

Article SubmissionRevised Article *Hydrology and Earth System Sciences*

Title:

Performance and Robustness of Probabilistic River Forecasts Computed with Quantile Regression based on Multiple Independent Variables in the North Central U.S.A.

Authors:

Frauke Hoss, Paul S. Fischbeck

Affiliation:

Carnegie Mellon University

Department of Engineering & Public Policy

5000 Forbes Avenue

Pittsburgh, PA 15213

Corresponding Author:

Frauke Hoss: fraukehoss@gmail.com

Performance and Robustness of Probabilistic River Forecasts Computed with Quantile Regression based on Multiple Independent Variables in the North Central U.S.A.

Abstract

This study ~~further develops the method of~~applies quantile regression (QR) to ~~the prediction of predict flood stage~~ exceedance probabilities ~~of flood stages by post-processing forecasts based on post-processing single-value flood stage forecasts~~. ~~A computationally cheap technique to predict forecast errors is valuable, because many national flood forecasting services, such as the National Weather Service (NWS), only publish deterministic single-value forecasts. Using data~~The study uses data from the 82 river gages, for which the ~~National Weather Service's~~NWS' North Central River Forecast Center issues forecasts daily, ~~this is the first QR application to U.S. American river gages~~. Archived forecasts for lead times up to six days from 2001-2013 were analyzed. ~~Earlier implementations of QR used the forecast itself as the only independent variable~~. ~~Besides the forecast itself, t~~This study ~~adds~~uses the ~~rise rate~~rate of rise of the river stage in the last 24 and 48 hours and the forecast error 24 and 48 hours ago ~~to as predictors in the QR model configurations~~. ~~Including those~~When compared to just using the forecast as ~~independent variable, adding the latter~~ four ~~variables~~predictors significantly improved the forecasts, as measured by the Brier Skill Score (BSS). Mainly, the resolution increases, as the ~~forecast-only original~~QR ~~implementation configuration~~ already delivered high reliability. Combining the forecast with the other four ~~variables~~predictors results in much less favorable BSSs. Lastly, the forecast performance does not ~~strongly~~ depend on the size of the training

22 | dataset, but on the year, the river gage, lead time and event threshold that are being forecast. We
23 | find that each event threshold requires a separate ~~model~~ configuration or at least calibration.

24 | **Keywords:** River forecasts, quantile regression, probabilistic forecasts, robustness

25 |

26 1 Introduction

27 River-stage forecasts ~~are inherently uncertain~~ are no crystal ball; the future remains uncertain.

28 The past has shown that unfortunate decisions have been made in ignorance of the potential
29 forecast errors (Pielke, 1999; Morss, 2010)~~(e.g., Pielke, 1999; Morss, 2010)~~. For many users,
30 such as emergency managers, forecasts are most important in ~~extreme~~ extreme situations, such as
31 droughts and floods. Unfortunately, it is exactly in those situations that forecast errors are
32 largest, due~~Due~~ to the ~~ir~~ infrequency of extreme events and the subsequent scarcity of data;
33 ~~forecasts have larger errors where accuracy has the most value~~. Additionally, users might only
34 experience such an event once or twice in their lifetime, so that they have no experience to what
35 extent they can rely on ~~deterministic~~ forecasts in such situations. Given the many sources and
36 complexity of uncertainty and the lacking user experience, it is easy to see how forecast users
37 find it difficult to estimate the forecast error. Including uncertainty in river forecast would
38 therefore be valuable, just as has been ~~weather forecasts has been strongly~~ recommended for
39 weather forecasts in general (e.g., National Research Council, 2006)~~(e.g., National Research~~
40 ~~Council, 2006)~~.

41 There are two types of approaches to ~~quantify~~ estimate forecast uncertainty (e.g., Leahy,
42 2007; Demargne et al., 2013; Regonda et al., 2013)~~(e.g., Leahy, 2007; Demargne et al., 2013;
43 ~~Regonda et al., 2013)~~: Those addressing ~~certain~~ major sources of uncertainty individually in the
44 output, e.g., input uncertainty and hydrological uncertainty, and those taking into account all
45 sources of uncertainty in a lumped fashion. Both approaches have their advantages. Modelling
46 each source separately can take into account that the different sources of uncertainty have
47 different characteristics (e.g., some sources of uncertainty depend on lead time, while others do
48 not). This approach is likely to result in better performing, more parsimonious~~

49 ~~model configurations~~. On the downside, ~~it the approach~~ is expensive to develop, maintain and
50 run. As an alternative, the lumped quantification of uncertainty is a less resource-intensive
51 approach (~~Regonda et al., 2013~~)(~~Regonda et al., 2013~~).

