

1 **Machine Learning Enhancement of Storm Scale Ensemble**

2 **Probabilistic Quantitative Precipitation Forecasts**

3 DAVID JOHN GAGNE II*

School of Meteorology, University of Oklahoma, Norman, Oklahoma

4 AMY MCGOVERN

School of Computer Science, University of Oklahoma, Norman, Oklahoma

5 MING XUE

Center for the Analysis and Prediction of Storms/School of Meteorology

University of Oklahoma, Norman, Oklahoma

* *Corresponding author address:* David John Gagne II, 120 David L. Boren Blvd., Suite 5900, Norman, OK 73072. Email: djgagne@ou.edu

ABSTRACT

7 Probabilistic quantitative precipitation forecasts challenge meteorologists due to the wide
8 variability of precipitation amounts over small areas and their dependence on conditions at
9 multiple spatial and temporal scales. Ensembles of convection-allowing numerical weather
10 prediction models offer a way to produce improved precipitation forecasts and estimates
11 of the forecast uncertainty. These models allow for the prediction of individual convective
12 storms on the model grid, but they often displace the storms in space, time, and inten-
13 sity, which results in added uncertainty. Machine learning methods can produce calibrated
14 probabilistic forecasts from the raw ensemble data that correct for systemic biases in the
15 ensemble precipitation forecast and incorporate additional uncertainty information from ag-
16 gregations of the ensemble members and additional model variables. This study utilizes the
17 2010 Center for the Analysis and Prediction of Storms Storm Scale Ensemble Forecast system
18 and the National Severe Storms Laboratory National Mosaic and Multisensor Quantitative
19 Precipitation Estimate as input data for training logistic regressions and random forests to
20 produce a calibrated probabilistic quantitative precipitation forecast. The reliability and
21 discrimination of the forecasts are compared through verification statistics and a case study.

22 1. Introduction

23 Most flooding fatalities occur due to flash floods, in which waters rise and fall rapidly
24 due to concentrated rainfall over a small area (Ashley and Ashley 2008). The first step
25 to anticipating flash floods is quantitative forecasting of the amount, location, and tim-
26 ing of precipitation. These quantitative precipitation forecasts are challenging due to the
27 wide variability of precipitation amounts over small areas, the dependence of precipitation
28 amounts on processes at a wide range of scales, and the dependence of extreme precipitation
29 on any precipitation actually occurring (Bremnes 2004; Ebert 2001; Doswell et al. 1996).
30 Recent advances in numerical modeling and machine learning are working to address these
31 challenges.

32 Numerical Weather Prediction (NWP) models are now being run experimentally at 4 km
33 horizontal grid spacing, or storm scale, allowing for the formation of individual convective
34 cells without a convective parameterization scheme. These models better represent storm
35 processes and output hourly predictions, but they have the challenge of correctly placing
36 and timing the precipitation compared to models with coarser grid spacing and temporal
37 resolution. An ensemble of storm scale NWP models can provide improved estimates of un-
38 certainty compared to coarser ensembles due to better sampling of the spatiotemporal errors
39 associated with individual storms (Clark et al. 2009). As each ensemble member produces
40 predictions of precipitation and precipitation ingredients, the question then becomes how
41 best to combine those predictions into the most accurate and useful consensus guidance.

42 The final product should highlight areas most likely to be impacted by heavy rain and
43 also provide an uncertainty estimate for that impact. Probabilistic Quantitative Precipi-

44 tation Forecasts (PQPFs) incorporate both of these qualities. A good PQPF should have
45 reliable probabilities, such that a 40% chance of rain verifies 40% of the time over a large
46 sample (Murphy 1977). PQPFs should also discriminate between extreme and trace precipi-
47 tation events consistently, so most extreme events occur with higher probabilities, and most
48 trace precipitation events are associated with low probabilities. Since individual ensemble
49 members may be biased in different situations, a simple count of the ensemble members that
50 exceed a threshold will often result in unreliable forecasts that discriminate poorly. Incor-
51 porating trends from past forecasts and additional information from other model variables
52 can offset these biases and produce an enhanced PQPF.

53 Ensemble weather prediction (Toth and Kalnay 1993; Tracton and Kalnay 1993; Molteni
54 et al. 1996) has required various forms of statistical post-processing to produce accurate pre-
55 cipitation forecasts and uncertainty estimates from the ensembles, but most previous studies
56 used coarser ensembles and longer forecast windows for calibration. The rank histogram
57 method (Hamill and Colucci 1997) showed that ensemble precipitation forecasts tended to
58 be underdispersive. Linear regression calibration methods have shown some skill improve-
59 ments in Hamill and Colucci (1998), Eckel and Walters (1998), Krishnamurti et al. (1999),
60 and Ebert (2001). Hall et al. (1999), Koizumi (1999), and Yuan et al. (2007) applied neural
61 networks to precipitation forecasts and found increases in performance over linear regression.
62 Logistic regression, a transform of a linear regression to fit an S-shaped curve ranging from 0
63 to 1, has shown more promise, as in Applequist et al. (2002), which tested linear regression,
64 logistic regression, neural networks, and genetic algorithms on 24-hour PQPF and found
65 that logistic regression consistently outperformed the other methods. Hamill et al. (2004,
66 2008) also utilized the logistic regression with an extended training period for added skill.

67 Storm scale ensemble precipitation forecasts have been post-processed with smoothing
68 algorithms that produce reliable probabilities within longer forecast windows. Clark and
69 Coauthors (2011) applied smoothing algorithms at different spatial scales to the SSEF pre-
70 cipitation forecasts and compared verification scores. Johnson and Wang (2012) compared
71 the skill of multiple calibration methods on neighborhood and object-based probabilistic
72 forecasts from the 2009 SSEF. Marsh et al. (2012) applied a Gaussian kernel density es-
73 timation function to the NSSL 4 km Weather Research and Forecasting (WRF) to derive
74 probabilities from deterministic forecasts. These methods are helpful for predicting larger
75 scale events but smooth out the threats from extreme precipitation in individual convective
76 cells. Gagne II et al. (2012) took the first step in examining how multiple machine learning
77 approaches performed in producing probabilistic, deterministic, and quantile precipitation
78 forecasts over the central United States at individual grid points.

79 The purpose of this paper is to analyze PQPF predictions produced by multiple machine
80 learning techniques incorporating data from the Center for the Analysis and Prediction of
81 Storms (CAPS) 2010 Storm Scale Ensemble Forecast (SSEF) system (Xue et al. 2011; Kong
82 et al. 2011). In addition to the choice of algorithm, some variations in the algorithm setups
83 are also examined. The strengths and weaknesses of the machine learning algorithms are
84 shown through the analysis of verification statistics, variables chosen by the machine learning
85 models, and a case study. This paper expands on the work presented in Gagne II et al. (2012)
86 by including statistics for the eastern US, a larger training set more representative of the
87 precipitation probabilities, a different case study day with comparisons of multiple runs,
88 multiple precipitation thresholds, and more physical justification for the performance of the
89 machine learning algorithms.

90 2. Data

91 a. *Ensemble Data*

92 The Center for the Analysis and Prediction of Storms (CAPS) 2010 Storm Scale Ensem-
93 ble Forecast (SSEF) system (Xue et al. 2011; Kong et al. 2011) provides the input data for
94 the machine learning algorithms. The 2010 SSEF consists of 19 individual model members
95 from the WRF Advanced Research WRF (ARW), 5 members from the WRF Nonhydrostatic
96 Mesoscale Model (NMM), and 2 members from the CAPS Advanced Research Prediction
97 System (ARPS; Xue et al. 2000, 2001, 2003). Each member has a varied combination of mi-
98 crophysics schemes, land surface models, and planetary boundary layer schemes. The SSEF
99 ran every weekday at 0000 UTC in support of the 2010 National Oceanic and Atmospheric
100 Administration/Hazardous Weather Testbed Spring Experiment (Clark et al. 2012), which
101 ran from May 3 to June 18, for a total of 34 runs. The SSEF provides hourly model output
102 over the contiguous United States at 4 km horizontal grid spacing out to 30 hours. Of the 26
103 members, the 14 members included initial condition perturbations derived from the National
104 Center for Environmental Prediction Short Range Ensemble Forecast (SREF; Stensrud et al.
105 1999) members, and only they are included in our post-processing procedure. The 12 other
106 members used the same control initial conditions although with different physics options and
107 models; designed to examine forecast sensitivities to physics parameterizations, they do not
108 contain the full set of initial conditions and model uncertainties and are therefore excluded
109 from our study.

110 *b. Verification Data*

111 A radar-based verification dataset was used as the verifying observations for the SSEF.
112 The National Mosaic and Multi-Sensor QPE (NMQ; Vasiloff et al. 2007) derives precipitation
113 estimates from the reflectivity data of the NEXRAD radar network. The estimates are made
114 on a grid with 1 km horizontal spacing over the continental US (CONUS). The original grid
115 has been bi-linearly interpolated to the same grid as the SSEF.

116 *c. Data Selection and Aggregation*

117 The relative performance of any machine learning algorithm is conditioned on the dis-
118 tribution of its training data. The sampling scheme for the SSEF is conditioned on the
119 constraints of 34 ensemble runs over a short, homogenous time period with 840,849 grid
120 points from each of the 30 time steps. The short training period and large number of grid
121 points preclude training a single model at each grid point, so a regional approach was used.

122 The SSEF domain was split into thirds (280,283 points per time step), and points were
123 selected with a uniform random sample from each subdomain in areas with quality radar
124 coverage. The gridded Radar Quality Index (RQI; Zhang et al. 2011) was evaluated at each
125 grid point to determine the trustworthiness of the verification data. Points with an $RQI > 0$
126 were located within the useful range of a NEXRAD radar and included in the sampling.
127 For points with precipitation values less than 0.25 mm, 0.04% were sampled, and for points
128 with more precipitation, 0.4% were sampled. Grid 0 corresponds to the western third of
129 the CONUS, Grid 1 corresponds to the central third, and Grid 2 corresponds to the eastern
130 third. Fig. 1 shows that the north central section of the United States and the states south

131 of the Great Lakes were most heavily sampled. Grid 0 was excluded from post-processing
132 due to the low frequency of heavy precipitation for most of the region.

133 A comparison of the sampled rainfall distributions and the full rainfall distributions
134 for each subgrid are shown in Fig. 2 and Table 1. Undersampling of the 0 and trace
135 (below 0.25 mm) precipitation points was necessary because of the large number of no-
136 precipitation events, which overwhelmed the signal from the actual precipitation events.
137 The random sampling of grid points helps reduce the chance of sampling multiple grid
138 points from the same storm without explicitly filtering subsets of the domain, which was
139 performed by Hamill et al. (2008). Filtering grid points surrounding each sampled point is
140 extremely computationally expensive because filtering requires multiple passes over the grid
141 while random sampling only requires 1 grid reordering.

142 Relevant model output variables, called predictors, (Table 2) were also extracted from
143 each ensemble member at each sampled point. These predictors captured additional infor-
144 mation about the mesoscale and synoptic conditions in each model. The predictors from
145 each ensemble member were then aggregated into 4 descriptive statistics for each predictor.
146 The mean provided the most likely forecast value, and the standard deviation estimated the
147 spread, and the minimum and maximum showed the extent of the forecasted values.

148 3. Methods

149 a. Machine Learning Methods

150 1) LOGISTIC REGRESSION

151 One of the goals of this study is to compare the skill of more advanced machine learning
152 methods with more traditional statistical methods. Logistic regression was used as the
153 baseline statistical method. Logistic regressions are linear regression models in which a logit
154 transformation is applied to the data so that the predicted values will range between 0 and 1.
155 Two formulations of logistic regression were used: simple and multiple. The first formulation
156 uses the ensemble mean precipitation forecast as expressed in Eqn. 1:

$$157 \quad p(R \geq t|x_1) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1)}} \quad (1)$$

158 in which R is the precipitation amount, t is the threshold, x_1 is the ensemble mean pre-
159 cipitation forecast, β_0 is the intercept term, and β_1 is the weight given to the ensemble
160 mean precipitation forecast. Observed precipitation with an amount greater than or equal
161 to the chosen threshold is assigned a probability of 1, and precipitation amounts below the
162 threshold are assigned a probability of 0. The β coefficients are estimated by iteratively
163 maximizing the log-likelihood of the transformed training data. This formulation will adjust
164 the probability of the ensemble forecast based on systematic biases in the mean but will not
165 change the areal coverage of the precipitation forecasts. The second formulation incorpo-
166 rates multiple variables with stepwise variable selection similar to standard Model Output

167 Statistics (Glahn and Lowry 1972) approaches (Eqn. 2):

$$168 \quad p(R \geq t|x_1, x_2, \dots, x_n) = \frac{1}{1 + \exp \left[- \left(\beta_0 + \sum_{n=1}^N \beta_n x_n \right) \right]} \quad (2)$$

169 in which x_1 through x_n are the predictors chosen from Table 2 and β_1 through β_n are the
170 weights for those predictors. The terms are determined through a forward selection process
171 that finds the set of up to 15 terms that minimize the Akaike Information Criterion (Akaike
172 1974), which rewards goodness of fit but penalizes for large numbers of terms. This approach
173 does produce the best fit regression model given the available parameters, but the searching
174 process can take extensive time given the large dimensionality of the training set. The
175 *glm* and *step* functions in R are used to generate the logistic regressions. A more detailed
176 description of logistic regression and the fitting process can be found in James et al. (2013).

