davidson RA, Apanasovich TV. Spatial generalized linear mixed models of electric power outages due to hurricanes and ice storms. *Reliability Engineering & System Safety*, 2008; 93(6):897-912.
5. Han SR, Guikema SD, Quiring SM, Lee KH, Rosowsky D, Davidson RA. Estimating the spatial distribution of power outages during hurricanes in the Gulf coast region. *Reliability Engineering & System Safety*, 2009; 94(2):199-210.
6. Guikema SD, Quiring SM, Han SR. Prestorm estimation of hurricane damage to electric power distribution systems. *Risk analysis*, 2010; 30(12):1744-1752.
7. Nateghi R, Guikema SD, Quiring SM. Comparison and validation of statistical methods for predicting power outage durations in the event of hurricanes. *Risk analysis*, 2011; 31(12):1897-1906.
8. Guikema SD, Quiring SM. Hybrid data mining-regression for infrastructure risk assessment based on zero-inflated data. *Reliability Engineering & System Safety*, 2012; 99:178-182.

- 566 9. Nateghi R, Guikema SD, Quiring SM. Forecasting hurricane-induced power outage dura-  
567 tions. *Natural Hazards*, 2014; 74(3):1795-1811.
- 568 10. Wanik DW, Anagnostou EN, Hartman BM, Frediani MEB, Astitha M. Storm outage  
569 modeling for an electric distribution network in Northeastern US. *Natural Hazards*, 2015.
- 570 11. Skamarock WC, Klemp JB, Dudhia J, Gill DO, Barker DM, Duda MG, Huang XY,  
571 Wang W, Powers JG. A description of the advanced research WRF version 3 [Internet].  
572 NCAR Technical Note, 2008 [cited 2015 Feb 15]; NCAR/TN-475+STR. Available from:  
573 [http://www2.mmm.ucar.edu/wrf/users/docs/arw\\_v3.pdf](http://www2.mmm.ucar.edu/wrf/users/docs/arw_v3.pdf).
- 574 12. National Centers for Environmental Prediction, National Weather Service, NOAA, U.S.  
575 Department of Commerce. NCEP global forecast system (GFS) analyses and forecasts  
576 [Internet]. Boulder (CO): Research Data Archive at the National Center for Atmospheric  
577 Research, Computational and Information Systems Laboratory; 2007 [cited 2015 Feb 15].  
578 Available from: <http://rda.ucar.edu/datasets/ds084.6/>.
- 579 13. National Centers for Environmental Prediction, National Weather Service, NOAA. NCEP  
580 ADP Global Upper Air and Surface Weather Observations (PREPBUFR format), May  
581 1997 - Continuing [Internet]. Boulder (CO): Research Data Archive at the National  
582 Center for Atmospheric Research, Computational and Information Systems Laboratory;  
583 2008 [cited 2015 Feb 15]. Available from: <http://rda.ucar.edu/datasets/ds337.0/>.
- 584 14. Sellier DD. Crown structure and wood properties: Influence on tree sway and response  
585 to high winds. *American Journal of Botany*, 2009; 96(5):885-896.
- 586 15. Ciftci C, Arwade SR, Kane B, Brena SF. Analysis of the probability of failure for open-  
587 grown trees during wind storms. *Probabilistic Engineering Mechanics*, 2014; 37:41-50.
- 588 16. Quiring SM, Zhu L, Guikema SD. Importance of soil and elevation characteristics for  
589 modeling hurricane-induced power outages. *Natural hazards*, 2011; 58(1):365-390.

- 590 17. Center for Land Use Education and Research, College of Agriculture and Natural Re-  
591 sources, University of Connecticut. Connecticut's Changing Landscape (CCL) Basic Land  
592 Cover Information, 1990 - Continuing [Internet]. Haddam (CT): Center for Land Use Ed-  
593 ucation and Research; 2006 [cited 2015 Feb 15]. Available from: [http://clear.uconn.](http://clear.uconn.edu/projects/landscape/)  
594 [edu/projects/landscape/](http://clear.uconn.edu/projects/landscape/).
- 595 18. Breiman L. Random forests. *Machine learning*, 2001; 45(1):5-32.
- 596 19. Meinshausen N. Quantile regression forests. *The Journal of Machine Learning Research*,  
597 2006; 7:983-999.
- 598 20. Meinshausen N. quantregForest: quantile regression forests. R package version 0.2-3  
599 [Internet]. Vienna (Austria): CRAN, Institute for Statistics and Mathematics, Vienna  
600 University of Economics and Business; 2012 [cited 2015 Feb 15]. Available from: [http:](http://cran.r-project.org/web/packages/quantregForest/)  
601 [//cran.r-project.org/web/packages/quantregForest/](http://cran.r-project.org/web/packages/quantregForest/).
- 602 21. Liaw A, Wiener M. The randomForest package. *R News*, 2002; 2(3):18-22.
- 603 22. Juban J, Fugon L, Kariniotakis G *et al*. Probabilistic short-term wind power forecasting  
604 based on kernel density estimators. *European Wind Energy Conference and Exhibition,*  
605 *EWEC*; 2007 May 7-10; Milan, Italy.
- 606 23. Francke T, López-Tarazón JA, Vericat D, Bronstert A, Batalla RJ. Flood-based analysis  
607 of high-magnitude sediment transport using a non-parametric method. *Earth Surface*  
608 *Processes and Landforms*, 2008; 33(13):2064-2077.
- 609 24. Zimmermann A, Francke T, Elsenbeer H. Forests and erosion: Insights from a study of  
610 suspended-sediment dynamics in an overland flow-prone rainforest catchment. *Journal of*  
611 *Hydrology*, 2012; 428:170-181.
- 612 25. Chipman HA, George EI, McCulloch RE. BART: Bayesian additive regression trees. Tech-  
613 nical Report, Department of Mathematics and Statistics, Acadia University, 2005 July.