52 The National Weather Service has chosen ~~for ensemble forecasting to quantify the~~
53 ~~uncertainty from major source~~to quantify the most significant sources of uncertainty using
54 ~~ensemble techniques~~ (~~Demargne et al., 2013~~)(~~Demargne et al., 2013~~). ~~As of today~~Currently, the
55 National Weather Service does not routinely publish uncertainty information along with their
56 short-term river-stage forecast (~~(Figure 1)~~). ~~Until the NWS has implemented probabilistic~~
57 ~~forecasting for short-term products (next few hours and days), the only way that users can get a~~
58 ~~sense of the uncertainty is by comparing the quantitative precipitation forecast (QPF) with the~~
59 ~~non-QPF forecast. The QPF forecast includes the precipitation predicted for the next 12 hours~~
60 ~~and zero precipitation for the forecasts beyond 12 hours.~~[†] ~~The non-QPF forecast assumes no~~
61 ~~precipitation. Combined, these two forecasts give an idea of how much difference (a short period~~
62 ~~of) precipitation would make for the stage height in the river. The non-QPF serves as a~~
63 ~~reasonable lower bound; however, the QPF forecast is not an upper bound (i.e., precipitation~~
64 ~~could exceed the forecast values).~~

65 As of today, only the “outlooks” produced by the Ensemble Streamflow Prediction part
66 of the NWS River Forecasting System are probabilistic, i.e., quantify uncertainty: an exceedance
67 curve for a period of three month and bar plots for each week of a three months period, see and.
68 These graphs can be used to determine with which probability each river stage will be exceeded
69 in those weeks or three months period. Although the short term weather forecasts for the next

[†]This practice differs from RFC to RFC and also over time. For the ABRFC Welles et al. report: ~1993-1994: zero QPF; ~1995-2000 24hr QPF for first 24hrs, zero QPF beyond 24hrs; ~2001-2003 12hr QPF for first 12hrs, zero QPF beyond 12hrs.

92 ~~few days are much used to prepare for flood events, they have remained deterministic, as shown~~
 93 ~~in.~~²

94 **Figure 11: Deterministic short-term weather forecast in six hour intervals as published by the NWS**
 95 **for Hardin, IL on 24 April 2014.**

96 **Source:**<http://water.weather.gov/ahps2/hydrograph.php?wfo=lsx&gage=hari2>.

97 ~~The Figure 12: Probabilistic long-term forecast as published by the NWS for Commerce, OK on 14~~
 98 ~~December 2012: Exceedance curve for three months period. (Not available for Hardin, IL). Source:~~
 99 ~~<http://water.weather.gov/ahps2/hydrograph.php?wfo=tsa&gage=como2>~~

100 ~~Figure 3: Probabilistic long-term forecast as published by the NWS for Commerce, OK on 14~~
 101 ~~December 2012: Bar plot for each week of a three months period. (Not available for Hardin, IL).~~
 102 ~~Source: <http://water.weather.gov/ahps2/hydrograph.php?wfo=tsa&gage=como2>~~

103 ~~NWS has developed the Hydrologic Ensemble Forecast Service (HEFS) in to be able to~~
 104 ~~provide also short-term and medium-term probabilistic forecasts. Its implementation at all 13~~
 105 ~~river forecasts center is planned to be completed in 2014 (Demargne et al., 2013)(Demargne et~~
 106 ~~al., 2013). HEFS includes two types of post-processors. The Hydrologic Model Output Statistics~~
 107 ~~(HMOS) Streamflow Ensemble Processor – which is also a module in NWS’ main forecast tool,~~
 108 ~~the Community Hydrologic Prediction System (CHPS) – corrects bias and evaluates the~~
 109 ~~uncertainty of each ensemble, while Hydrologic Ensemble Post-Processing (EnsPost) corrects~~
 110 ~~bias and lumps the set of ensembles into one uncertainty estimate (Demargne et al., 2013; Seo,~~
 111 ~~2008). HMOS performs a similar task as the QR approach presented here, but with two major~~
 112 ~~differences. First, it relies on linear regression based on streamflows at various times as~~
 113 ~~predictor, instead of using QR with several types of independent variables. Second, it does not~~

²~~The deterministic forecasts are also available as text or tables.~~

114 compute distributions of water levels from which confidence intervals or exceedance
115 probabilities of flood stages can be derived, but generates ensembles (Regonda et al., 2013).

116 In contrast to ~~the an~~ ensemble approach ~~chosen by the NWS~~ such as HEFS, the statistical
117 post-processing ~~method that is further developed~~ in this paper —quantile regression— does not
118 distinguish between sources of uncertainty, but studies the overall uncertainty in a lumped
119 fashion. ~~This choice is motivated by the fact that the total predictive uncertainty, rather than its~~
120 ~~different sources, are relevant for decision making . To further strengthen the main advantage of~~
121 ~~this method, i.e., requiring relatively little resources, To make this approach useful for actors~~
122 with limited resources, we exclusively use publicly available data to ~~build our models~~ define our
123 configurations.

124 Most previously developed post-processors to generate probabilistic forecasts share the
125 overall set-up but differ in their implementation. Explanatory-Independent variables such as the
126 forecasted and observed river stage, river flow or precipitation, and previous forecast errors are
127 used to predict the forecast error, conditional probability distribution of the forecast error or
128 other ~~metries~~ measures of uncertainty for various lead times (e.g., Kelly and Krzysztofowicz,
129 1997; Montanari and Brath, 2004; Montanari and Grossi, 2008; Regonda et al., 2013; Seo et al.,
130 2006; Solomatine and Shrestha, 2009; Weerts et al., 2011)(e.g., ~~Kelly and Krzysztofowicz, 1997;~~
131 ~~Montanari and Brath, 2004; Montanari and Grossi, 2008; Regonda et al., 2013; Seo et al., 2006;~~
132 ~~Solomatine and Shrestha, 2009; Weerts et al., 2011). Among others, T~~ these methodtechniques
133 differ ~~in their mathematical methods~~ in a number of ways, including their sub-setting of data, and
134 the output ~~metri~~ e. Please see Regonda et al. ~~(2013)~~(2013) and Solomatine & Shrestha
135 ~~(2009)~~(2009) for a summary of each methodtechnique. In a meta-analysis of four different post-
136 processing methodtechniques to generate confidence intervals, the quantile regression

137 ~~method~~technique was one of the two most reliable ~~method~~techniques (~~Solomatine and Shrestha,~~
138 ~~2009~~)(~~Solomatine and Shrestha, 2009~~), while being the mathematically least complicated ~~method~~
139 and requiring few assumptions.