177 2) RANDOM FOREST

178 An non-parametric, nonlinear alternative to linear regression is the classification and re-
179 gression decision tree (Breiman 1984). Decision trees recursively partition a multidimensional
180 dataset into successively smaller subdomains by selecting variables and decision thresholds
181 that maximize a dissimilarity metric. At each node in the tree, every predictor is evaluated
182 with the dissimilarity metric, and the predictor and threshold with the highest metric value
183 are selected as the splitting criteria for that node. After enough partitions, each subdomain
184 is similar enough that the prediction can be approximated with a single value. The primary
185 advantages of decision trees are that they can be human-readable and perform variable se-

186 lection as part of the model growing process. The disadvantages lie in the brittleness of the
187 trees. Trees can undergo significant structural shifts due to small variations in the training
188 data, which results in large error variance.

189 Random forests (Breiman 2001) consist of an ensemble of classification and regression
190 trees (Breiman 1984) with two key modifications. First, the training data cases are bootstrap
191 resampled with replacement for each tree in the ensemble. Second, a random subset of the
192 predictors are selected for evaluation at each node. The final prediction from the forest is
193 the mean of the predicted probabilities from each tree. Random forests can produce both
194 probabilistic and regression predictions through this method. The random forest method
195 contains a few advantages that often lead to performance increases over traditional regres-
196 sion methods. The averaging of the results from multiple trees produces a smoother range
197 of values than individual decision trees while also reducing the sensitivity of the model pre-
198 dictions to minor differences in the training set (Strobl et al. 2008). The random selection of
199 predictors within the tree-building process allows for less optimal predictors to be included
200 in the model and increases the likelihood of the discovery of interaction effects among pre-
201 dictors that would be missed by the stepwise selection method used in logistic regression
202 (Strobl et al. 2008). Random forests have been shown to improve predictive performance
203 on multiple problem domains in meteorology, including storm classification (Gagne II et al.
204 2009), aviation turbulence (Williams 2013), and wind energy forecasting (Kusiak and Verma
205 2011). For this project, we used the R *randomForest* library, which implements the original
206 approach (Breiman 2001). For the parameter settings, we chose to use 100 trees, a minimum
207 node size of 20, and the default values for all other parameters. A more detailed description
208 of random forest can be found in James et al. (2013).

209 In addition to gains in performance, the random forest can also be used to rank the
 210 importance of each input variable (Breiman 2001). Variable importance is computed by
 211 first calculating the accuracy of each tree in the forest on classifying the cases that were not
 212 selected for training, known as the out-of-bag cases. Within the out-of-bag cases, the values
 213 of each variable are randomly rearranged, or permuted, and those cases are then re-evaluated
 214 by each tree. The mean variable importance score is then the difference in prediction accuracy
 215 on the out-of-bag cases averaged over all trees. Variable importance scores can vary randomly
 216 among forests trained on the same dataset, so the variable importance scores from each of
 217 the 34 forests trained for cross-validation were averaged together for a more robust ranking.

218 *b. Evaluation Methods*

219 Two scores were used to assess the probabilistic forecasts. The Brier Skill Score (BSS)
 220 (BSS; Brier 1950) is one method used to evaluate probabilistic forecasts. The Brier Skill
 221 Score can be decomposed into three terms (Murphy 1973), as shown in Eq. 3:

$$222 \quad BSS = \frac{\frac{1}{N} \sum_{k=1}^K n_k (\bar{o}_k - \bar{o})^2 - \frac{1}{N} \sum_{k=1}^K n_k (p_k - \bar{o}_k)^2}{\bar{o}(1 - \bar{o})} \quad (3)$$

223 N is the number of forecasts, K is the number of probability bins, n_k is the number of
 224 forecasts in each probability bin, \bar{o}_k is the observed relative frequency for each bin, \bar{o} is the
 225 climatological frequency, and p_k is the forecast probability for a particular bin k . The first
 226 term in the numerator describes the resolution of the forecast probability, which should be
 227 maximized and increases as the observed relative frequency differs more from climatology.

228 The second term in the numerator describes the reliability of the forecast probability, which
229 should be minimized and decreases with smaller differences between the forecast probability
230 and observed relative frequency. The denominator term is the uncertainty, which is based
231 on the climatological probability of precipitation and cannot be reduced through calibration.
232 BSS increases with skill from 0. The components of the BSS can be displayed graphically
233 with an attributes diagram (Wilks 2011), in which the observed relative frequency of binned
234 probability forecasts are plotted against lines showing perfect reliability, no skill where the
235 reliability and resolution are equal, and no resolution where the observed relative frequency
236 and probability equal climatology.

237 The Area Under the Relative Operating Characteristic (ROC) curve, or AUC (Mason
238 1982) evaluates how well a probabilistic forecast correctly identifies heavy precipitation
239 events compared to identifying the events based on random chance. To calculate AUC,
240 first, the decision probability threshold is varied from 0 to 1 in equal steps. At each step, a
241 contingency table is constructed by splitting the probabilities into two categories with the
242 decision threshold (Table 3). Using Table 3, the following scores can be computed (Table 4).
243 Probability of detection (POD) is the ratio of hits to the total number of observed events.
244 The False Alarm Ratio (FAR) accounts for the number of false alarms compared to total
245 yes forecasts. The Probability of False Detection is the ratio of false alarms to total no
246 forecasts. The Equitable Threat Score (ETS) is the skill score officially used to assess the
247 skill of 24 hour precipitation forecasts, and is sensitive to the climatological base rate of the
248 validation data. The Peirce Skill Score (PSS) is another skill score that is insensitive to the
249 climatological base rate but is sensitive to hedging by over forecasting rare events (Doswell
250 et al. 1990). The bias ratio compares determines if false alarms or misses are more prevalent.

251 From the contingency table, the POD and POFD are calculated and plotted against each
252 other, forming a ROC curve. The AUC is the area between the right and bottom sides of the
253 plot and the curve itself. AUC with values above 0.5 has positive skill. AUC only determines
254 how well the forecast discriminates between two categories, so it does not take the reliability
255 of the forecast into account.

256 The ROC curve also can be used to determine an “optimal” decision threshold to convert
257 a probabilistic forecast to a deterministic forecast. As the decision probability increases, the
258 POD decreases from 1 to 0 while 1-POFD increases from 0 to 1. At some decision threshold,
259 POD and 1-POFD should be equal. At this optimal threshold (OT), the user has an equal
260 chance of correctly detecting rain events and no-rain events. This point occurs where the
261 ROC curve is farthest from the positive diagonal “no-skill” line. In addition, the vertical
262 distance between the ROC curve and the positive diagonal is the Peirce Skill Score (PSS;
263 Peirce 1884; Hansen and Kuipers 1965) since $PSS = POD - POFD$ (Manzato 2007). Since the
264 PSS can be calculated from the contingency table at each decision threshold, finding the
265 maximum PSS also finds the OT. The performance of the threshold choice can be validated
266 by treating the predictions as a binary classification problem and using the binary verification
267 statistics derived from the contingency table (Table 4).

268 *c. Experimental Procedure*

269 Each model was trained and evaluated using a leave-one-day-out cross validation pro-
270 cedure. For the 34 daily model runs in the training set, each machine learning model was
271 trained on 33 days and tested on 1. Separate models were trained for each 6 hour forecast

272 period and each subgrid to better capture the trends in the weather and ensemble disper-
273 siveness. Models were trained at two thresholds, 0.25 mm h^{-1} and 6.35 mm h^{-1} . The 0.25
274 mm h^{-1} threshold models were used to determine if rain was to occur at a point or not,
275 and the 6.35 mm h^{-1} models were trained and evaluated only on points at which either
276 rain occurred or the ensemble mean precipitation forecast predicted rain. Deterministic rain
277 forecasts were derived from the 0.25 mm h^{-1} models with the ROC curve optimal threshold
278 technique. These deterministic forecasts were used to mask 0-precipitation areas in the case
279 study. Additional models were trained at the 2.54 and 12.70 mm h^{-1} thresholds to evaluate
280 the skill in predicting lighter and heavier precipitation, but only the 6.35 mm h^{-1} models
281 were physically evaluated with variable importance and the case study. The 0.25 , 2.54 , 6.35 ,
282 and 12.70 mm h^{-1} were chosen because they correspond to the 0.01 , 0.1 , 0.25 , and 0.5 in
283 h^{-1} thresholds, respectively, that are used for determining trace and heavy precipitation
284 amounts.

285 4. Results

286 a. *Deterministic Rain Model Evaluation*

287 Evaluation of the deterministic forecasts of precipitation location shows that the multiple
288 logistic regression and random forest do add skill compared to the uncalibrated ensemble
289 probability and the simple logistic regression. Fig. 3 shows only slight variations in the opti-
290 mal threshold with forecast hour and that the thresholds for each machine learning algorithm
291 are similar. The low threshold for the raw ensemble indicates that the best precipitation cov-

292 erage is found when any ensemble member forecasts rain. PSS and ETS show similar hourly
293 trends, but PSS is higher than ETS since ETS tends to maximize at a higher probability
294 threshold. The multi-predictor machine learning algorithms (random forest and multiple
295 logistic regression) provide the most improvement for the first 15 hours of the forecast with
296 only a slight improvement for the afternoon and evening. The improvement comes from
297 an increased POD without also increasing the FAR. Similar results are seen in grid 2 in
298 terms of thresholds and score values, and the multi-predictor algorithms are able to provide
299 consistent improvement over all forecast hour (Fig. 4). Similar results to these are found in
300 the SSEF verification from Kong et al. (2011) although the ETS is lower in that paper due
301 to that study verifying against all grid points, including points with no radar coverage.

302 *b. Probabilistic Threshold Exceedance Model Evaluation*

303 1) EVALUATION OVER ALL FORECAST HOURS

304 The attributes diagrams for the raw ensemble and machine learning algorithms show
305 how the forecasts were altered to improve their calibration and skill. Fig. 5 contains the
306 attributes diagrams for the ensemble and algorithms applied to grids 1 and 2 at the 6.35 mm
307 h^{-1} threshold for all forecast hours. The raw ensemble in both domains tends to be over-
308 confident with probabilistic forecasts above climatology and under confident with forecasts
309 below climatology, resulting in a negative BSS. In spatial terms, the ensemble probabilities
310 are too high for areas where it forecasts heavy rain, and it is missing some areas where rain
311 actually occurred. The simple logistic regression improves the ensemble forecast by rescaling
312 the probabilities based on the ensemble mean rain amount. In grid 1 and grid 2 this gener-

313 ally results in all of the probabilities being scaled closer to the climatological probability to
314 address the general overconfidence of the ensemble. This scaling brings the observed relative
315 frequencies above climatology into the positive skill zone, but they are still overconfident.

316 The multiple logistic regression and random forest incorporate additional ensemble vari-
317 ables to add more information to the probability estimates. This does result in a large
318 improvement in BSS over the simple logistic regression. The multiple logistic regression
319 performs the best in grid 1 because it produces nearly perfect reliability between 0 and 50%
320 while maintaining positive skill from 60 to 90%. The random forest also places its forecasts
321 close to the perfect reliability line but slightly further away than the stepwise logistic re-
322 gression, and the forecast probabilities do not exceed 70%, so the random forest has a lower
323 maximum resolution than the multiple logistic regression.

324 In grid 2, the BSS for all three models decreases. The observed relative frequency for
325 the random forest lies just above the no skill line and then increases above 50%, which may
326 be a due to chance from the low sample size in the higher percentage range. The multiple
327 logistic regression performs better with the lower probability forecasts but has negative skill
328 with the higher probabilities. The weaker performance may be due to having fewer heavy
329 rain events in grid 2 during the study time period.

330 The ROC curves show that the multiple-predictor machine learning algorithms enhance
331 the discrimination abilities of the ensemble. The raw ensemble in both subgrids has slightly
332 positive skill in terms of AUC (Fig. 6). At the optimal PSS threshold, the raw ensemble
333 detects only 67% of heavy rain events (POD) in grid 1, and 86% of its positive forecasts are
334 false alarms (FAR). In grid 2, the detection ability is worse with a 66% POD and a 91%
335 FAR. The Bias scores greater than 1 indicate a larger proportion of false alarms than misses.

336 Since the simple logistic regression rescales the predictions over a smaller probability range,
337 it has a slightly lower AUC than the raw ensemble. Its optimal threshold is higher than the
338 raw ensemble, so it has a correspondingly lower POD, POFD, and FAR. It also has a smaller
339 bias score at the optimal threshold. The multiple logistic regression and random forest have
340 a much larger AUC and POD and similar FAR compared to the raw ensemble. The multiple
341 logistic regression and random forest have very similar scores with only slight differences in
342 the POFD, FAR, and bias. The same relationships hold in grid 2 except that there is a
343 slightly larger difference in the scores of the multiple logistic regression and random forest.
344 The overall scores are also slightly lower for grid 2, which is likely due to the lower frequency
345 of convective precipitation events in the training data.