- 614 26. Chipman HA, George EI, McCulloch RE *et al.* BART: Bayesian additive regression trees.  
615 The Annals of Applied Statistics, 2010; 4(1):266-298.
- 616 27. Hastie T, Tibshirani R *et al.* Bayesian backfitting (with comments and a rejoinder by the  
617 authors). Statistical Science, 2000; 15(3):196-223.
- 618 28. Chipman HA, George EI, McCulloch RE. Bayesian CART model search. Journal of the  
619 American Statistical Association, 1998; 93(443):935-948.
- 620 29. Friedman JH. Greedy function approximation: a gradient boosting machine. Annals of  
621 Statistics, 2001; 1189-1232.
- 622 30. Chipman HA, McCulloch RE. BayesTree: Bayesian methods for tree based models. R  
623 package version 0.3-1 [Internet]. Vienna (Austria): CRAN, Institute for Statistics and  
624 Mathematics, Vienna University of Economics and Business; 2009 [cited 2015 Feb 15].  
625 Available from: <http://cran.r-project.org/web/packages/BayesTree/>.
- 626 31. Pratola MT, Chipman HA, Gattiker JR, Higdon DM, McCulloch RE, Rust WN. Parallel  
627 Bayesian additive regression trees. Journal of Computational and Graphical Statistics,  
628 2014; 23(3):830-852.
- 629 32. Kapelner A, Bleich J. bartMachine: Machine Learning with Bayesian additive regression  
630 trees. R package version 1.1.1 [Internet]. Vienna (Austria): CRAN, Institute for Statistics  
631 and Mathematics, Vienna University of Economics and Business; 2014 [cited 2015 Feb  
632 15]. Available from: <http://cran.r-project.org/web/packages/bartMachine/>.
- 633 33. Blattenberger G, Fowles R. Avalanche forecasting: using Bayesian additive regression  
634 trees (BART). Demand for communications services—insights and perspectives. New York:  
635 Springer; 2014. Part III, Empirical Applications: Other Areas; p. 211-227.
- 636 34. Gehlke CE, Biehl K. Certain effects of grouping upon the size of the correlation coef-  
637 ficient in census tract material. Journal of the American Statistical Association, 1934;

- 638 29(185A):169-170.
- 639 35. Openshaw S, Taylor PJ. A million or so correlation coefficients: three experiments on  
640 the modifiable areal unit problem. *Statistical Applications in the Spatial Sciences*, 1979;  
641 21:127-144.
- 642 36. Nash JE, Sutcliffe JV. River flow forecasting through conceptual models part IA discussion  
643 of principles. *Journal of Hydrology*, 1970; 10(3):282-290.
- 644 37. Moriasi DN, Arnold JG, Van Liew MW, Bingner RL, Harmel RD, Veith TL. Model  
645 evaluation guidelines for systematic quantification of accuracy in watershed simulations.  
646 *Transactions of the ASABE*, 2007; 50(3):885-900.
- 647 38. Özkaynak H, Frey HC, Burke J, Pinder RW. Analysis of coupled model uncertainties in  
648 source-to-dose modeling of human exposures to ambient air pollution: A PM 2.5 case  
649 study. *Atmospheric Environment*, 2009; 43(9):1641-1649.
- 650 39. Hamill TM. Interpretation of rank histograms for verifying ensemble forecasts. *Monthly*  
651 *Weather Review*, 2001; 129(3):550-560.
- 652 40. Bickel PJ, Freedman DA. Some asymptotic theory for the bootstrap. *The Annals of*  
653 *Statistics*, 1981:1196-1217.
- 654 41. Reilly A, Guikema SD. Bayesian Multiscale Modeling of Spatial Infrastructure Perfor-  
655 mance Predictions with an Application to Electric Power Outage Forecasting. *Journal of*  
656 *Infrastructure Systems*, 2014; 21(2):04014036.
- 657 42. NASA Earth Observations (NEO). Leaf Area Index (LAI), 2000 - Continuing [Inter-  
658 net]. Greenbelt (MD): EOS Project Science Office, NASA Goddard Space Flight Center;  
659 2015 [cited 2015 Sep 15]. Available from: [http://neo.sci.gsfc.nasa.gov/view.php?](http://neo.sci.gsfc.nasa.gov/view.php?datasetId=MOD15A2_M_LAI)  
660 [datasetId=MOD15A2\\_M\\_LAI](http://neo.sci.gsfc.nasa.gov/view.php?datasetId=MOD15A2_M_LAI).