140 This paper further develops one of the ~~method~~techniques mentioned above: the Quantile
141 Regression ~~method~~approach to post-process river forecasts ~~first~~ introduced by ~~Wood et al.~~
142 ~~(2009)~~ and further elaborated by ~~Weerts et al. (2011)~~(~~2011~~) and ~~López López et al. (2014)~~. ~~–The~~
143 ~~Weerts~~at study achieved impressive results in estimating the 50% and 90% confidence interval
144 of river-stage forecasts for three case studies in England and Wales using QR with calibration
145 and validation datasets spanning two years each. ~~This paper combines elements of the studies~~
146 ~~mentioned above.~~ ~~–In some aspects, our approach differs from the original approach by Weerts~~
147 ~~et al. and López López et al.~~~~those three studies.~~ We predict the ~~probabilities that flood stages~~
148 ~~are exceeded~~exceedance probabilities of flood stages rather than uncertainty bounds, ~~because~~
149 ~~the former are more relevant to decision-making. In an attempt to balance missed alarms and~~
150 ~~false alarms, decision-makers are likely to resort to the best estimate (i.e., the deterministic~~
151 ~~forecast) rather than basing actions on the 50% or 90% confidence interval. Additionally,~~
152 ~~predicting the probability of an event corresponds with other forecasts with which users have~~
153 ~~much experience, e.g., the probability of precipitation. Morss et al. found in a survey of the~~
154 ~~general U.S. public that most people are able to base decisions on those forecasts.~~ Additionally,
155 we are fortunate to have a much larger dataset ~~than the three earlier studies~~, consisting of
156 archived forecasts for 82 river gages covering 11 years ~~available~~. ~~The study does not add to the~~
157 mathematical technique of quantile regression itself.

158 In this paper, the QR ~~method~~technique is applied to the 82 river gages of the North
159 Central River Forecast Center (NCRFC) encompassing (parts of) Illinois, Michigan, Wisconsin,
160 Minnesota, Indiana, North Dakota, Iowa, and Missouri.³

161 Identifying the best-performing set of independent variables is central to this paper. To
162 our knowledge, this paper is the first application of the QR method to the U.S. American context.

163 All possible combinations of the following predictors have been studied: forecast, the

164 The method is further developed by demonstrating the benefit—measured by an increase
165 in Brier Skill Score (BSS)—of including the rise raterate of rises of water levels in past hours,

166 and the past forecast errors as independent variables into the quantile regression. The

167 performance of these joint predictors has been measured and compared using the Brier Skill

168 Score (BSS). For extremely high water levels the variable combination has to be customized for

169 each river gage. For those, sets of few independent variables work best. Variable combinations

170 for other event thresholds should include as many dependent variables as possible. Using the

171 same combination for all of them works satisfactorily. Furthermore, it is found that the forecast—

172 the only independent variable in the original QR method—is difficult to combine with the other

173 dependent variables. Last, the method is shown to be robust to the size of the training dataset.

174 However, the forecast performance does vary significantly across locations, lead times, water

175 levels, and forecast year. This exercise has been repeated for various water levels and lead times.

176 Additionally, the robustness of the resulting QR configurations across different sizes of training

177 datasets, locations, lead times, water levels, and forecast year has been assessed.

178 The paper is structured as follows. The Method section ~~summarizes the additions that this~~

179 ~~paper makes to the quantile regression method introduced by Weerts et al. . It reviews the~~

³-As of spring 2014, the NCRFC does not publish any sort of probabilistic forecasts.

180 ~~method~~ quantile regression, ~~explains the additions~~, introduces the performance ~~metri~~ measure,
181 and discusses the ~~computations~~ performed analyses and data. The Results section first reviews
182 the overall forecast error for the dataset. ~~It then compares the proposed method to the original~~
183 ~~quantile regression as demonstrated for river gages in Wales and England~~. It then describes the
184 results of identifying the best-performing set of independent variables. Finally, it discusses the
185 robustness of the ~~proposed method~~ studied QR configurations. The fourth and last section
186 presents the conclusions and proposes further research ideas.