346 2) EVALUATION BY FORECAST HOUR

347 Comparisons of BSS and AUC by hour and by sub-grid show additional trends in the
348 probability of precipitation forecasts. The raw ensemble consistently has the worst BSS (Fig.
349 7). For all models, the best performance occurs at forecast hour 1 then decreases sharply at
350 hour 2 before stabilizing. This initial decrease is likely due to the radar data assimilation
351 placing the storms in the same place initially and then having the individual storms diverge
352 from the predicted storm motions. There is a slight increase in performance between hours
353 6 and 12 for grid 2. This increase may be due to the diurnal cycle of convection resulting in
354 larger storm clusters during this time period, or it may be due to the model fully spinning
355 up. There is another major decrease in performance in grid 1 between hours 12 and 24.
356 This time period is when the greatest uncertainty exists due to convective initiation and

357 the tendency for initial convection to be isolated. The multiple logistic regression and the
358 random forest do not have any statistically significant ($\alpha < 0.05$) differences in their AUC
359 and BSS. The biggest departure from the ensemble mean in terms of AUC occurred in the
360 18 to 24 hour range for both grids, which corresponds with peak convective activity. The
361 simple logistic regression did still improve on the ensemble forecast in terms of reliability
362 but not to the same extent as the multiple-predictor models, and it made the AUC worse.
363 The temporal trends in the BSS match those found in the 2009 SSEF by Johnson and Wang
364 (2012). The multiple-predictor methods appeared to produce similar increases in BSS to the
365 single predictor neighborhood and object-based methods.

366 *c. Evaluation of Multiple Precipitation Thresholds*

367 Machine learning models were trained at 3 precipitation thresholds (2.54, 6.35, and 12.70
368 mm h⁻¹) in order to evaluate how skill varies with precipitation intensity. Fig. 8 shows the
369 random forest BSS by hour for the three precipitation thresholds. All three follow the same
370 diurnal patterns, but the 2.54 mm h⁻¹ forecasts have consistently much higher skill than
371 the other thresholds. While the 2.54 and 6.35 mm h⁻¹ show skill for all forecast hours, the
372 12.70 mm h⁻¹ forecasts do not show any skill from 14 to 26 hours. The decreasing skill with
373 threshold size is likely due to the smaller heavy precipitation areas and spatial and timing
374 errors in placement of convection in the ensembles. Similar results were found for the other
375 machine learning methods (not shown).

376 *d. Variable Importance*

377 Variable importance scores were calculated and averaged over each random forest and
378 sub-grid to determine if the random forests were choosing relevant variables and how the
379 choice of variables was affected by region. Variable importance is indicative of how random-
380 izing the values of each variable affects the random forest performance. It accounts for how
381 often a variable is used in the model, the depth of the variable in the tree, and the num-
382 ber of cases that transit through the branch containing that variable, but the importance
383 score cannot be decomposed into those factors. The top 8 variable importance scores for the
384 random forests trained on each 6-hour period in Grid 1 are shown in Table 5. In the first
385 6 hours, all of the top 5 variables are aggregations of the hour precipitation or precipitable
386 water with the rest being hour max upward wind. In this time period, the ensemble members
387 are very similar, so there is great overlap among the precipitation regions. By hours 7-12,
388 the max upward wind becomes more dominant in the rankings, although the precipitation
389 max and standard deviation are still found among the top 5 variables. Vertical velocities
390 become the most common feature in hours 13 through 18 with only precipitable water and
391 specific humidity contributing moisture information. For hours 19 through 30, the standard
392 deviation of Surface-Based CAPE appears, which is likely associated with the presence of
393 nearby boundaries.

394 The grid 2 variable importances highlight the greater importance of moisture and lower
395 importance of vertical velocities in the Eastern United States. The predicted precipitation is
396 again only important in hours 1 through 12, but the precipitable water and specific humidity
397 show high importance through the entire forecast period. Upward and Downward winds are

398 in the rankings for each time period, but they tend to be on toward the bottom of the top
399 8. Surface-Based CAPE is also important in the latter hours of the forecasts, but the mean
400 and max are selected instead of the standard deviation.

401 *e. Case Study: 13 May 2010*

402 The case of 13 May 2010 illustrates the spatial characteristics, strengths, and weaknesses
403 of the precipitation forecasts from the SSEF and the machine learning methods. Since the
404 SSEF is run through forecast hour 30, the last six forecast hours of one run overlap with the
405 first six hours of the next run. This overlap allows for the comparison of two runs on the same
406 observations and illustrates the effects of lead time on the forecast probabilities. Fig. 9 shows
407 the distribution of the observed 1-hour precipitation at 02 UTC on May 13. The bulk of the
408 precipitation originates from a broken line of discrete supercell thunderstorms positioned
409 along a stationary front stretching from northern Oklahoma into northern Missouri with
410 additional storms in Iowa and ahead of the dry line in Oklahoma and Texas. Additional
411 precipitation is falling in Virginia from a mesoscale convective system. A comparison of
412 the 2-hour and 26-hour dew point and temperature SSEF forecasts (Fig. 10) shows major
413 differences in the placement and strength of the fronts and dry line. The cold front is further
414 north and more diffuse in the 26 hour forecast while the dry line is further east, moving from
415 the central Texas panhandle to western Oklahoma. The differences in the placement of the
416 surface boundaries also affect the placement of the precipitation areas in both forecasts.

417 Comparisons of the raw ensemble probabilities and the simple logistic regression (Fig. 11)
418 show the effects of calibrating the probabilities only on precipitation. In the 2 hour forecast,

419 there is little dispersion among the ensemble members, so there is very high confidence in
420 the raw ensemble. In this case though, the high confidence areas are displaced west from
421 the actual heavy precipitation regions, resulting in some missed precipitation areas. The
422 simple logistic regression corrects for this overconfidence by lowering the SSEF probabilities.
423 The area with a greater than 10% probability is noticeably smaller and approximates the
424 area occupied by the observed rain areas. Since it rescales the probabilities, it does not
425 translate the predicted rain areas, resulting in more misses. At 26 hours, the raw ensemble
426 probabilities are lower and spread over a wider area, indicating greater ensemble dispersion.
427 The probabilities in Oklahoma and Kansas are shifted further east due to the positioning of
428 the surface boundaries. In this case, it results in the observed precipitation being captured
429 better than they were in the 2 hour forecast with fewer misses but more false alarms. The
430 simple logistic regression again reduces the probabilities and the area covered by them so
431 that it does not have significantly more false alarms than at 2 hours, but it misses some
432 precipitation that was captured by the raw ensemble. Of the calibration methods, the
433 simple logistic is best for accurately depicting areal coverage, but it overcorrected in terms
434 of downscaling the raw probabilities.

435 The results from the multiple-predictor methods show the advantages and limitations of
436 incorporating additional predictors. In Fig. 12, the multiple logistic regression is compared
437 with the random forest at 2 and 26 hours. In the 2 hour forecasts, both the multiple logistic
438 regression and random forest expand their probabilities over a wider area and shift them
439 slightly east compared to the raw ensemble, enabling them to better capture the heavy
440 precipitation regions. The random forest also captures the precipitation areas north of the
441 cold front better than the logistic regression, which may weight SBCAPE too highly and not

442 handle cold sector precipitation as well. At 26 hours, the probability areas increase and the
443 probability values are generally lower. Both algorithms capture almost all of the precipitation
444 areas, but they also produce many more false alarms in the process. Precipitation did extend
445 through central Texas into Mexico ahead of the dry line, but none of that precipitation was
446 heavy. The multiple logistic regression highlighted that precipitation as well as ahead of
447 the dry line. Unlike the 2 hour forecast, it did cover the precipitation north of the front in
448 Iowa. The random forest covered all of the precipitation areas but also had precipitation
449 probabilities along the dry line and squall line in Texas and had some 10% probabilities in
450 New Mexico where no precipitation occurred. There were also some spotty probabilities in
451 the Gulf of Mexico.

452 The extended probability areas of the random forest and multiple logistic regression
453 can be explained by examining some of the variables deemed important by the random
454 forest. In Fig. 13, the Hour Max Upward Wind (13a), SBCAPE (13b), Precipitable Water
455 (13c), and 850 mb Specific Humidity (13d) are shown. The Hour Max Upward Wind was
456 the most important variable for the 25-30 hour random forest (Table 5), and its spatial
457 distribution closely matches that of the random forest, including the positive areas in New
458 Mexico and southwest Texas. It also supports the precipitation north of the cold front in
459 Iowa. The SBCAPE was highest in central Texas in the area where both machine learning
460 models placed probabilities but no heavy rain occurred. While SBCAPE is important, there
461 are many situations in this map in which either SBCAPE is high and no precipitation
462 occurred or SBCAPE was low and precipitation occurred, especially north of the front.
463 Precipitable water and 850 mb specific humidity were also high over a wide area and are
464 generally correlated well with precipitation but occur over a much larger area than the actual

465 precipitation. Reliance on these predictors likely contributed to the false alarms. There
466 was a maximum in SBCAPE and 850 mb specific humidity in the Gulf of Mexico, which
467 possibly led to random forest probabilities there. Because of their larger areal coverage, the
468 predictors may be more useful for producing 3- or 6-hour precipitation predictions, which
469 are less dependent on the locations of individual storms.

470 5. Discussion

471 The results of the machine-learning post-processing of storm scale ensemble precipita-
472 tion forecasts displayed not only improvements to the forecasts but also limitations of the
473 ensemble, the algorithms, and grid-point-based framework. First, the post-processing model
474 performance is constrained by the information available from the ensemble. If most of the
475 ensemble members are predicting precipitation in the wrong place or not at all, and the
476 environmental conditions are also displaced, then the machine learning algorithm will not
477 be able to provide much additional skill. Second, the machine learning model will only make
478 predictions based on the range of cases it has previously seen. For higher rain thresholds, the
479 algorithms will need more independent samples in order to make skilled predictions. Third,
480 the grid-point framework does not fully account for the spatial information and error in the
481 ensemble. The spatial error could be incorporated further by smoothing the verification grid
482 with a neighborhood filter (Ebert 2009) or warping the ensemble forecast to more closely
483 match the observed precipitation (Gilleland et al. 2010) and then training. Machine learn-
484 ing algorithms could also be used to improve the calibration from object-based frameworks
485 (Johnson and Wang 2012) and could incorporate information from the additional predictors

486 within the bounds of the objects.

487 **6. Conclusions**

488 Multiple machine learning algorithms were applied to the 2010 CAPS Storm Scale Ensem-
489 ble Forecast (SSEF) system in order to improve the calibration and skill of its probabilistic
490 heavy precipitation forecasts. Two types of machine learning methods were compared over
491 a period from May 3 through June 18 with both verification statistics and a case study.
492 Verification statistics showed that all of the machine learning methods improved the cali-
493 bration of the SSEF precipitation forecasts but only the multiple-predictor methods were
494 able to calibrate the models better and discriminate more skillfully between light and heavy
495 precipitation cases. Hourly performance varied with diurnal storm cycles and the increasing
496 dispersiveness of the ensemble members. Comparisons of the rankings of predictors indi-
497 cated that the ensemble predicted precipitation was only important for the first 12 hours
498 of the model runs. After that period, the upward wind and atmospheric moisture variables
499 became better indicators of the placement of precipitation. The case study showed that the
500 multiple-predictor machine learning methods could shift the probability maxima to better
501 match the actual precipitation areas, but they would also produce more false alarm areas
502 in the process. For shorter-term forecasts, the false corrections were made without a sig-
503 nificant increase in the false alarm area. Calibrating the probabilities with only ensemble
504 rainfall predictions results in predicted areas that are too small and still displaced from
505 the observed precipitation. The multiple-predictor machine learning algorithms did prove
506 especially beneficial in that situation. Ultimately, machine learning techniques can provide

507 an enhancement to precipitation forecasts by consistently maximizing the potential of the
508 available information.

509 *Acknowledgments.*

510 Special thanks go to my Master’s committee members Fanyou Kong and Michael Rich-
511 man. Zac Flamig provided assistance with the NMQ data. CAPS SSEF forecasts were
512 supported by a grant (NWSPO-2010-201696) from the NOAA Collaborative Science, Tech-
513 nology, and Applied Research (CSTAR) Program and the forecasts were produced at the
514 National Institute for Computational Science (<http://www.nics.tennessee.edu/>). Scientists
515 at CAPS, including Fanyou Kong, Kevin Thomas, Yunheng Wang, and Keith Brewster con-
516 tributed to the design and production of the CAPS ensemble forecasts. This study was
517 funded by the NSF Graduate Research Fellowship under Grant 2011099434 and by NSF
518 grant AGS-0802888.

519

520 **REFERENCES**

- 521 Akaike, H., 1974: A new look at the statistical model identification. *IEEE Transactions on*
522 *Automatic Control*, **19**, 716–723.
- 523 Applequist, S., G. E. Gahrs, R. L. Pfeffer, and X. Niu, 2002: Comparison of methodologies
524 for probabilistic quantitative precipitation forecasting. *Wea. Forecasting*, **17**, 783–799.

- 525 Ashley, S. T. and W. S. Ashley, 2008: Flood fatalities in the united states. *J. Appl. Meteorol.*,
526 **47**, 805–818.
- 527 Breiman, L., 1984: *Classification and regression trees*. Wadsworth International Group, 358
528 pp.
- 529 Breiman, L., 2001: Random forests. *Mach. Learn.*, **45**, 5 – 32.
- 530 Bremnes, J. B., 2004: Probabilistic forecasts of precipitation in terms of quantiles using
531 NWP model output. *Mon. Wea. Rev.*, **132**, 338–347.
- 532 Brier, G. W., 1950: Verification of forecasts expressed in terms of probability. *Mon. Wea.*
533 *Rev.*, **78**, 1–3.
- 534 Clark, A. J. and Coauthors, 2011: Probabilistic precipitation forecast skill as a function of
535 ensemble size and spatial scale in a convection-allowing ensemble. *Mon. Wea. Rev.*, **139**,
536 1410–1418.
- 537 Clark, A. J., W. A. Gallus, M. Xue, and F. Kong, 2009: A comparison of precipitation
538 forecast skill between small convection-allowing and large convection-parameterizing en-
539 sembles. *Wea. Forecasting*, **24**, 1121–1140.
- 540 Clark, A. J., S. J. Weiss, J. S. Kain, and Coauthors, 2012: An overview of the 2010 hazardous
541 weather testbed experimental forecast program spring experiment. *Bull. Amer. Meteor.*
542 *Soc.*, **93**, 55–74.
- 543 Doswell, C. A., H. E. Brooks, and R. A. Maddox, 1996: Flash flood forecasting: An
544 ingredients-based methodology. *Weather and Forecasting*, **11**, 560–581.