Table I: Explanatory Variables Included in Modeling

<b>Variable</b>	<b>Description</b>	<b>Type</b>	<b>Unit</b>
<b>Wind10m</b>	Sustained wind speed at 10 meters	Numerical	m/s
<b>Gust</b>	Wind speed of gust at 10 meters	Numerical	m/s
<b>WStress</b>	Wind stress	Numerical	-
<b>wgt9</b>	Duration of 10m wind greater than 9m/s	Numerical	hour
<b>ggt18</b>	Duration of gust greater than 18m/s	Numerical	hour
<b>ggt27</b>	Duration of gust greater than 27m/s	Numerical	hour
<b>PreRate</b>	Precipitation rate	Numerical	mm/hr
<b>TotPrec</b>	Total accumulated precipitation	Numerical	mm
<b>Temp</b>	Temperature	Numerical	°C
<b>SoilMst</b>	Soil moisture	Numerical	kg/kg
<b>SnoWtEq</b>	Snow water equivalent (only for winter)	Numerical	kg/m <sup>2</sup>
<b>sumAssets</b>	Sum of assets (infrastructure)	Numerical	count
<b>PercDeveloped</b>	Percentage of urban area	Numerical	%
<b>PercConif</b>	Percentage of coniferous trees	Numerical	%
<b>PercDecid</b>	Percentage of deciduous trees	Numerical	%
<b>seasoncat</b>	Leaves on (summer: Jun to Sep); Leaves off (winter: Dec to Mar); Transition (Oct, Nov, Apr and May).	Categorical	-

Table II: Comparison of QRF and BART with Hurricane Validation Data

<b>Hurricane (Outages)</b>	<b>Model</b>	<b>Predicted</b>	<b>Pred. Interval</b>	<b>MAE (By Town)</b>	<b>RMSE (By Town)</b>
<b>Irene (4890)</b>	QRF	4542	(4311, 4666)	8.86	15.70
	BART	4795	(4688, 4898)	6.12	9.43
	MEAN	5200	-	23.95	35.25
<b>Sandy (5052)</b>	QRF	5060	(4674, 5110)	6.91	11.02
	BART	5171	(5039, 5302)	5.50	8.02
	MEAN	5121	-	28.89	52.03

Figure 1: Spatial Resolutions: 2-km Grid Cell, Town, Division and Territory. (Grid cells without infrastructure or outside the service territory are excluded.)

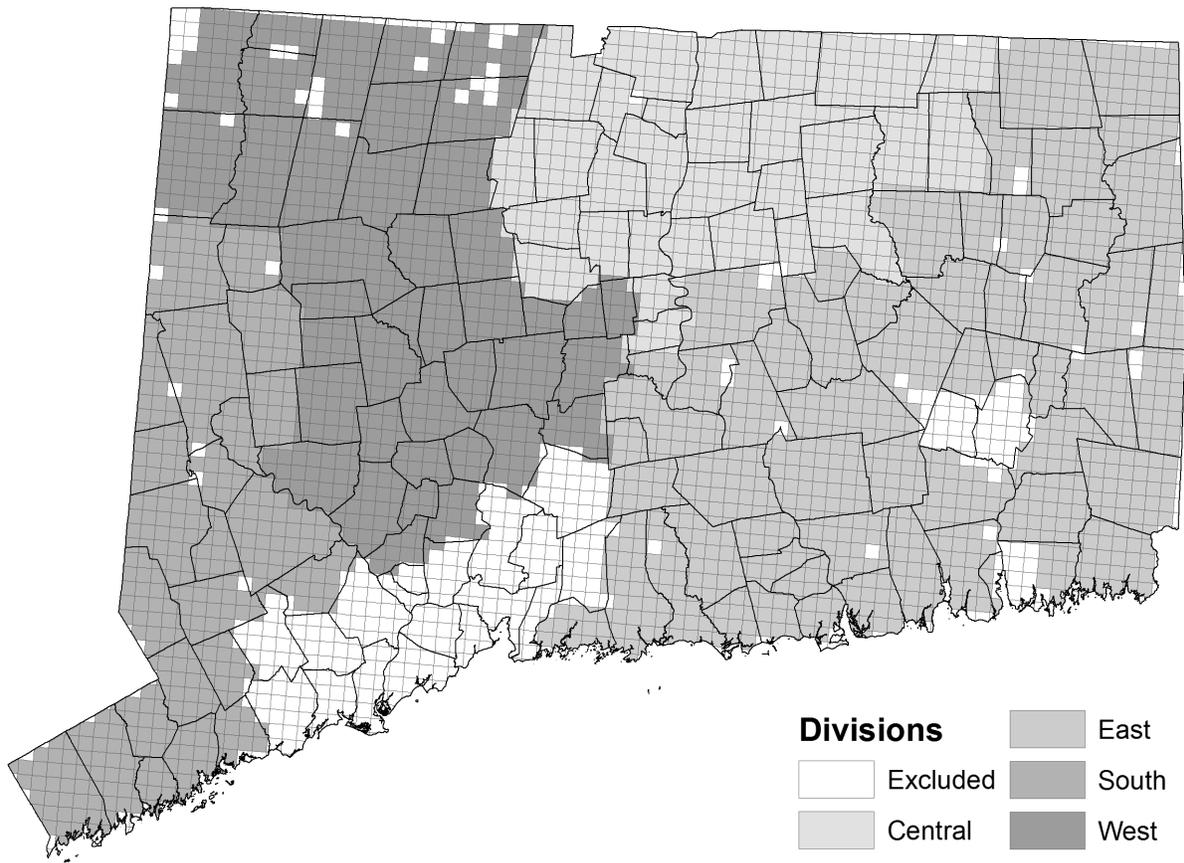


Figure 2: Storm Events Simulation: Weather Research and Forecasting Model Nested Domains in 18 km, 6 km and 2 km grids.