187 2 Method

188 The use of quantile regression to quantify estimate the error distribution of river-stage forecasts
189 has first been ~~presented~~ introduced by Woods et al. (2009) for the Lewis River in Washington
190 State. Later, by Weerts et al. (2011)(2011) applied it to ~~for~~ river catchments in ~~the~~ England and
191 Wales. ~~In this paper, we further develop Weerts' original method in three ways: a) by including~~
192 ~~additional variables instead of using only the forecast itself as an independent variable; elements~~
193 of both studies are combined. However, our predictand is the probability of exceeding flood
194 stages rather than confidence bounds. Additionally, this study tests ~~b) by testing~~ the robustness of
195 the ~~method~~ technique across locations, lead times, event thresholds, forecast years, and the size of
196 training dataset is tested. ~~; c) by estimating the more decision-relevant probability of exceeding~~
197 ~~flood stages rather than confidence bounds~~. To develop the different QR configurations ~~of~~
198 ~~quantile regression~~ and to compare their performance, the Brier Skill Score (BSS) is used.

199 In the following, ~~the~~ quantile regression itself ~~and~~ the proposed addition to the
200 method analysis to identify the best-performing set of independent variables ~~and the undertaken~~
201 ~~computations~~ are explained.

2.1 Quantile Regression

In the context of river forecasts, linear quantile regression has been used to estimate the distribution of forecast errors as a function of the forecast itself. Weerts et al. (2011)(2014) summarize this stochastic approach as follows:

“[It] estimates effective uncertainty due to all uncertainty sources. The approach is implemented as a post-processor on a deterministic forecast. [It] estimates the probability distribution of the forecast error at different lead times, by conditioning the forecast error on the predicted value itself. Once this distribution is known, it can be efficiently imposed on forecast values.”

Quantile Regression was first introduced by Koenker (2005; 1978)(2005; 1978). It is different from ordinary least square regression in that it predicts percentiles rather than the mean of a dataset. Koenker and Machado (Koenker and Machado, 1999, p.1305)(Koenker and Machado, 1999, p.1305) and Alexander et al. (2011)(2014) demonstrate that studying the coefficients and their uncertainty for different percentiles generates new insights, especially for non-normally distributed data. For example, using quantile regression to analyze the drivers of international economic growths, Koenker and Machado (1999)(1999) find that benefits of improving the terms of trade show a monotonously increasing trend across percentiles, thus benefitting faster-growing countries proportionally more.

~~In its original application to river forecasts by~~ When applying QR to river forecasts, Weerts et al. (2011)(2014) ~~transformed~~; the forecast values and the corresponding forecast errors ~~are transformed~~ into the Gaussian domain using Normal Quantile Transformation (NQT) to account for heteroscedasticity. Detailed instructions to perform NQT can be found in, ~~as instructed by~~ Bogner et al. (2012)(2012). ~~to account for heteroscedasticity.~~ Building on this study, López

243 López et al. ~~(2014)~~(2014) compare different configurations of QR with the forecast as the only
 244 independent variable, including configurations omitting NQT. They find that no configuration
 245 was consistently superior for a range of forecast quality ~~metries-measures~~ (López López et al.,
 246 ~~2014~~)(López López et al., 2014). To be able to combine ~~predictors-variables~~ of different nature,
 247 we ~~build-a model~~-based our QR configuration on untransformed ~~variables~~predictors. The reason
 248 to do so will be discussed and illustrated later (see Figure 11 and Figure 12).

249 ~~Using the transformed data,~~ A quantile regression is run for each lead time and desired
 250 percentile with the forecast error as the dependent variable and the forecast and other variables as
 251 ~~the~~ independent variables.⁴ To prevent the quantile regression lines from crossing each other, a
 252 fixed effects model is implemented below a certain forecast value. Weerts et al. ~~(2011)~~(2011)
 253 give a detailed mathematical description for applying QR to river forecasts. Mathematically, the
 254 approach is formulated as follows (with and without NQT):

255 **Equation 1: ~~Original QR implementation-configuration~~ with NQT**, with percentiles of the forecast
 256 error as the dependent variable and the ~~only one~~ independent variable ~~being the forecast itself~~, but
 257 transformed into the normal domain.

$$F_{\tau}(t) = f_{cst}(t) + NQT^{-1}[a_{\tau} * V_{NQT}(t) + b_{\tau}]$$

258 **Equation 2: QR ~~implementation-configuration~~ without NQT**, with percentiles of the forecast error
 259 as the dependent variable and multiple independent variables.

$$F_{\tau}(t) = f_{cst}(t) + \sum_i^I a_{i,\tau} * V_i(t) + b_{\tau}$$

260 with $F_{\tau}(t)$ – estimated forecast associated with percentile τ and time t

⁴~~As mentioned in Weerts et al. (2011), our quantile regression models have likewise a higher predictive capacity, if the forecast error rather than the forecast itself is used as the dependent variable.~~

261	$f_{cst}(t)$	– original forecast at time t
262	$V_i(t)$	– the independent variable i (e.g., the original forecast) at time t
263	$V_{i,NQT}(t)$	– the independent variable I transformed by NQT at time t
264	$a_{i,\tau}, b_\tau$	– model configuration coefficients
265		

266 The second part of the equations stands for the error estimate based on the quantile regression
 267 ~~model configuration~~ for each percentile τ and lead time. In Equation 1, that was used ~~in the~~
 268 ~~original QR method proposed~~ by Weerts et al. ~~(2011)(2014)~~, this estimation was executed in the
 269 Gaussian domain using only the forecast as independent variable. Our study mainly uses
 270 Equation 2, i.e., it does not transform the predictors and the predictand. All quantile regressions
 271 were done using the command $rq()$ in the R-package “quantreg” (Koenker, 2013).⁵