545 Doswell, C. A., R. Davies-Jones, and D. L. Keller, 1990: On summary measures of skill in
546 rare event forecasting based on contingency tables. *Wea. Forecasting*, **5**, 576–585.

547 Ebert, E., 2001: Ability of a poor man’s ensemble to predict the probability and distribution
548 of precipitation. *Mon. Wea. Rev.*, **129**, 2461–2480.

549 Ebert, E. E., 2009: Neighborhood verification: A strategy for rewarding close forecasts.
550 *Weather and Forecasting*, **24 (6)**, 1498–1510.

551 Eckel, F. A. and M. K. Walters, 1998: Calibrated probabilistic quantitative precipitation
552 forecasts based on the MRF ensemble. *Wea. Forecasting*, **13**, 1132–1147.

553 Gagne II, D. J., A. McGovern, and J. Brotzge, 2009: Classification of convective areas using
554 decision trees. *Journal of Atmospheric and Oceanic Technology*, **26 (7)**, 1341–1353.

555 Gagne II, D. J., A. McGovern, and M. Xue, 2012: Machine learning enhancement of storm
556 scale ensemble precipitation forecasts. *Proceedings of the Conference on Intelligent Data
557 Understanding (CIDU)*, Boulder, CO, IEEE-CIS, 39–46.

558 Gilleland, E., et al., 2010: Spatial forecast verification: image warping. Tech. Rep.
559 NCAR/TN-482+STR, NCAR.

560 Glahn, H. R. and D. A. Lowry, 1972: The use of model output statistics (MOS) in objective
561 weather forecasts. *J. Appl. Meteor.*, **11**, 1203–1211.

562 Hall, T., H. E. Brooks, and C. A. Doswell, 1999: Precipitation forecasting using a neural
563 network. *Wea. Forecasting*, **14**, 338–345.

- 564 Hamill, T. M. and S. J. Colucci, 1997: Verification of Eta-RSM short-range ensemble fore-
565 casts. *Mon. Wea. Rev.*, **125**, 1312–1327.
- 566 Hamill, T. M. and S. J. Colucci, 1998: Evaluation of Eta-RSM ensemble probabilistic pre-
567 cipitation forecasts. *Mon. Wea. Rev.*, **126**, 711–724.
- 568 Hamill, T. M., R. Hagedorn, and J. S. Whitaker, 2008: Probabilistic forecast calibration
569 using ECMWF and GFS ensemble reforecasts. Part II: Precipitation. *Mon. Wea. Rev.*,
570 **136**, 2620–2632.
- 571 Hamill, T. M., J. S. Whitaker, and X. Wei, 2004: Ensemble reforecasting: Improving
572 medium-range forecast skill using retrospective forecasts. *Mon. Wea. Rev.*, **132**, 1434–
573 1447.
- 574 Hansen, A. W. and W. J. A. Kuipers, 1965: On the relationship between the frequency of
575 rain and various meteorological parameters. *Meded. Verhand.*, **81**, 2–15.
- 576 James, G., D. Witten, T. Hastie, and R. Tibshirani, 2013: *An Introduction to Statistical*
577 *Learning with Applications in R*. Springer, 430 pp.
- 578 Johnson, A. and X. Wang, 2012: Verification and calibration of neighborhood and object-
579 based probabilistic precipitation forecasts from a multimodel convection-allowing ensem-
580 ble. *Mon. Wea. Rev.*, **140**, 3054–3077.
- 581 Koizumi, K., 1999: An objective method to modify numerical model forecasts with newly
582 given weather data using an artificial neural network. *Wea. Forecasting*, **14**, 109–118.
- 583 Kong, F., et al., 2011: Evaluation of CAPS multi-model storm-scale ensemble forecast for

584 the NOAA HWT 2010 spring experiment. *24th Conf. Wea. Forecasting/20th Conf. Num.*
585 *Wea. Pred.*, Seattle, WA, Amer. Meteor. Soc., Paper 457.

586 Krishnamurti, T. N., C. M. Kishtawal, T. E. LaRow, D. R. Bachiochi, Z. Zhang, C. E.
587 Williford, S. Gadgil, and S. Surendran, 1999: Improved weather and seasonal climate
588 forecasts from multimodel superensemble. *Science*, **285**, 1548–1550.

589 Kusiak, A. and A. Verma, 2011: Prediction of status patterns of wind turbines: A data-
590 mining approach. *J. Sol. Energy Eng.*, **133**, doi:10.1115/1.4003188.

591 Manzato, A., 2007: A note on the maximum Peirce skill score. *Wea. Forecasting*, **22**, 1148–
592 1154.

593 Marsh, P. T., J. S. Kain, V. Lakshmanan, A. J. Clark, N. Hitchens, and J. Hardy, 2012:
594 A method for calibrating deterministic forecasts of rare events. *Wea. Forecasting*, **27**,
595 531–538.

596 Mason, I., 1982: A model for assessment of weather forecasts. *Aust. Meteor. Mag.*, **30**,
597 291–303.

598 Molteni, F., R. Buizza, T. N. Palmer, and T. Petroliagis, 1996: The ECMWF ensemble
599 prediction system: Methodology and validation. *Quart. J. Roy. Meteor. Soc.*, **122**, 73–
600 119.

601 Murphy, A. H., 1973: A new vector partition of the probability score. *J. Appl. Meteor.*, **12**,
602 595–600.

- 603 Murphy, A. H., 1977: The value of climatological, categorical, and probabilistic forecasts in
604 the cost-loss ratio situation. *Mon. Wea. Rev.*, **105**, 803–816.
- 605 Peirce, C. S., 1884: The numerical measure of the success of predictions. *Science*, **4**, 453–454.
- 606 Stensrud, D. J., H. E. Brooks, J. Du, S. Tracton, and E. Rogers, 1999: Using ensembles for
607 short-range forecasting. *Mon. Wea. Rev.*, **127**, 433–446.
- 608 Strobl, C., A.-L. Boulesteix, T. Kneib, T. Augustin, and A. Zeileis, 2008: Conditional
609 variable importance for random forests. *BMC Bioinformatics*, **9** (1), 307, doi:10.1186/
610 1471-2105-9-307, URL <http://www.biomedcentral.com/1471-2105/9/307>.
- 611 Toth, Z. and E. Kalnay, 1993: Ensemble forecasting at NMC: The generation of perturba-
612 tions. *Bull. Amer. Meteor. Soc.*, **74**, 2317–2330.
- 613 Tracton, M. S. and E. Kalnay, 1993: Operational ensemble prediction at the National Me-
614 teorological Center: Practical aspects. *Wea. Forecasting*, **8**, 379–398.
- 615 Vasiloff, S., et al., 2007: Improving QPE and very short term QPF: An initiative for a
616 community-wide integrated approach. *Bull. Amer. Meteor. Soc.*, **88**, 1899–1911.
- 617 Wilks, D. S., 2011: *Statistical Methods in the Atmospheric Sciences*. 3d ed., Academic Press,
618 676 pp.
- 619 Williams, J. K., 2013: Using random forests to diagnose aviation turbulence. *Mach. Learn.*,
620 doi:10.1007/s10994-013-5346-7.
- 621 Xue, M., K. K. Droegemeier, and V. Wong, 2000: The Advanced Regional Prediction System

622 (ARPS) - A multiscale nonhydrostatic atmospheric simulation and prediction model. Part
623 I: Model dynamics and verification. *Meteor. Atmos. Phys.*, **75**, 161–193.

624 Xue, M., D. Wang, J. Gao, K. Brewster, and K. K. Droegemeier, 2003: The Advanced
625 Regional Prediction System (ARPS), storm-scale numerical weather prediction and data
626 assimilation. *Meteor. Atmos. Phys.*, **82**, 139–170.

627 Xue, M., et al., 2001: The Advanced Regional Prediction System (ARPS) - A multiscale
628 nonhydrostatic atmospheric simulation and prediction model. Part II: Model physics and
629 applications. *Meteor. Atmos. Phys.*, **76**, 134–165.

630 Xue, M., et al., 2011: CAPS realtime storm scale ensemble and high resolution forecasts for
631 the NOAA hazardous weather testbed 2010 spring experiment. *24th Conf. Wea. Forecast-*
632 *ing/20th Conf. Num. Wea. Pred.*, Seattle, WA, Amer. Meteor. Soc., PAPER 9A.2.

633 Yuan, H., X. Gao, S. L. Mullen, S. Sorooshian, J. Du, and H. H. Juang, 2007: Calibration
634 of probabilistic quantitative precipitation forecasts with an artificial neural network. *Wea.*
635 *Forecasting*, **22**, 1287–1303.

636 Zhang, J., Y. Qi, K. Howard, C. Langston, and B. Kaney, 2011: Radar quality index (RQI)—
637 a combined measure of beam blockage and vpr effects in a national network. *Proc. Eighth*
638 *Int. Symp. on Weather Radar and Hydrology*, Exeter, United Kingdom, Royal Meteorolo-
639 gical Society.

640 List of Tables

641	1	Frequencies of total and sampled grid points for the different precipitation	
642		thresholds.	33
643	2	The names and descriptions of model predictors sampled from the SSEF runs.	
644		CAPE is Convective Available Potential Energy, and CIN is Convective Inhi-	
645		bition.	34
646	3	An example binary contingency table for whether or not rain is forecast.	35
647	4	The scores calculated from the binary contingency table for use with the	
648		optimal threshold predictions.	36
649	5	Top 8 variable importance scores for the random forests trained over 6-hourly	
650		periods in Grid 1.	37
651	6	Top 8 variable importance scores for the random forests trained over 6-hourly	
652		periods in Grid 2.	38

Table 1: Frequencies of total and sampled grid points for the different precipitation thresholds.

Grid	0-0.25 mm	0.25-6 mm	6 mm
Grid 1 Total	218,052,271	13,062,623	1,884,725
Grid 1 Sampled	86,704	51,802	7,490
Grid 2 Total	194,343,572	10,923,189	1,137,694
Grid 2 Sampled	77,210	4,3218	4,497

Table 2: The names and descriptions of model predictors sampled from the SSEF runs. CAPE is Convective Available Potential Energy, and CIN is Convective Inhibition.

Variable	Levels
Hour Precipitation	Surface
Surface-Based CAPE	Surface
Surface-Based CIN	Surface
Dewpoint	2 m
Pressure	Mean Sea Level
Composite Radar Reflectivity	Column Maximum
Precipitable Water	Column Sum
Height	700 mb
U-Wind	700 mb
V-Wind	700 mb, 500 mb
Specific Humidity	850 mb, 700 mb, 500 mb
Temperature	2 m, 850 mb, 700 mb, 500 mb
Hour Max Reflectivity	Column and Time Maximum
Hour Max Upward Wind	Column and Time Maximum
Hour Max Downward Wind	Column and Time Maximum

Table 3: An example binary contingency table for whether or not rain is forecast.

		Observed	
		Yes	No
Forecast	Yes	a (hit)	b (false alarm)
	No	c (miss)	d (true negative)

Table 4: The scores calculated from the binary contingency table for use with the optimal threshold predictions.

Score	Formula
POD	$\frac{a}{a+c}$
FAR	$\frac{b}{a+b}$
POFD	$\frac{b}{b+d}$
ETS	$\frac{a - a_{random}}{a + b + c - a_{random}}$
a_{random}	$\frac{(a+c)(a+b)}{a+b+c+d}$
PSS	$\frac{a}{a+c} - \frac{b}{b+d}$
Bias	$\frac{a+b}{a+c}$

Table 5: Top 8 variable importance scores for the random forests trained over 6-hourly periods in Grid 1.

Variable	Mean	Variable	Mean
Grid 1 Hours 1-6		Grid 1 Hours 7-12	
Hour Precipitation SD	0.0118	Hour Max Upward Wind Mean	0.0129
Precipitable Water Mean	0.0107	Hour Max Upward Wind Max	0.0126
Hour Precipitation Mean	0.0105	Hour Max Upward Wind SD	0.0113
Precipitable Water Max	0.0103	Hour Precipitation SD	0.0105
Hour Precipitation Max	0.0096	Hour Precipitation Max	0.0088
Hour Max Upward Wind Mean	0.0088	Hour Max Downward Wind SD	0.0087
Hour Max Upward Wind Max	0.0088	Hour Max Downward Wind Min	0.0084
Hour Max Upward Wind SD	0.0081	Specific Humidity 700 mb Max	0.0076
Grid 1 Hours 13-18		Grid 1 Hours 19-24	
Hour Max Upward Wind Mean	0.0096	Hour Max Downward Wind Min	0.0114
Hour Max Downward Wind SD	0.0088	Surface-Based CAPE SD	0.0112
Hour Max Downward Wind Mean	0.0086	Hour Max Downward Wind Mean	0.0110
Hour Max Downward Wind Min	0.0084	Hour Max Downward Wind SD	0.0105
Hour Max Upward Wind SD	0.0071	Specific Humidity 850 mb Mean	0.0092
Precipitable Water Max	0.0070	Temperature 700 mb Min	0.0089
Specific Humidity 700 mb Max	0.0069	Hour Max Upward Wind Mean	0.0088
Hour Max Upward Wind Max	0.0062	Hour Max Upward Wind Max	0.0084
Grid 1 Hours 25-30			
Hour Max Upward Wind Mean	0.0123		
Hour Max Downward Wind Mean	0.0119		
Hour Max Upward Wind SD	0.0112		
Hour Max Upward Wind Max	0.0108		
Hour Max Downward Wind SD	0.0105		
Hour Max Downward Wind Min	0.0104		
Specific Humidity 850 mb Mean	0.0078		
Surface-Based CAPE SD	0.0077		

Table 6: Top 8 variable importance scores for the random forests trained over 6-hourly periods in Grid 2.