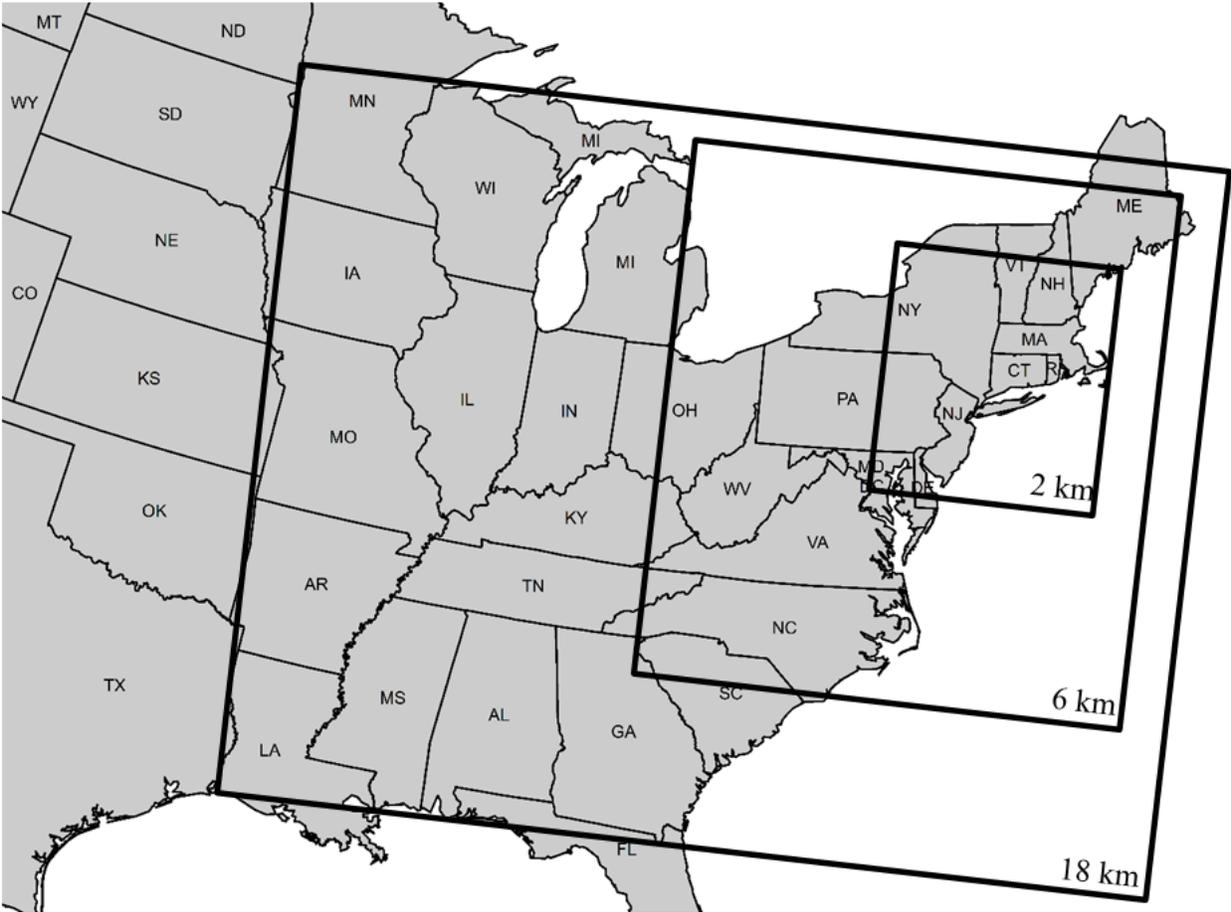


Figure 3: Density of Outages in Different Resolutions: (a) Grid Cell, (b) Town, (c) Division, (d) Territory. (Dotted lines stand for gaussian kernel density estimations.)