272 2.2 Brier Skill Score

273 The ~~original-QR implementation configuration~~ by Weerts et al. ~~(2011)(2014)~~ was evaluated by
 274 determining the fraction of observations that fell into the confidence intervals predicted by the
 275 QR ~~model configuration~~; i.e., ideally, ~~9080~~% of the observations should be larger than the
 276 predicted 10th percentile for that day, and smaller than the predicted 90th percentile. López López
 277 et al. ~~(2014)(2014)~~ used a number of ~~metrics-measures~~ to assess ~~model configuration~~
 278 performance, e.g., the Brier Skill Score (BSS), the mean continuous ranked probability (skill)
 279 score (RPSS), the relative operating characteristic (ROC), and reliability diagrams to compare
 280 QR configurations.

281 We use the Brier Skill Score ~~– first introduced by Brier (1950) – to compare-assess the~~
 282 ~~different versions of the-QR model configurations-proposed in this paper. We chose to optimize~~

⁵ ~~All quantile regressions were done using the command $rq()$ in the R-package “quantreg” (Koenker, 2013).~~

297 ~~our QR models based on the BSS, first introduced by Brier for two two reasons. First, to be able~~
 298 ~~to optimize model performance it is best to choose a single measure. First, for decision-making~~
 299 ~~the probability with which a certain water level, e.g., a flood stage, is exceeded is more useful~~
 300 ~~than confidence intervals. Second~~Second, out of the available measures the Brier Score is
 301 attractive, because it can be decomposed into two different measures of forecast quality (see
 302 Equation 3): Reliability and resolution. The third component is uncertainty, which is a
 303 hydrological characteristic inherent to the river gage. This uncertainty is different than the
 304 forecast uncertainty that the technique studied in this paper estimates. Besides the uncertainty
 305 that can be mathematically explained, it also includes natural variability. ThusIn sum, the BS'
 306 uncertainty term is not subject to the forecast quality. Equation 3 gives the definition of the (de-
 307 composed) Brier Score (e.g., Jolliffe and Stephenson, 2012; Wikipedia, 2014; WWRP/WGNE,
 308 2009)(e.g., Jolliffe and Stephenson, 2012; Anon, 2014; WWRP/WGNE, 2009).⁶

309 **Equation 3: Brier Score; de-composed into three terms: reliability, resolution and uncertainty.**

$$BS = \frac{1}{N} \sum_{k=1}^K n_k (f_k - \bar{o}_k)^2 - \frac{1}{N} \sum_{k=1}^K n_k (\bar{o}_k - \bar{o})^2 + \bar{o}(1 - \bar{o}) = \frac{1}{N} \sum_{t=1}^N (f_t - o_t)^2$$

310 with BS – Brier Score

⁶~~Bröcker (2012)(2012) showed that the conventional decomposition of the Brier Score is biased for finite sample sizes. It systematically overestimates reliability, under- or overestimates resolution, and underestimates uncertainty. Several authors proposed less biased decompositions (e.g., Bröcker, 2012; Ferro and Fricker, 2012)(e.g., Bröcker, 2012; Ferro and Fricker, 2012). Additionally, Stephenson et al. (2008)(2008) proved that the Brier Score has two additional components when it is computed based on bins, as is usually done. Nonetheless, we chose to stick to the conventional decomposition and using bins, as implemented in the R-package “verification” (NCAR Research Applications Laboratory, 2014; Wilks, 1995)(NCAR Research Applications Laboratory, 2014; Wilks, 1995) to ensure that our results can be readily compared to other studies like López López et al. (2014)(2014). After all, the Score is mainly used to compare model configurations, rather than establishing the absolute performance of each model configuration.~~

311	N	– number of forecasts
312	K	– the number of bins for forecast probability of binary event occurring on each
313	day	
314	n_k	– the number of forecasts falling into each bin
315	\bar{o}_k	– the frequency of binary event occurring on days in which forecast falls into bin
316	k	
317	f_k	– forecast probability
318	\bar{o}	– frequency of binary event occurring
319	f_t	– forecast probability at time t
320	o_t	– observed event at time t (binary: 0 – event did not happen, 1 – event happened)

321 The Brier Score pertains to binary events, e.g., the exceedance of a certain river stage or
322 flood stage. Reliability compares the estimated probability of such an event with its actual
323 frequency. For example, perfect reliability means that on 60% of all days for which it was
324 predicted that the water level would exceed flood stage with a 60% probability, it actually does
325 so. ~~A forecast with~~ The reliability curve for the forecast representing perfect reliability would
326 follow the diagonal in Figure 2, i.e., the area in Figure 2a representing reliability would equal
327 zero (Jolliffe and Stephenson, 2012; Wikipedia, 2014; WWRP/WGNE, 2009)(e.g., Jolliffe and
328 Stephenson, 2012; Anon, 2014; WWRP/WGNE, 2009). The configuration by López López et al.
329 (2014)(2014) performs well in terms of reliability. When estimating confidence intervals, Weerts
330 et al. (2011)(2011) achieved good results especially for the more extreme percentiles (i.e., 10th
331 and 90th).