Variable	Mean	Variable	Mean
Grid 2 Hours 1-6		Grid 2 Hours 7-12	
Hour Precipitation Mean	0.0116	Precipitable Water Mean	0.0111
Precipitable Water Mean	0.0114	Precipitable Water Min	0.0085
Precipitable Water Min	0.0104	Precipitable Water Max	0.0079
Precipitable Water Max	0.0097	Hour Max Upward Wind Mean	0.0077
Hour Max Upward Wind Mean	0.0094	Hour Max Upward Wind SD	0.0066
Hour Max Reflectivity Mean	0.0091	Hour Precipitation Mean	0.0065
Hour Max Upward Wind Max	0.0082	Hour Precipitation SD	0.0062
Hour Precipitation Max	0.0080	Hour Max Upward Wind Max	0.0062
Grid 2 Hours 13-18		Grid 2 Hours 19-24	
Precipitable Water Mean	0.0077	Specific Humidity 850 mb Mean	0.0162
Hour Max Upward Wind Mean	0.0075	Specific Humidity 850 mb Max	0.0128
Specific Humidity 850 mb Mean	0.0065	Hour Max Upward Wind Mean	0.0110
Precipitable Water Min	0.0065	Surface-Based CAPE Max	0.0106
Precipitable Water Max	0.0056	Surface-Based CAPE Mean	0.0103
Temperature 700 mb Mean	0.0056	Hour Max Upward Wind SD	0.0101
Hour Max Upward Wind Max	0.0055	Hour Max Downward Wind Mean	0.0098
Temperature 500 mb Min	0.0054	Hour Max Upward Wind Max	0.0090
Grid 2 Hours 25-30			
Specific Humidity 850 mb Mean	0.0114		
Hour Max Upward Wind Mean	0.0094		
Precipitable Water Mean	0.0084		
Specific Humidity 850 mb Max	0.0083		
Hour Precipitation SD	0.0072		
Hour Max Upward Wind Max	0.0071		
Hour Max Downward Wind Mean	0.0070		
Surface-Based CAPE Max	0.0069		

653 List of Figures

- 654 1 Map of the number of grid points sampled within a 400 km² box spatially
655 averaged with a 300 km radius median filter. The domain sub grids are also
656 shown and labeled. 41
- 657 2 Histogram comparing the relative frequencies of the full precipitation distri-
658 bution for the each subgrid to the sampled rainfall distributions. 42
- 659 3 For a precipitation threshold of 0.25 mm h⁻¹, the optimal probability thresh-
660 olds (OT) and binary verification statistics calculated at the OT for the SSEF
661 and each machine learning model in grid 1. 43
- 662 4 For a precipitation threshold of 0.25 mm h⁻¹, the optimal probability thresh-
663 olds (OT) and binary verification statistics calculated at the OT for the SSEF
664 and each machine learning model in grid 2. 44
- 665 5 Attributes diagrams for the raw ensemble and each machine learning model.
666 The solid line with circles is the reliability curve, which indicates the observed
667 relative frequency of rain events above the 6.35 mm threshold for each prob-
668 ability bin. The dot-dashed line with triangles shows the relative frequency
669 of all forecasts that fall in each probability bin. If a point on the reliability
670 curve falls in the gray area, it contributes positively to the BSS. The horizon-
671 tal and vertical dashed lines are located at the climatological frequency for a
672 particular sub grid. The diagonal dashed line indicates perfect reliability. 45
- 673 6 Relative Operating Characteristic (ROC) curves for the raw ensemble and
674 each machine learning model. The diagonal lines indicate PSS. 46

675	7	BSS and AUC comparisons by hour.	47
676	8	Evaluation of random forests trained with 3 precipitation thresholds at each	
677		forecast hour. Brier Skill Score (BSS) is the evaluation metric.	48
678	9	Observed 1-hour precipitation on 13 May 2010 at 02 UTC.	49
679	10	Filled contours of 2 m dew point (DEWP2M) and 2 m temperature (TEMP2M)	
680		for the 2-hour and 26-hour SSEF forecasts valid on 13 May 2010 at 0200 UTC.	50
681	11	The 2-hour (left) and 26-hour (right) forecasts of the SSEF ensemble prob-	
682		ability (top) and the simple logistic regression (bottom). The green areas	
683		indicate the observed areas of 1-hour precipitation greater than 6.35 mm h^{-1} .	51
684	12	The 2-hour (left) and 26-hour (right) forecasts of the multiple logistic regres-	
685		sion (top) and the random forest (bottom). The green areas indicate the	
686		observed areas of 1-hour precipitation greater than 6.35 mm h^{-1} .	52
687	13	Maps of the SSEF 26-hour forecasts of predictors considered important by the	
688		random forest model.	53

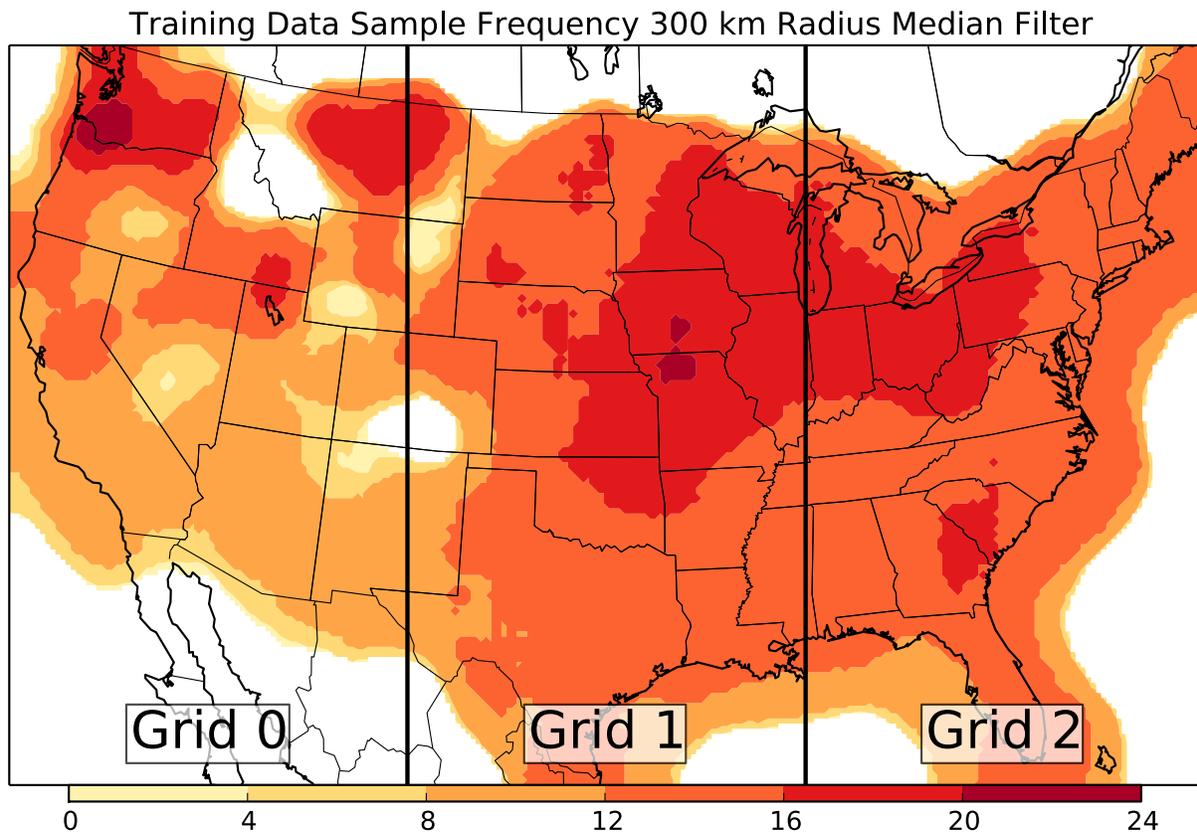


Figure 1: Map of the number of grid points sampled within a 400 km^2 box spatially averaged with a 300 km radius median filter. The domain sub grids are also shown and labeled.

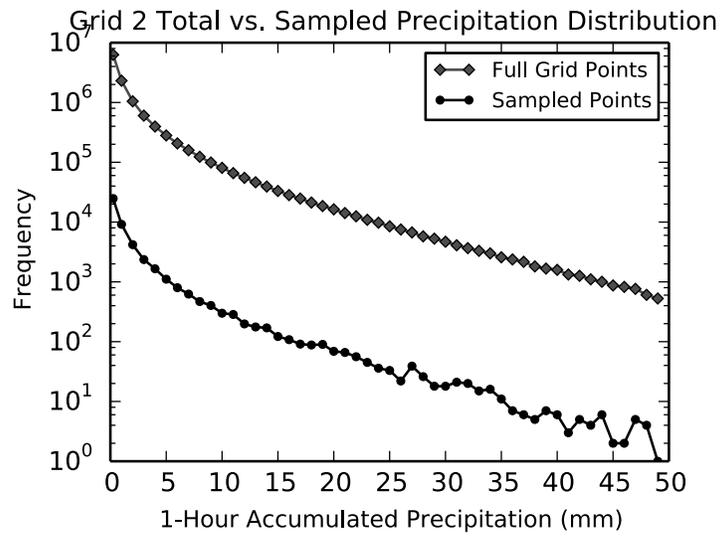
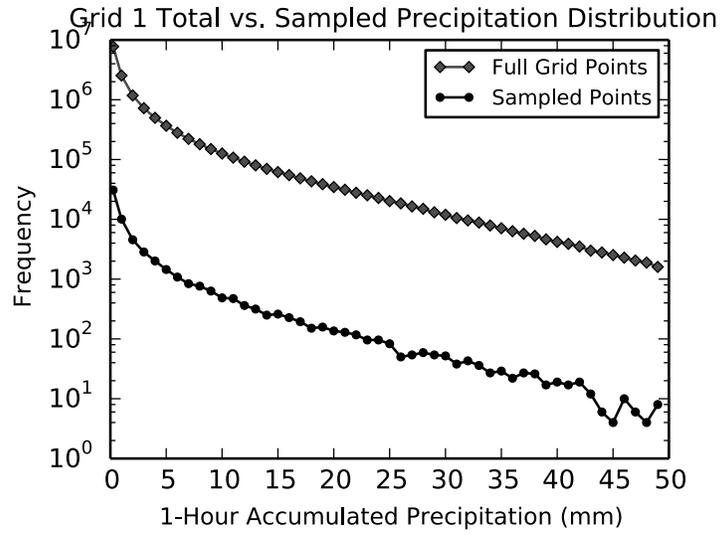


Figure 2: Histogram comparing the relative frequencies of the full precipitation distribution for the each subgrid to the sampled rainfall distributions.