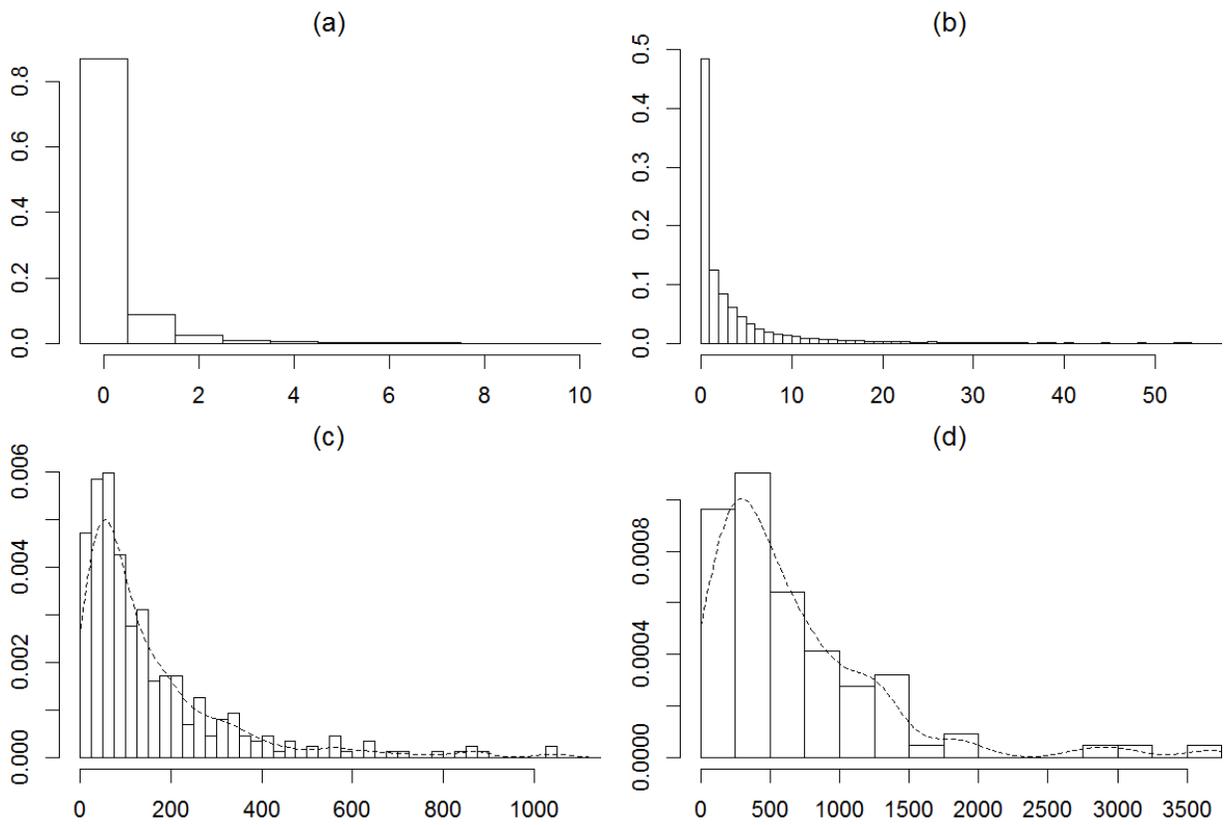


Figure 4: Comparison of QRF and BART in Town Resolution: (a) Irene, (b) Sandy. (80% prediction intervals are given, as well as their coverage rates.)

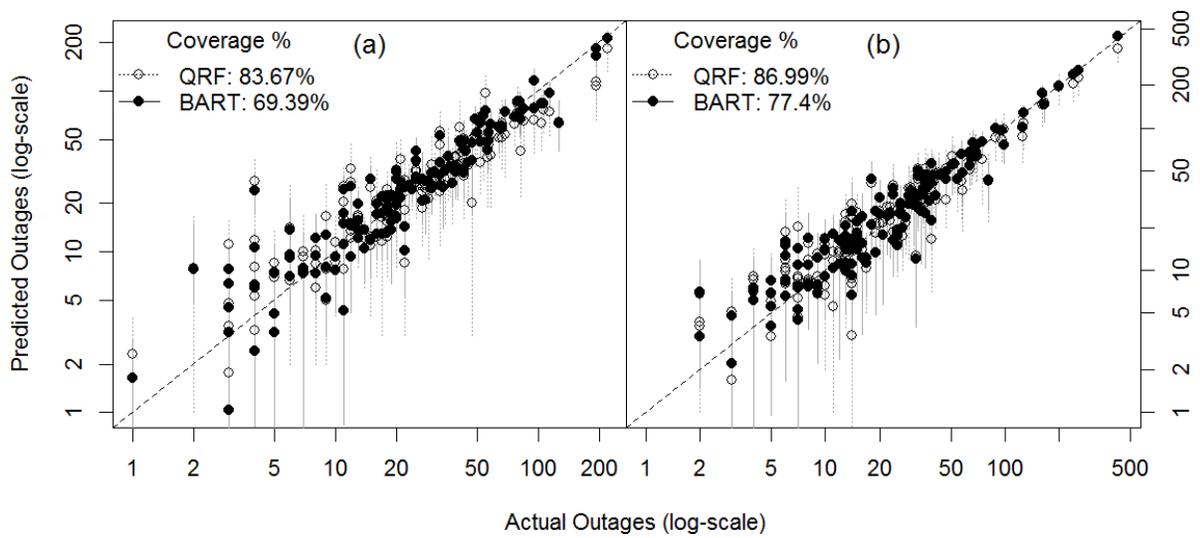


Figure 5: Comparison of QRF and BART in Modeling Spatial Variability: Total Number of Outages by Town for (a) Actual Number of Irene, (b) Actual Number of Sandy, (c) QRF Validation of Irene, (d) QRF Validation of Sandy, (e) BART Validation of Irene and (f) BART Validation of Sandy.



Figure 6: Predictions for Each Storm in Validation Dataset. (80% prediction intervals are given, as well as their coverage rates.)