332 Figure 2: Theory behind Brier Skill Score illustrated for an imaginary forecast (red line): (a)
333 reliability and resolution; (b) skill. In figure a, the area representing reliability should be as small,
334 and for resolution as large as possible. The forecast has skill (BSS > 0), i.e., performs better than the
335 reference forecast, if it is inside the shaded area in the figure b. Ideally, the forecast would follow
336 the diagonal (BSS=1). (Adapted from Hsu and Murphy, 1986; Wilson, n.d.).

337 ~~Figure 4: Theory behind Brier Skill Score illustrated for an imaginary forecast (red line): (a)~~
338 ~~reliability and resolution; (b) skill. In figure a, the area representing reliability should be as small,~~
339 ~~and for resolution as large as possible. The forecast has skill ($BSS > 0$), i.e. performs better than~~
340 ~~random guessing, if it is inside the shaded area in the figure b. Ideally, the forecast would follow the~~
341 ~~diagonal ($BSS=1$). (Adapted from Hsu and Murphy, 1986; Wilson, n.d.).~~

342 Resolution ~~pertains to how much better the forecast performs than taking the historical~~
343 ~~frequency (climatology) as a forecast.~~measures the difference between the predicted probability
344 of an event on a given day and the observed average probability. When calculated for a time
345 period longer than a day, the forecast performs better if the resolution term is higher. -For
346 example, for a gage where flood stage is exceeded on 5% of the days in a year, simply using the
347 historical frequency as the forecast would mean forecasting that the probability of the water level
348 exceeding flood stage is 5% on any given day. The accumulated difference between the
349 predicted frequency and the historical average across a time period of several days would then be
350 zero (e.g., Jolliffe and Stephenson, 2012; Wikipedia, 2014; WWRP/WGNE, 2009)(e.g., Jolliffe
351 and Stephenson, 2012; Anon, 2014; WWRP/WGNE, 2009). In Figure 2, the curve for a a
352 forecast with good resolution would be steeper than the dashed line that represents climatology,
353 i.e., the area in aFigure 2a representing resolution would be maximized. In absolute terms, the
354 resolution can never exceed the third term in Equation 3 representing the uncertainty inherent to
355 the river gage. Through the resolution component, the Brier Score is related to the area under the
356 relative operating characteristic (ROC) curve (for more detail, see Ikeda et al., 2002)(for more
357 detail, see Ikeda et al., 2002). The latter likewise quantifies how much better ~~a forecast is~~ than
358 ~~random guessing~~the reference forecast (i.e., climatology) a forecast is -in detecting a binary
359 event; though unlike the Brier Score it focuses on the ratios of false and missed alarms (e.g.,

360 Jolliffe and Stephenson, 2012; Wikipedia, 2014; WWRP/WGNE, 2009)(e.g., Jolliffe and
361 Stephenson, 2012; Anon, 2014; WWRP/WGNE, 2009).

362 A forecast possesses skill, i.e., performs better than ~~random guessing or climatology~~the
363 reference forecast (in this case climatology), if it is inside the shaded area in Figure 2b. The
364 Brier *Skill* Score (BSS) equals the Brier Score normalized by climatology to make the score
365 comparable across gages with different frequencies of a binary event. Equation 4 defines the
366 BSS' decomposition into the resolution and reliability components described above (Brown and
367 Seo, 2013).⁷-The BSS can range from minus infinity to one. A BSS below zero indicates no
368 skill; the perfect score is one (e.g., Jolliffe and Stephenson, 2012; Wikipedia, 2014;
369 WWRP/WGNE, 2009)(e.g., Jolliffe and Stephenson, 2012; Anon, 2014; WWRP/WGNE, 2009).
370 All measures of forecast quality were computed using the R-package “verification” (NCAR,
371 2014).

372 **Equation 4: Decomposition of Brier Skill Score**

373
$$BSS = 1 - \frac{BS}{\bar{o}(1-\bar{o})} = \frac{RES}{\bar{o}(1-\bar{o})} - \frac{REL}{\bar{o}(1-\bar{o})}$$

374 with BSS – Brier Skill Score
375 BS – Brier Score
376 RES – Resolution
377 REL – Reliability
378 \bar{o} – Frequency of binary event occurring
379 $\bar{o}(1 - \bar{o})$ – Climatological variance
380

⁷-All measures of forecast quality were computed using the R-package “verification” (NCAR,
2014).

381 **2.3 ~~Proposed addition: More than one independent variable~~Identifying the best-performing**
382 **sets of independent variables**

383 ~~Intuitively, more information should lead to better prediction of the distribution of the forecast~~
384 ~~error, because the regression models would be based on more data~~The challenge is to identify a
385 well-performing set of predictors that is both parsimonious and comprehensive. Wood et al.
386 (2009) found rate of rise and lead time to be informative independent variables. Weerts et al.
387 (2011) achieved good results using only the forecast itself as predictor. Besides these variables,
388 ~~t~~The most obvious ~~variables predictors~~ to include ~~besides the forecast itself~~ are the observed
389 water level 24 and 48 hours ago, ~~the observed rise in water level in the last 24 and 48 hours~~
390 (~~called rise rate hereafter~~), the forecast error 24 and 48 hours ago (i.e., the difference between the
391 current water level at issue time of the forecast and the forecast that was produced 24/48 hours
392 ago), or the time of the year, e.g., using month or season as categorical predictors. Other
393 Additional potential ~~variables independent variables~~ are the water levels observed up- and
394 downstream at various times, the precipitation upstream of the catchment area, and the
395 precipitation forecast. However, requesting the corresponding precipitation and precipitation
396 forecast requires an extensive effort or direct access to the database. ~~these latter variables are~~
397 ~~much more difficult to gather because of the way data is archived~~database at the National
398 Climatic Data Center (NCDC).⁸