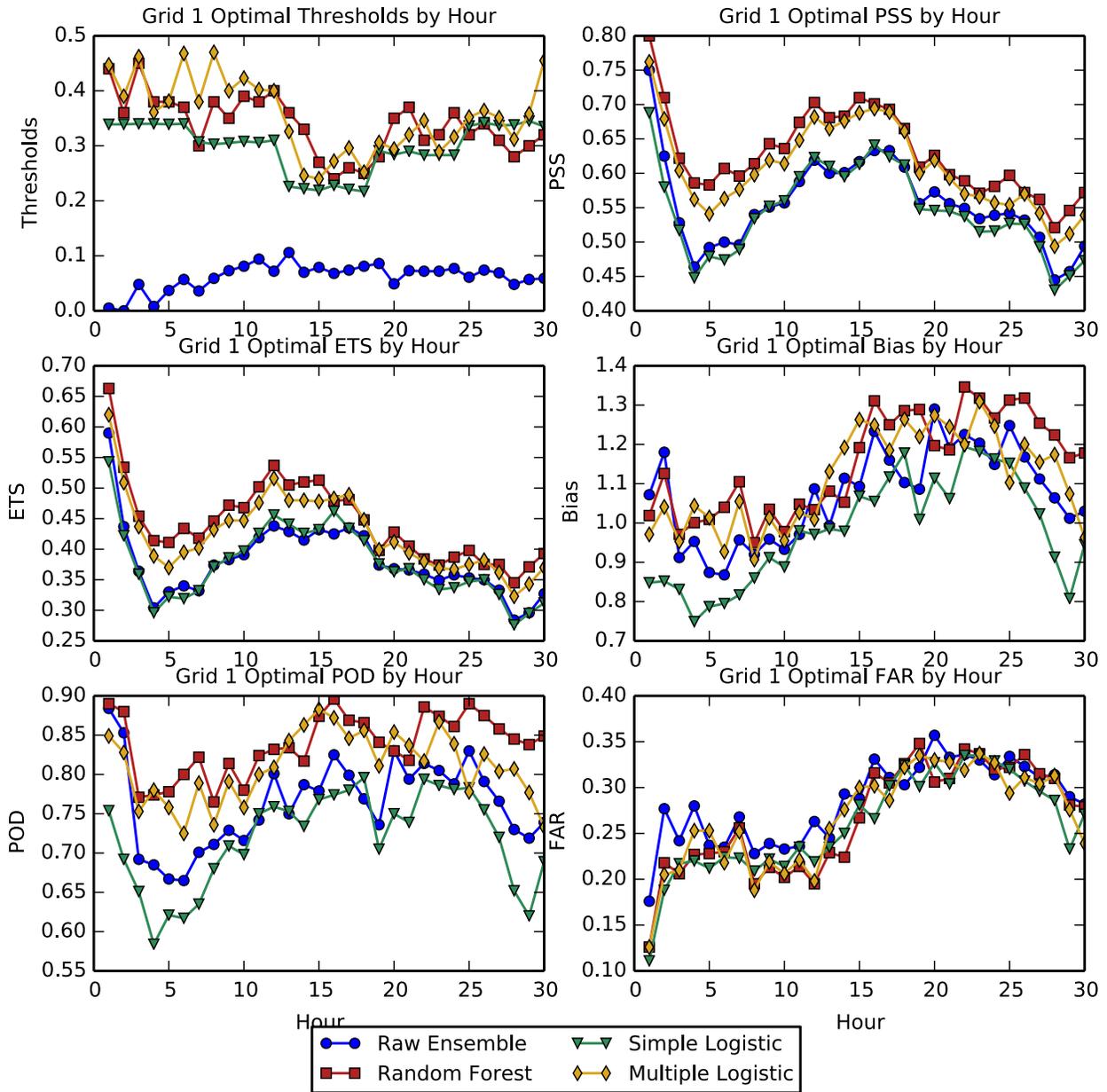


Figure 3: For a precipitation threshold of 0.25 mm h^{-1} , the optimal probability thresholds (OT) and binary verification statistics calculated at the OT for the SSEF and each machine learning model in grid 1.

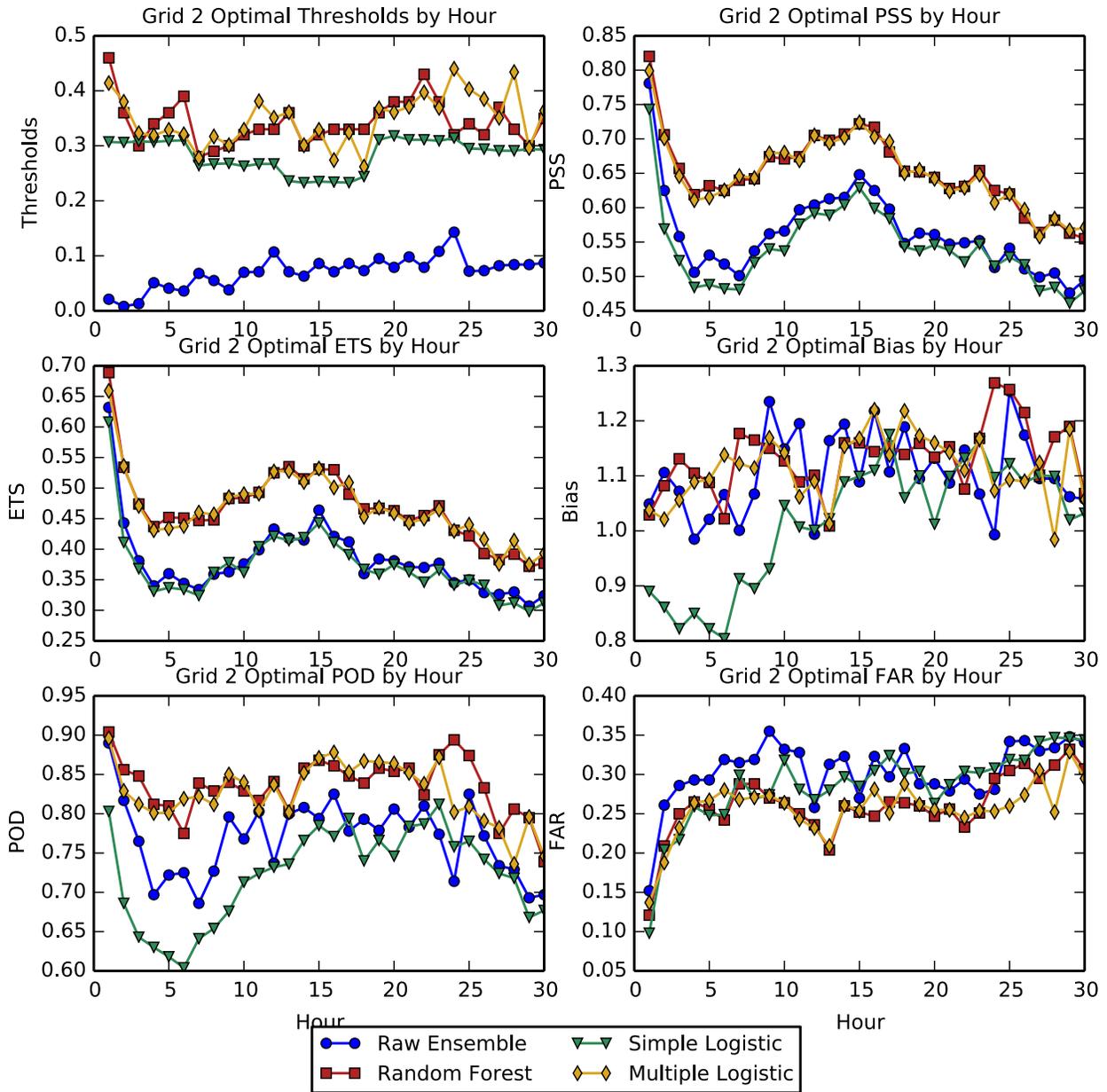


Figure 4: For a precipitation threshold of 0.25 mm h^{-1} , the optimal probability thresholds (OT) and binary verification statistics calculated at the OT for the SSEF and each machine learning model in grid 2.

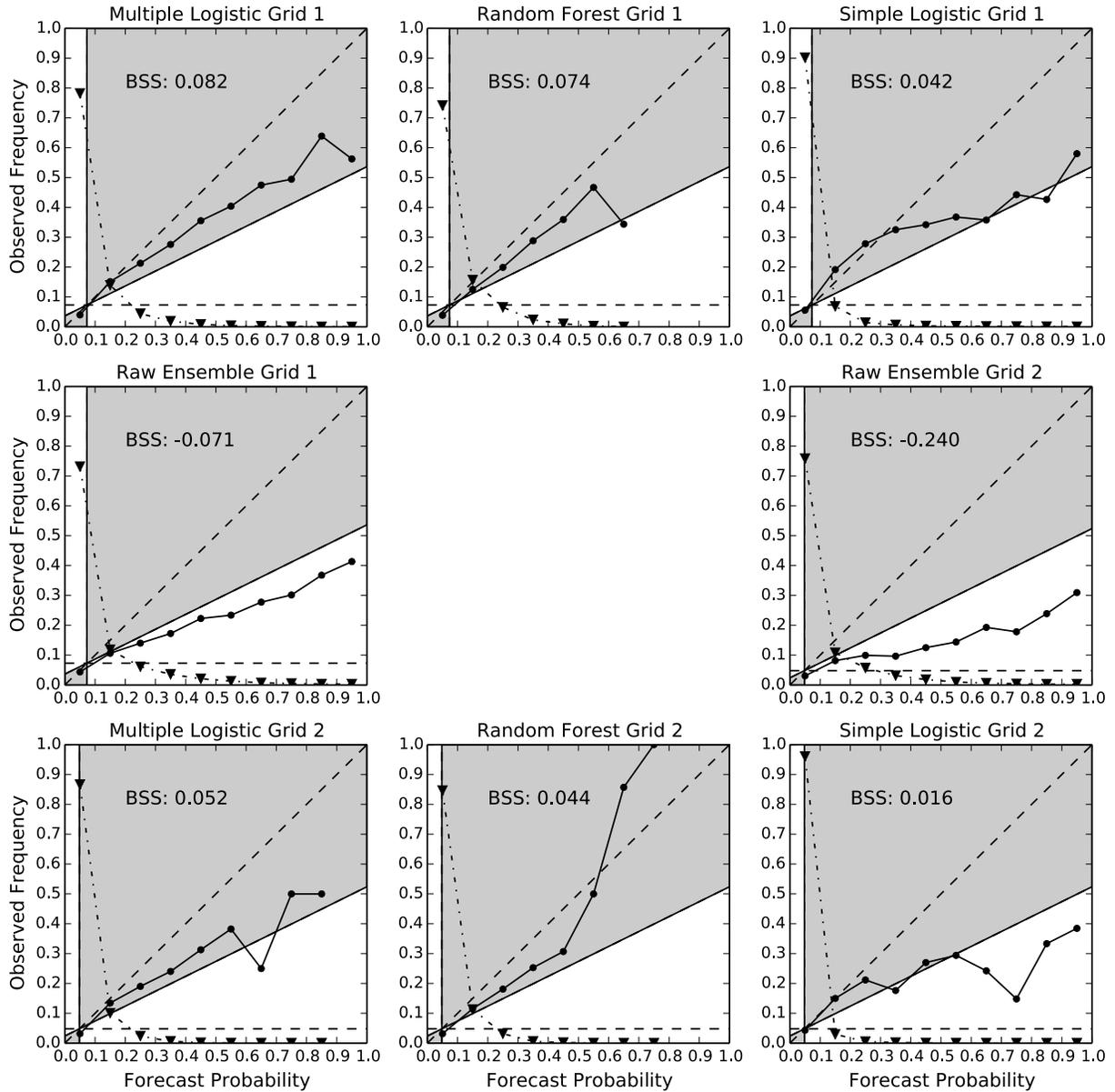


Figure 5: Attributes diagrams for the raw ensemble and each machine learning model. The solid line with circles is the reliability curve, which indicates the observed relative frequency of rain events above the 6.35 mm threshold for each probability bin. The dot-dashed line with triangles shows the relative frequency of all forecasts that fall in each probability bin. If a point on the reliability curve falls in the gray area, it contributes positively to the BSS. The horizontal and vertical dashed lines are located at the climatological frequency for a particular sub grid. The diagonal dashed line indicates perfect reliability.

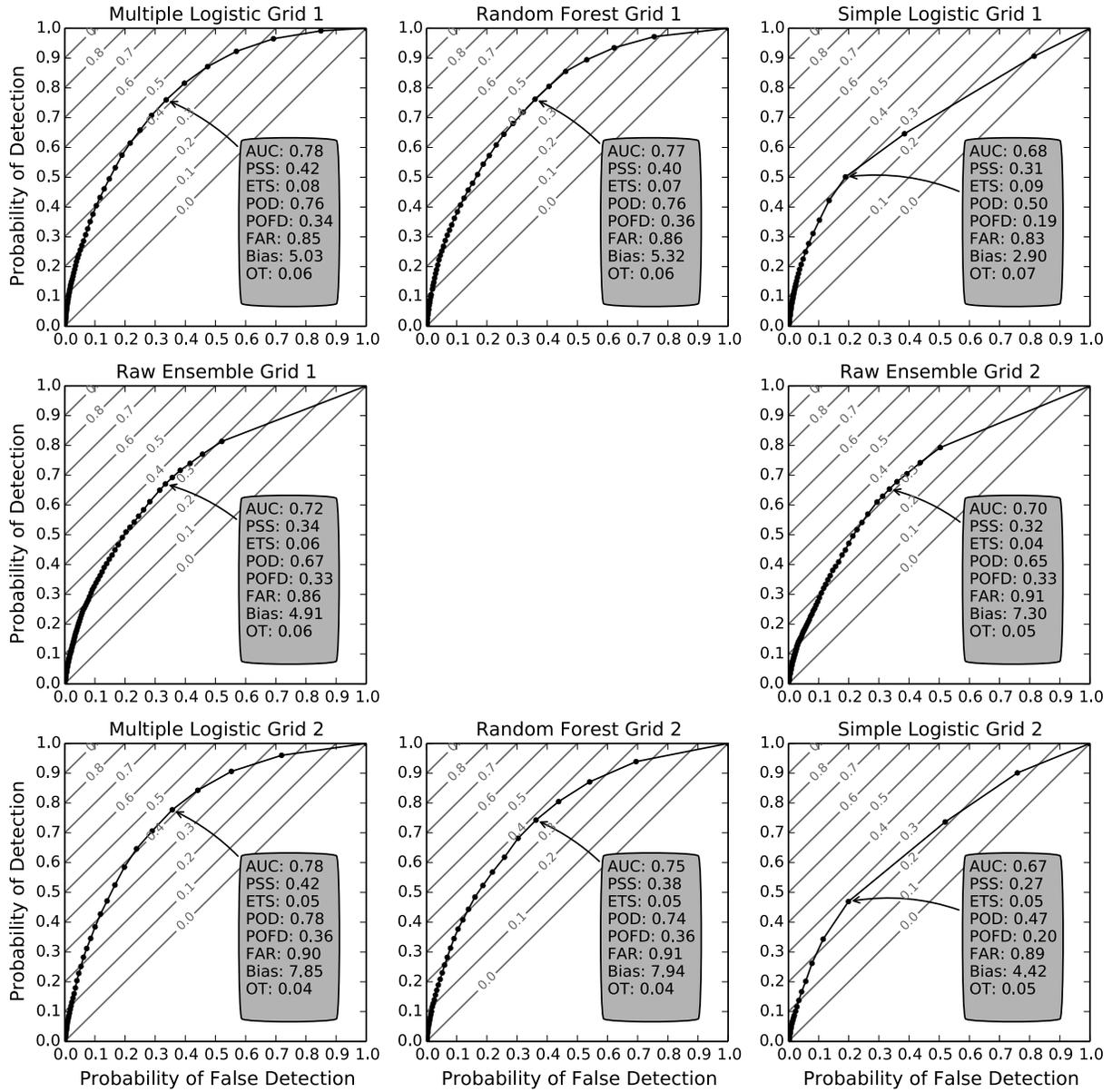


Figure 6: Relative Operating Characteristic (ROC) curves for the raw ensemble and each machine learning model. The diagonal lines indicate PSS.