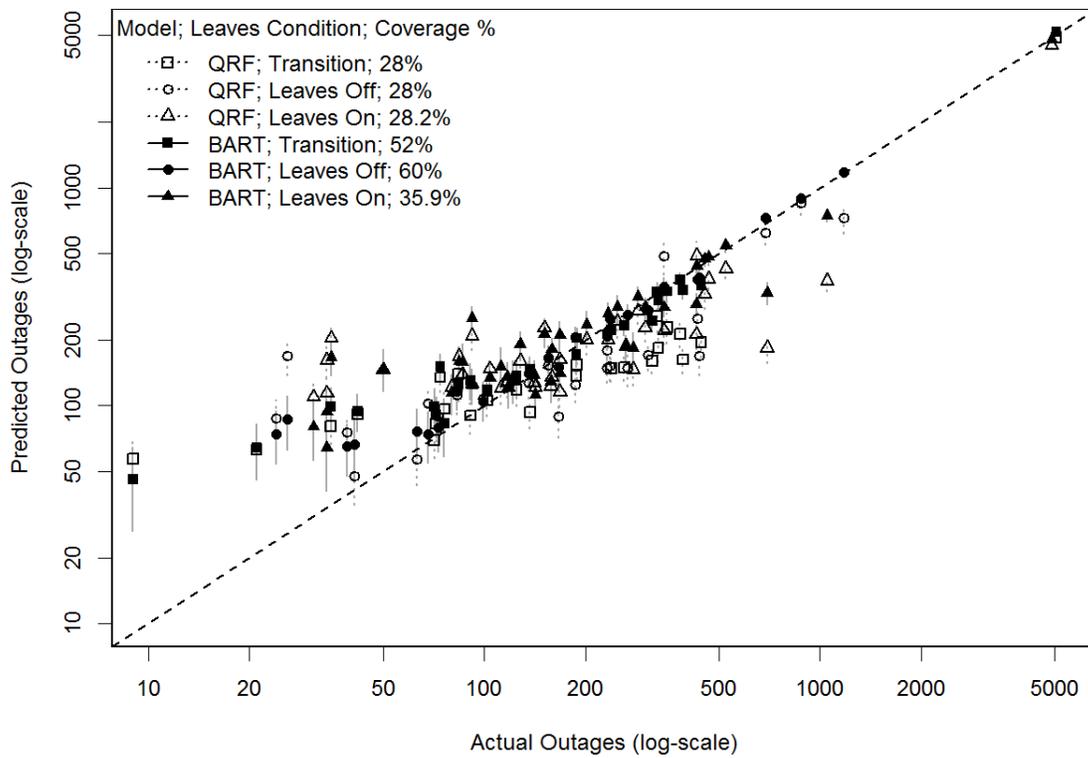


Figure 7: Nash-Sutcliffe Efficiency (NSE) or Relative Error (RE) for Each Storm in Different Resolutions: (a) NSE for Grid Cells, (b) NSE for Towns, (c) NSE for Divisions, (d) RE for Territory.

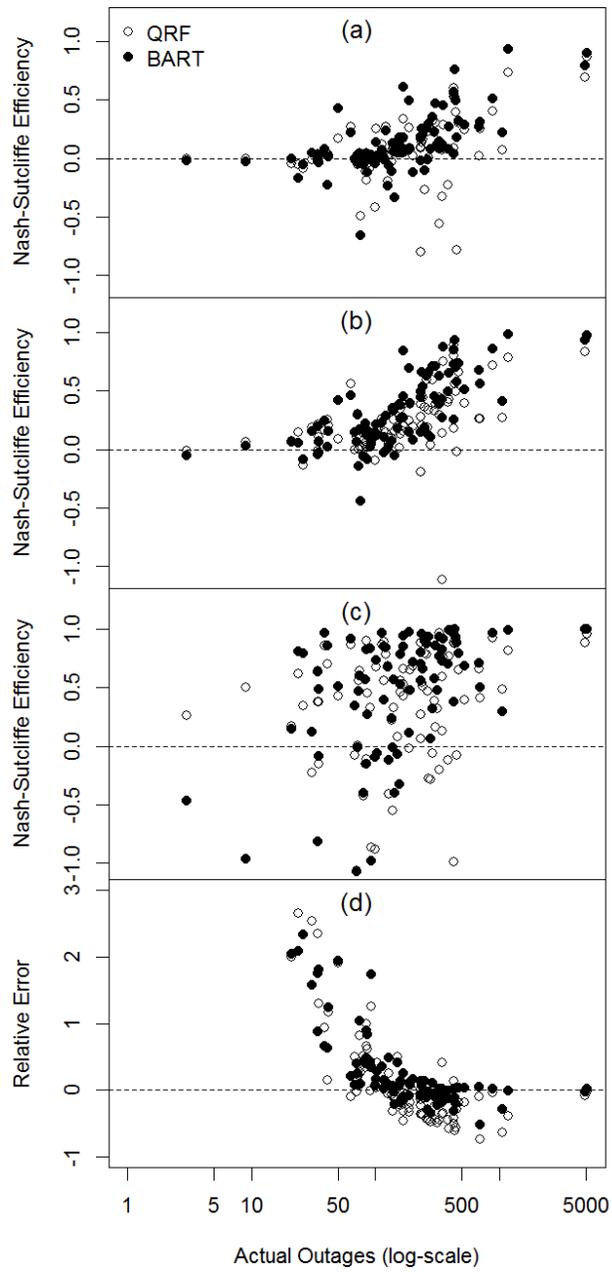


Figure 8: Nash-Sutcliffe Efficiency (NSE) or Relative Error (RE) for Each Season in Different Resolutions: (a) NSE for Grid Cells, (b) NSE for Towns, (c) NSE for Divisions, (d) RE for Territory.