⁸~~For the NCRFC, the river forecast and the observed water levels are saved in the same text product available at [last accessed July 2014]:~~
~~<http://edo.ncdc.noaa.gov/pls/plhas/HAS.FileAppSelect?datasetname=9957ANX>. (Station ID: KMSR, Bulletin ID: FGUS5). Requesting the corresponding precipitation and precipitation forecast requires an extensive effort or direct access to the database.~~

399 **Table : Variable Combinations**

400 In preliminary trials on two case studies (gages HARI2 and HYNI2), it was found that the
401 rates of rise and the forecast errors are better predictors than the water levels observed in
402 previous days. After all, the observed water levels are used to compute the rates of rise and
403 forecast errors, so that these latter variables include the information of the former variable. It was
404 also found that season and months are not significant in quantile regression configurations to
405 predict the quantiles of the forecast error. Probably, the time of the year is already reflected in
406 the observed water levels and forecast errors in the previous days. In preliminary trials on two
407 case studies (gages HARI2 and HYNI2), it was found that season and months are not significant
408 in quantile regression models to predict the quantiles of the forecast error. It was also found that
409 the rise rates and the forecast errors are better predictors than the water levels observed in
410 previous days. After all, the observed water levels are used to compute the rise rates and forecast
411 errors, so that these latter variables include the information of the former variable.

412 To determine which set of predictors performs best in generating probabilistic forecasts,
413 all 31 possible combinations of the forecast (fcst), the rate of rise in the last 24 and 48 hours
414 (rr24, rr48), and the forecast error 24 and 48 hours ago (err24, err48) – see Equation 5 – were
415 tested for 82 gages that the NCRFC issues forecasts for every morning (Table 1). Based on the
416 Bier Skill Score, it was determined which joint predictor on average and most often leads to the
417 best out-of-sample results for various lead times and water levels.

418 Equation 5: QR configuration without NQT, with percentiles of the forecast error as the dependent
 419 variable and varying combinations of the five independent variables. This equation was used to
 420 predict the water level distribution for each day at 82 gages with different lead times.

$$F_{\tau}(t) = fcst(t) + a_{fcst,\tau} * fcst(t) + a_{rr24,\tau} * rr24(t) + a_{rr48,\tau} * rr48(t) \\ + a_{err24,\tau} * err24(t) + a_{err48,\tau} * err48(t) + b_{\tau}$$

421 with $F_{\tau}(t)$ – estimated forecast associated with percentile τ and time t
 422 $fcst(t)$ – original forecast at time t
 423 $rr24(t), rr48(t)$ – rates of rise in the last 24 and 48 hours at time t
 424 $err24(t), err48(t)$ – forecast errors 24 and 48 hours ago (e.g., the original forecast) at
 425 time t
 426 $a_{xx,\tau}, b_{\tau}$ – configuration coefficients; forced to be zero if the predictor is
 427 excluded from the joint predictor that is being studied.

428

429 ~~To determine which set of variables preforms best in generating probabilistic forecasts, all 31~~
 430 ~~possible combinations of the forecast (fcst), the rise rate in the last 24 and 48 hours (rr24, rr48), and~~
 431 ~~the forecast error 24 and 48 hours ago (err24, err48) were tested for 82 gages that the NCRFC~~
 432 ~~issues forecasts for every morning (-). Based on the Bier Skill Score, a metric of forecast quality~~
 433 ~~explained below, it was determined which variable combination on average and most often leads to~~
 434 ~~the best out-of-sample results for various lead times and water levels. Table 1: Joint predictors.~~

435

436 2.4 Computations

437 The output of our QR application to river forecasts is the probability that a certain water level in
 438 the river or flood stage is exceeded on a given day, e.g., “On the day after tomorrow, the
 439 probability that the river exceeds 15 feet at location X is 60%.” This is done in two steps. First, a
 440 training dataset (first half of the data) is used to build-define one quantile regression
 441 modelconfiguration for each-each of the following percentiles: $\pi \equiv [0.05, 0.1, 0.15, \dots, 0.85,$

442 0.90, 0.95] and each lead time.- The dependent variable is the water level. As described ~~above~~in
443 Equation 5, the forecast itself, the ~~rise rates~~rates of rise and forecast errors serve as independent
444 variables.

445 In the second step, these QR ~~model~~configurations are used to predict the water levels
446 corresponding with each ~~model's~~ percentile on each day in the verification dataset (the second
447 half of the dataset). Effectively, for each day in the verification dataset, a discrete probability
448 distribution of water levels is predicted. Each predicted QR model percentile π contributes one
449 point to that distribution.