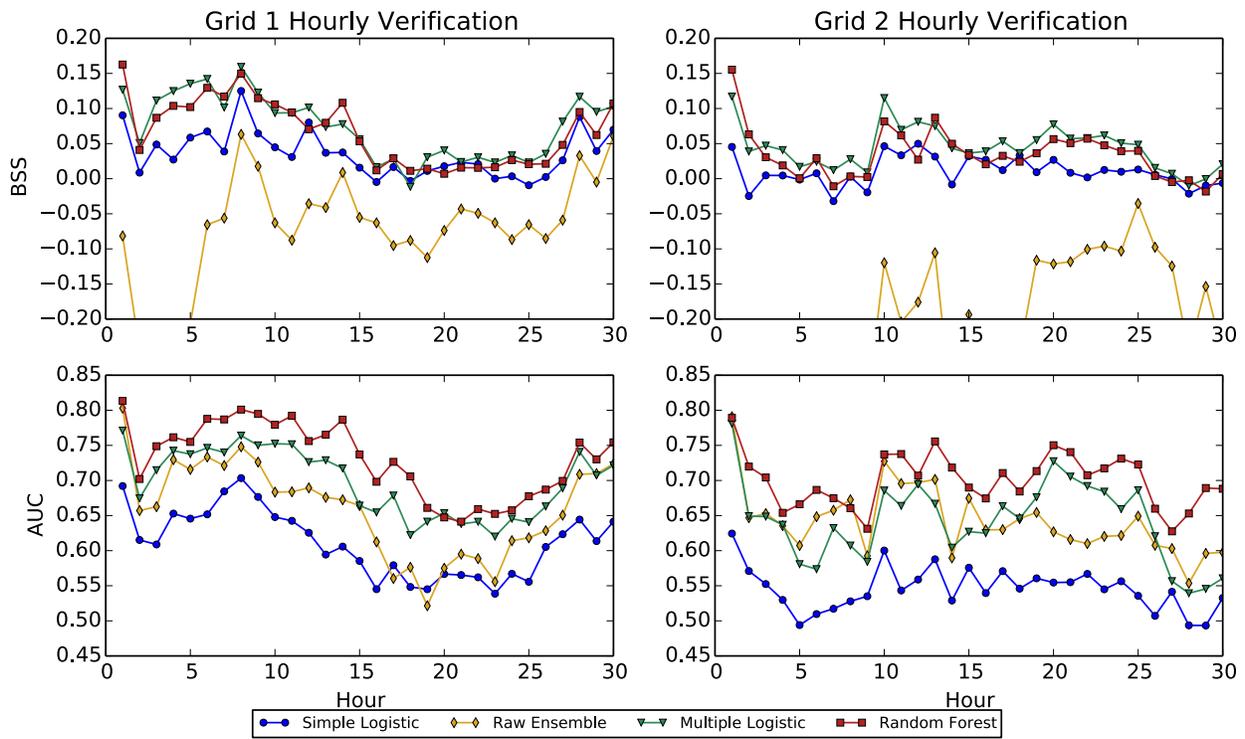


Figure 7: BSS and AUC comparisons by hour.

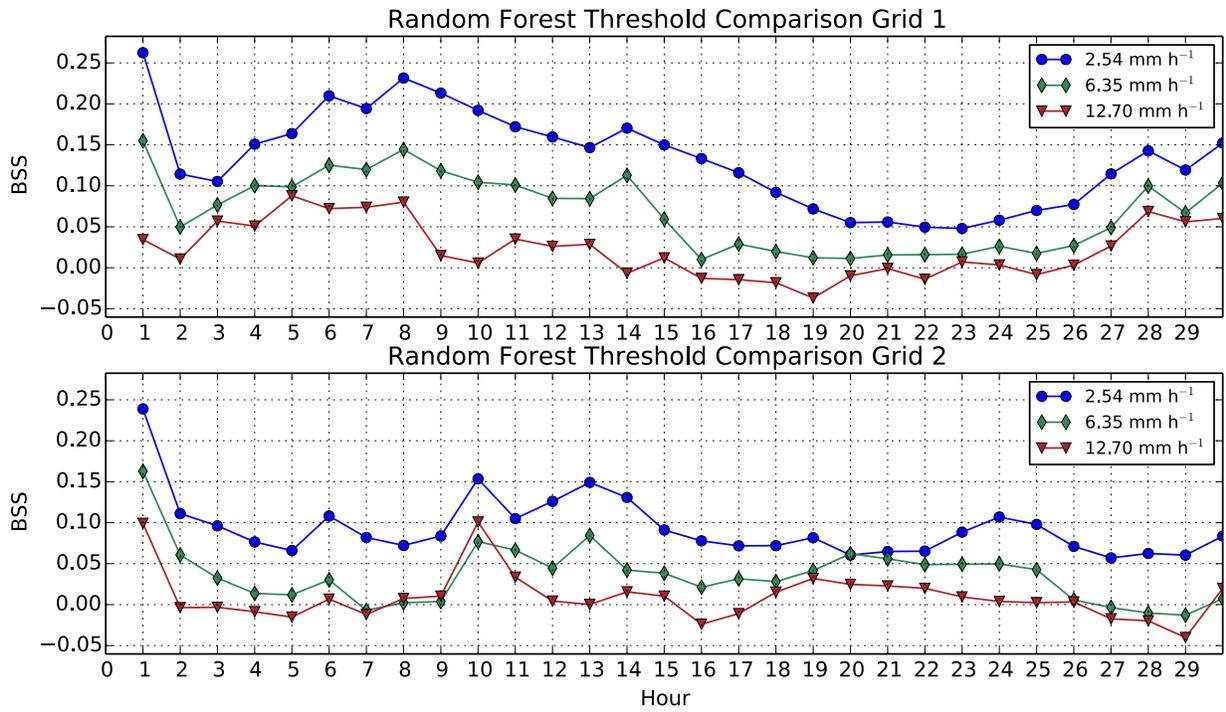


Figure 8: Evaluation of random forests trained with 3 precipitation thresholds at each forecast hour. Brier Skill Score (BSS) is the evaluation metric.

Observed Precipitation Valid at 13 May 2010 0200 UTC

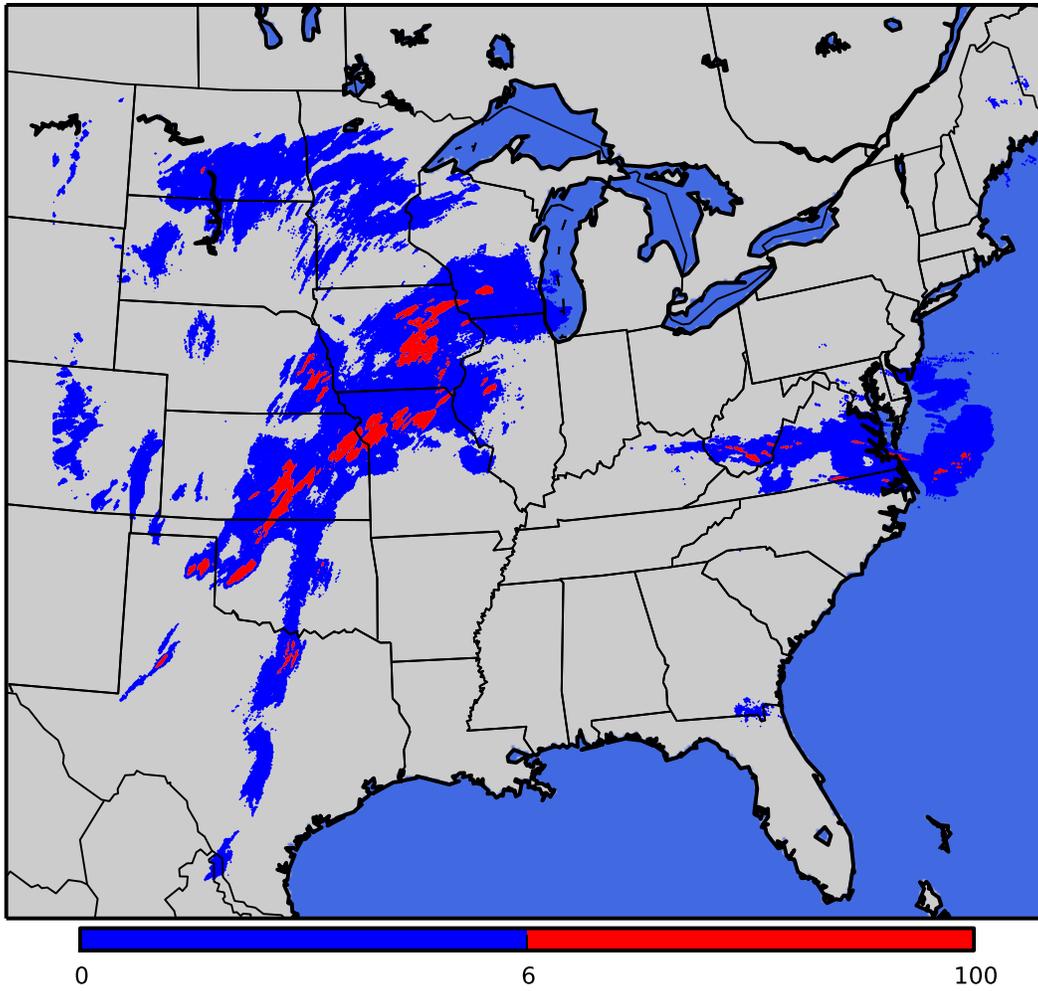
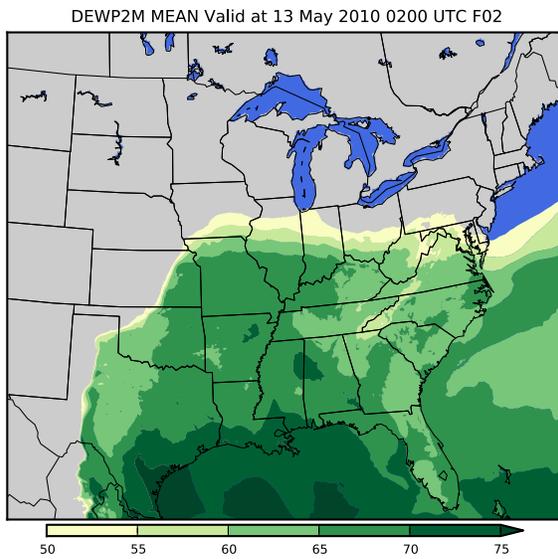
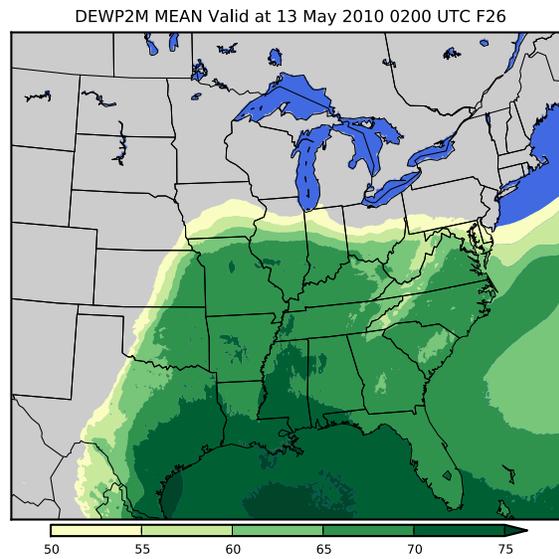


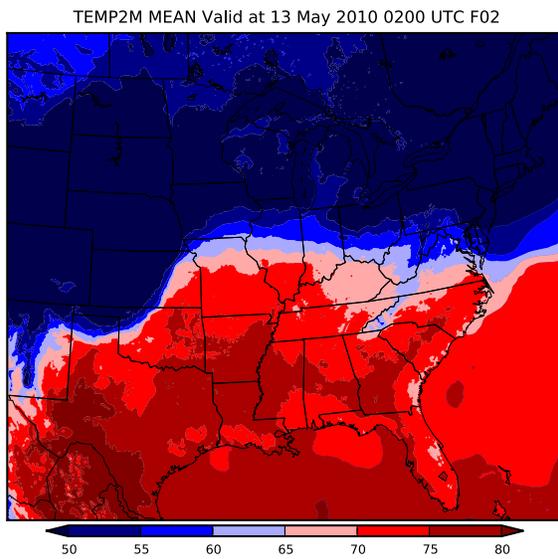
Figure 9: Observed 1-hour precipitation on 13 May 2010 at 02 UTC.