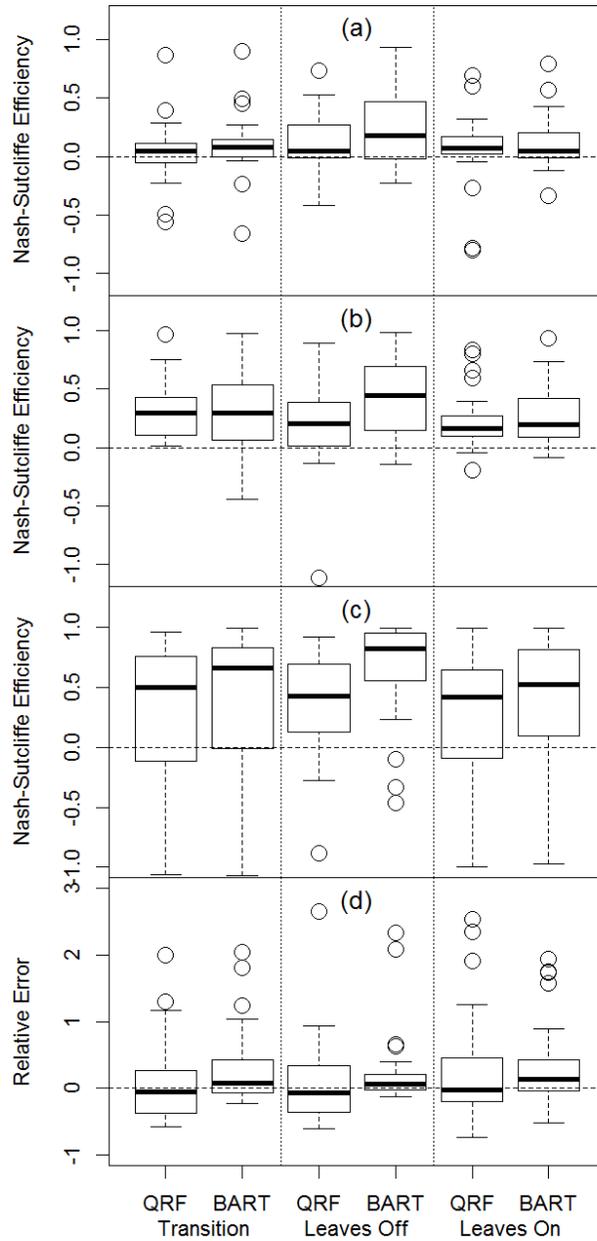


Figure 9: Uncertainty Ratio (UR) for Each Storm in Different Resolutions: (a) Grid Cell, (b) Town, (c) Division, (d) Territory.

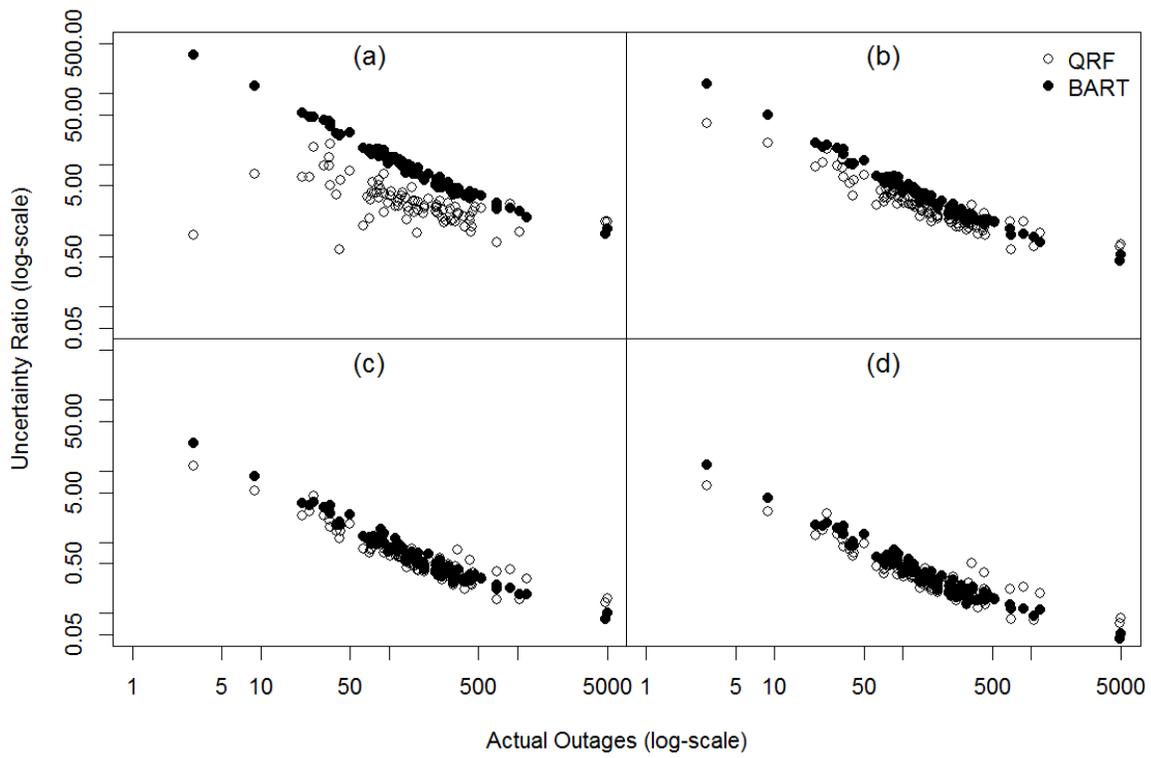


Figure 10: Exceedance Probability (EP) for Each Storm in Different Resolutions: (a) Grid Cell, (b) Town, (c) Division, (d) Territory. (A small amount of noise is add in (c) and (d) to avoid overlap of points.)