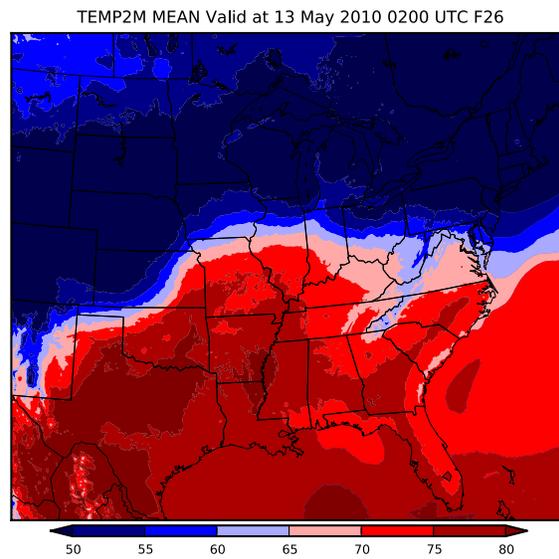
(a) 2 m Dew point (°F)



(b) 2 m Dew point (°F)



(c) 2 m Temperature (°F)



(d) 2 m Temperature (°F)

Figure 10: Filled contours of 2 m dew point (DEWP2M) and 2 m temperature (TEMP2M) for the 2-hour and 26-hour SSEF forecasts valid on 13 May 2010 at 0200 UTC.

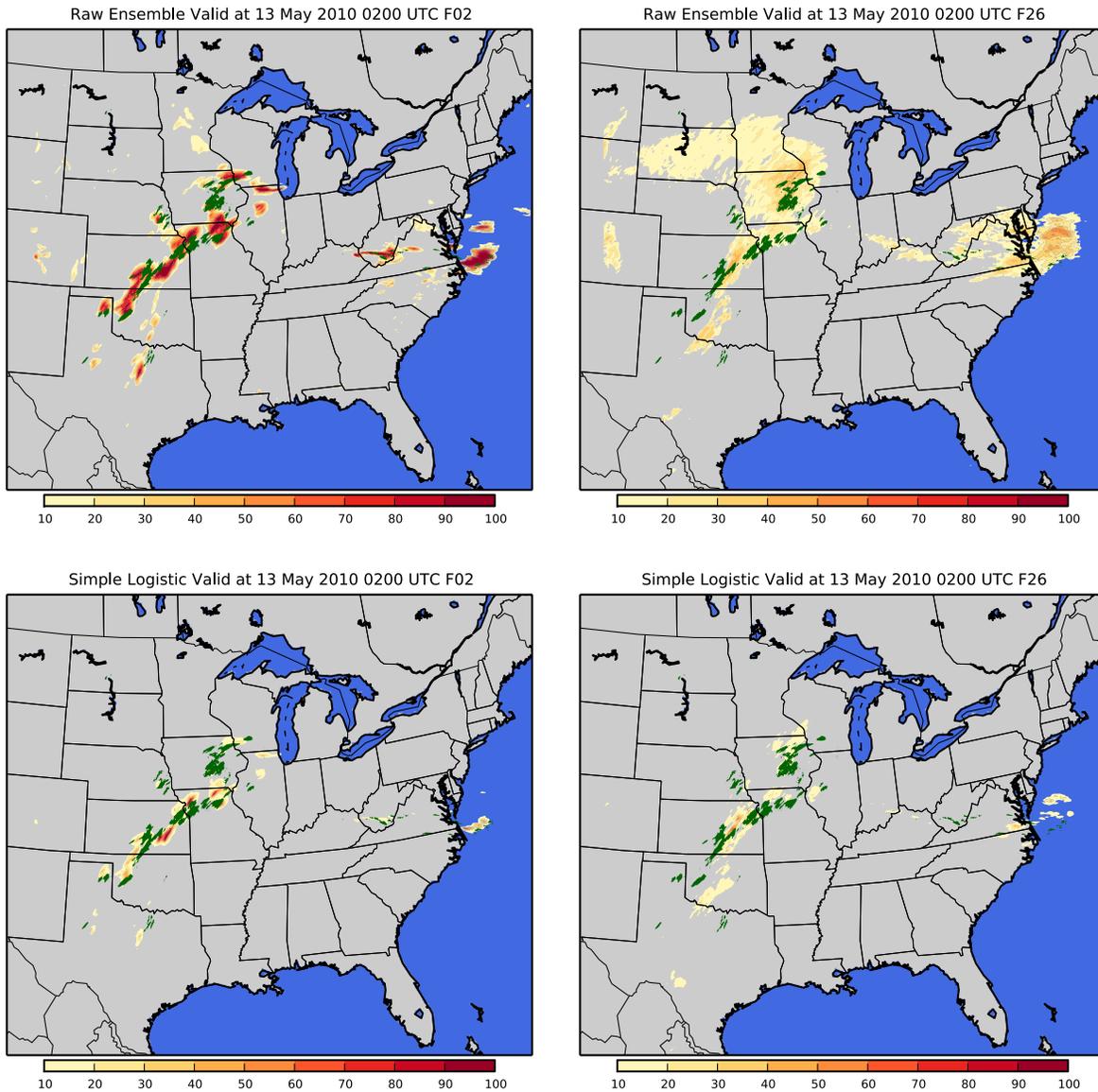


Figure 11: The 2-hour (left) and 26-hour (right) forecasts of the SSEF ensemble probability (top) and the simple logistic regression (bottom). The green areas indicate the observed areas of 1-hour precipitation greater than 6.35 mm h^{-1} .

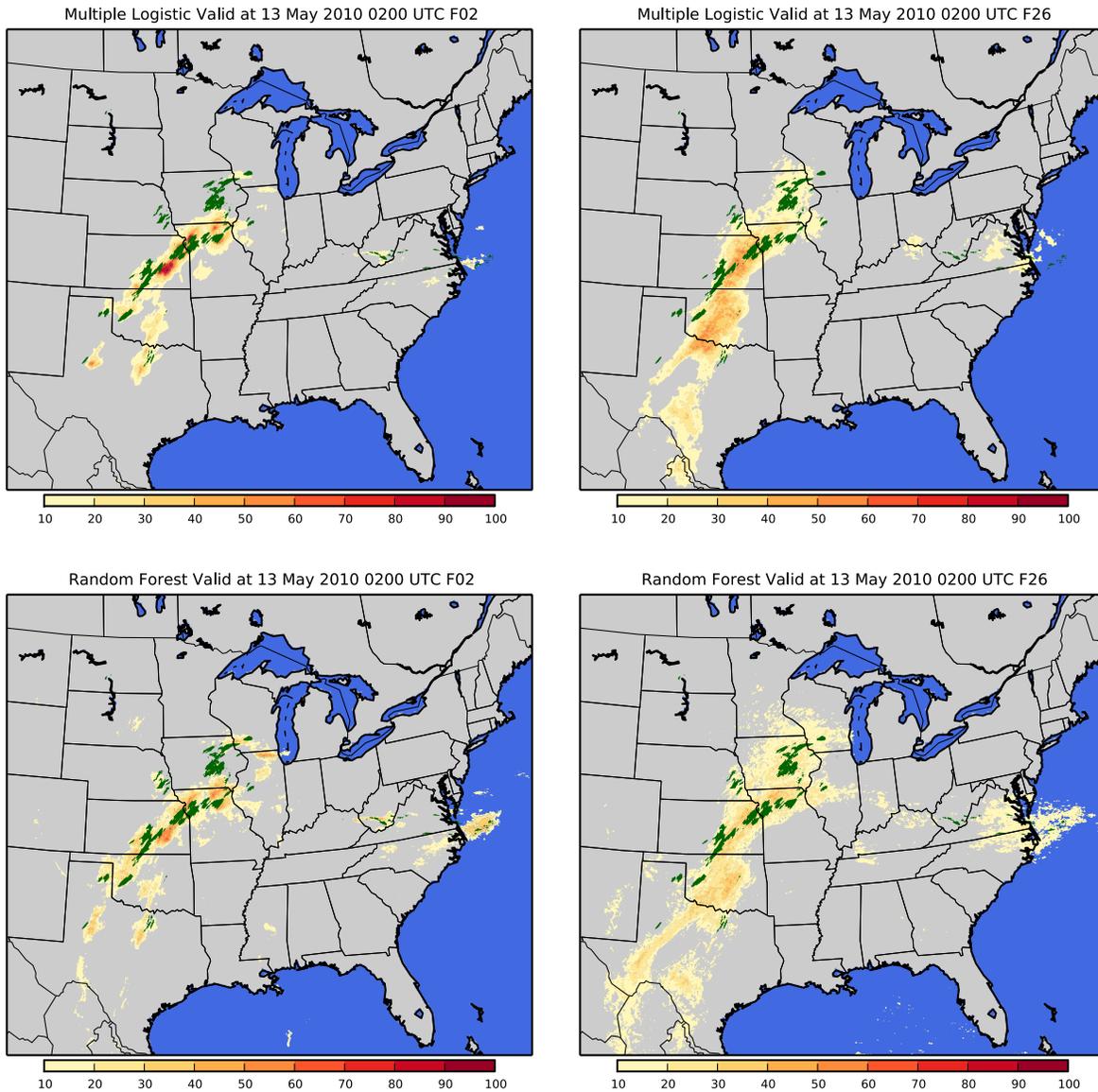


Figure 12: The 2-hour (left) and 26-hour (right) forecasts of the multiple logistic regression (top) and the random forest (bottom). The green areas indicate the observed areas of 1-hour precipitation greater than 6.35 mm h^{-1} .

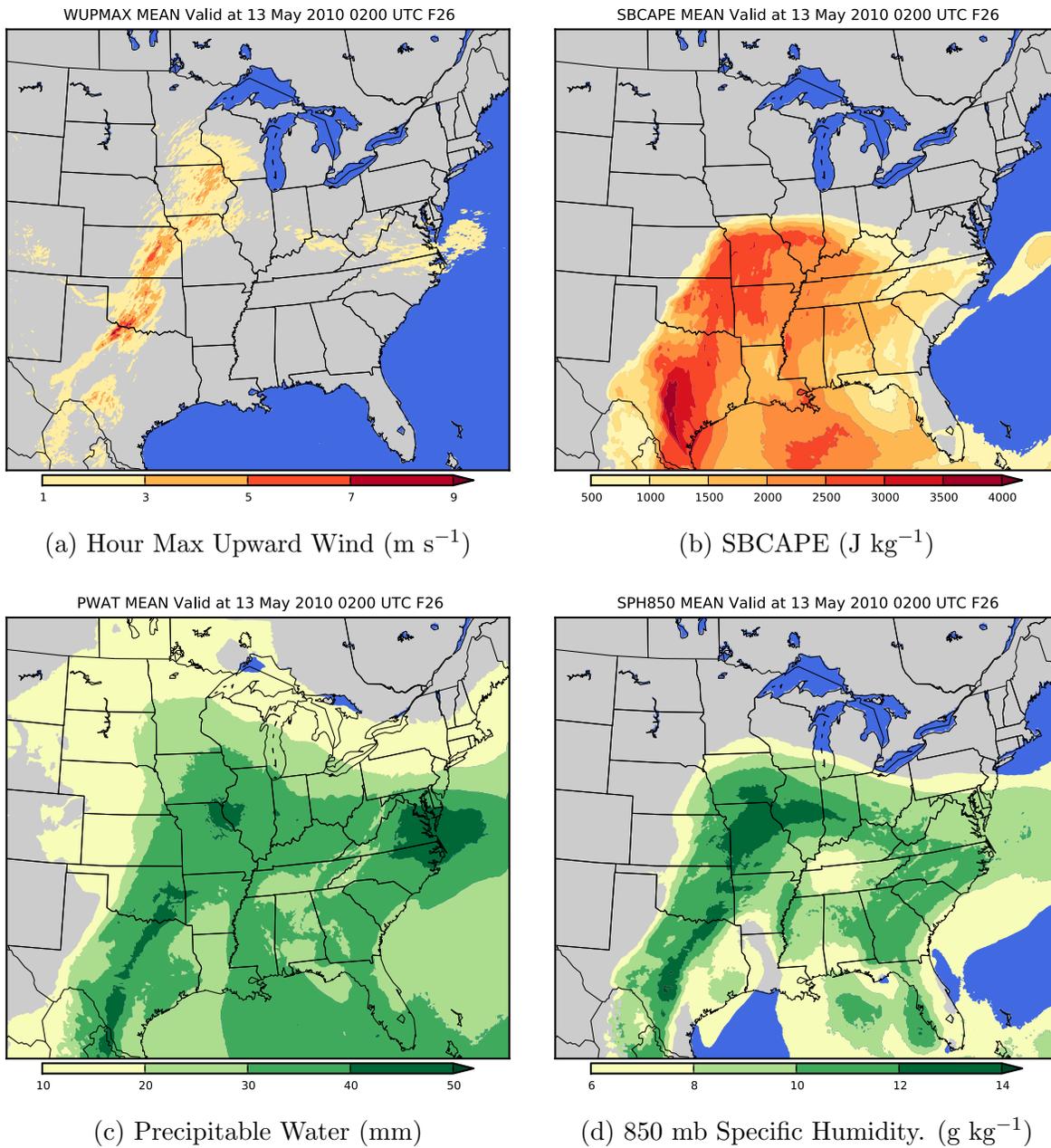


Figure 13: Maps of the SSEF 26-hour forecasts of predictors considered important by the random forest model.