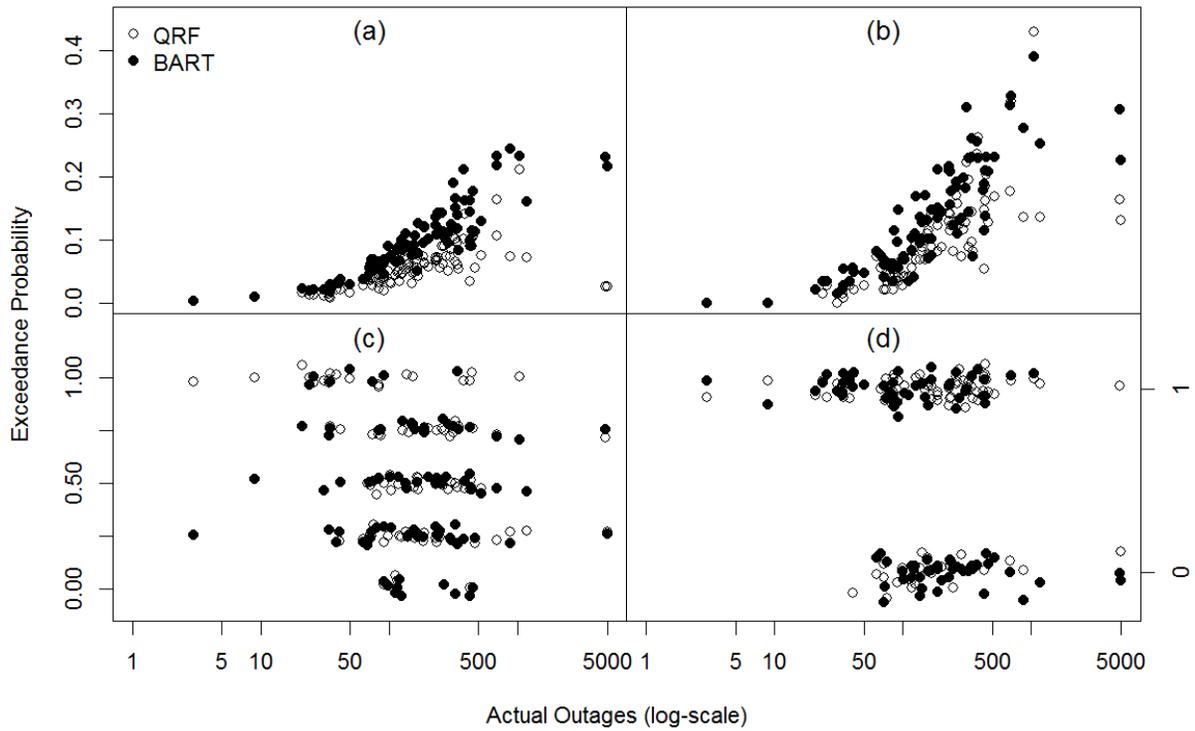


Figure 11: Rank Histogram of QRF Predictions in Different Resolutions: (a) Grid Cell, (b) Town, (c) Division, (d) Territory.

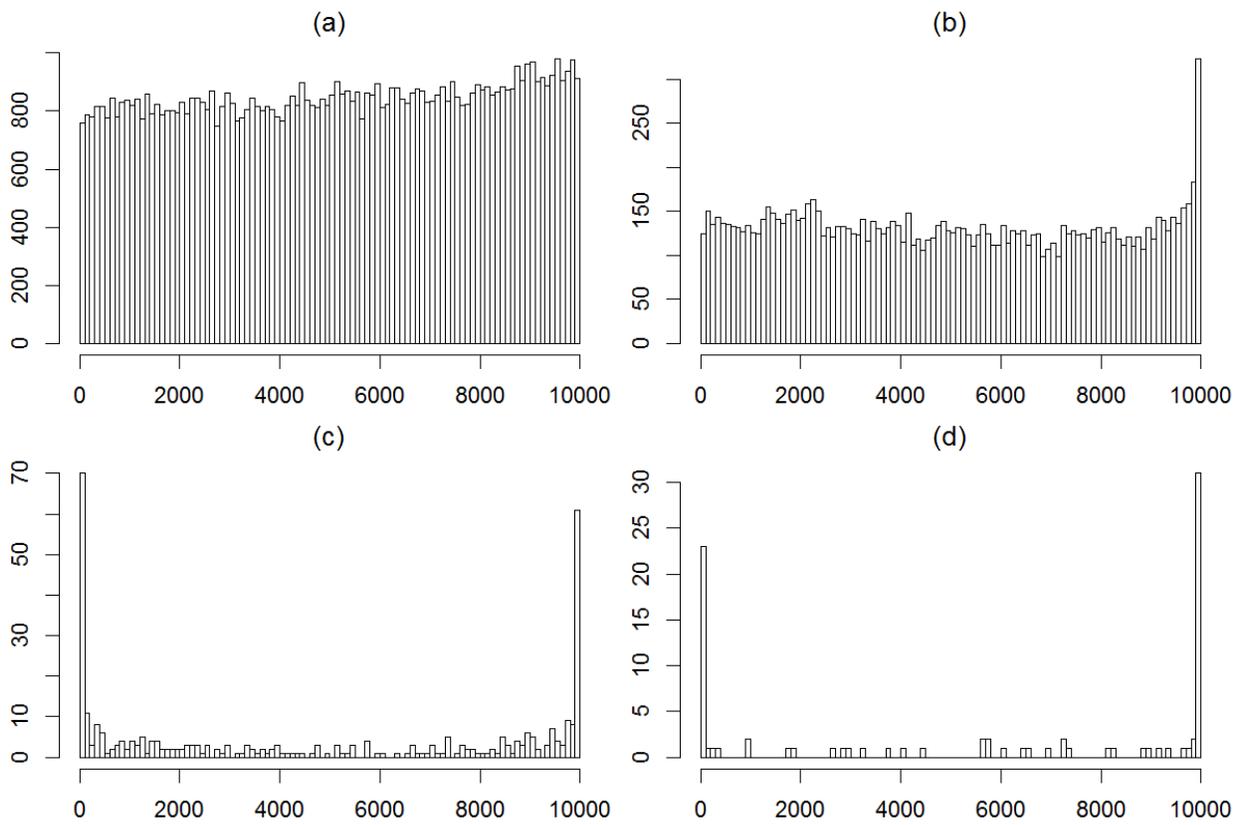


Figure 12: Rank Histogram of BART Predictions in Different Resolutions: (a) Grid Cell, (b) Town, (c) Division, (d) Territory.