450 ~~In our opinion, this probability distribution of water levels is too much information to~~
451 ~~efficiently make decisions. The model performance should be assessed for a decision relevant~~
452 ~~output. Therefore~~Then, -we calculate the probability with which various water levels (called
453 event thresholds hereafter) will be exceeded. The probability of exceeding each water level is
454 computed by linearly interpolating between the points of the discrete probability distribution that
455 was computed in the previous step.⁹

456 To be able to compare various ~~model~~-configurations, the Brier Skill Score is determined
457 ~~across all the days in~~based on forecast exceedance probability for all days in the verification
458 dataset. As explained above, the BSS is based on the difference between the predicted
459 exceedance probability and the observed exceedance (binary) averaged across all days in the
460 verification dataset.

461 To study whether the various combinations of ~~variables~~predictors perform equally well
462 for high and low thresholds, these last computational steps (i.e., interpolating to determine the

⁹ ~~Using the command “approx(x, y, xout, yleft=1, yright=0, ties=mean)” in the R package “stats”~~
~~(R Core Team, 2014).~~

463 exceedance probability for a certain water level and calculating the BSS) were done for the 10th,
464 25th, 75th, and 90th percentile of observed water levels and the ~~decision-relevant~~ four decision-
465 relevant flood stages (action stage, and minor, moderate, and major flood stage) of each gage.
466 Flood stages indicated when material damage or substantial hinder is caused by high water
467 levels. Therefore, the flood stages correspond with different percentiles at different river gages.

468 To determine the ~~optimal~~best-performing set of independent variables, the entire procedure is
469 repeated for each of the 31 ~~variable-combination~~joint predictors in Table 1, thus using a different
470 set of independent variables each time. To test the robustness of this approach, the procedure was
471 also repeated for each river gage and for several lead times. The result is 31 BSSs for 82 river
472 gages for four different lead times (one to four days) and for ~~different~~eight event thresholds (i.e.,
473 flood stages or percentiles of the observed water level).

474

475 **2.5 Data**

476 The National Weather Service (NWS) ~~issues river stage forecasts for ~4,000 river gages every~~
477 ~~day. Such's~~ daily published-short-term river forecasts predict the stage height in six-hour
478 intervals for up to five days ahead (20 6-hour intervals).⁴⁰ When floods occur and increased
479 information is needed, the local river forecast center (RFC) can decide to publish river-stage
480 forecasts more frequently and for more locations. Welles et al.- ~~(2007)~~(2007) provides a detailed
481 description of the forecasting process.

⁴⁰~~The river stage forecasts are produced by one of NWS' thirteen river forecasts centers (RFCs). Every morning the forecasts are forwarded to one of NWS's 122 local weather forecast offices (WFOs), who then disseminate the information to the public through a variety of media channels or by issuing warnings.~~

482 For this paper, all forecasts published by the North Central River Forecast Center
483 (NCRFC) between 1 May 2001 and 31 December 2013 were requested from the NCDC's HDSS
484 Access System ([National Climatic Data Center, 2014; Station ID: KMSR, Bulletin ID:
485 FGUS5](#)).^{††} In total, the NCRFC produces forecasts for 525 gages. For 82 of those gages,
486 forecasts have been published daily for a sufficient number of years, and are not inflow forecasts.
487 The latter have been excluded from the forecast error analysis because they forecast discharge
488 rather than water level. About half of the analyzed gages are along the Mississippi River ([Figure
489 3](#)). The Illinois River and the Des Moines River are two other prominent rivers in the region. The
490 drainage areas of the 82 river gages average 61,500 square miles (minimum 200 sq.miles;
491 maximum 708,600 sq.miles). [Figure 4 shows an empirical cumulative density function of
492 drainage areas sizes.](#)

493 [Figure 3: River gages for which the North Central River Forecast Centers publishes forecasts daily.
494 Henry \(HYN12\) and Hardin \(HARI2\) are indicated by the upper and lower red arrow respectively.
495 For gages indicated by black dots the basin size is missing.](#)

496 [Figure 4: Empirical cumulative density function \(ecdf\) of sizes of drainage area for the river gages
497 that are being forecasted daily by the NCRFC.](#)

498
499 Two river gages serve as an illustration for the points made throughout this paper.
500 Hardin, IL is just upstream [of](#) the confluence of the Illinois River and the Mississippi River
501 ([Figure 3](#)). Therefore, it probably experiences high water levels through backwatering, when the
502 high water levels in the Mississippi River prevent the Illinois River from draining. Henry, IL is

^{††} [URL \[last accessed July 2014\]:
http://edo.ncdc.noaa.gov/pls/plhas/HAS.FileAppSelect?datasetname=9957ANX; Station ID:
KMSR, Bulletin ID: FGUS5.](#)

521 | located ~200 miles (~~~320 km~~) upstream of Hardin, having a difference in elevation of ~25 feet.
522 | (~~~7.6 m~~). The Illinois River is ~330 miles (~~~530 km~~) long (Illinois Department of Natural
523 | Resources, 2011),¹² draining an area of ~13,500 square miles (~~~35,000 km²~~) at Henry (USGS,
524 | 2015a)¹³ and ~28,700 square miles (~~~72,000 km²~~) at Hardin (USGS, 2015b).¹⁴

525 | **Figure 5: Portion of the North Central River Forecast Centers river gages with Henry (HYN12) and**
526 | **Hardin (HARI2) indicated by the upper and lower red arrow respectively. Source:**
527 | **<http://www.erh.noaa.gov/ncrfc/>**

Field Cod

528 | 3 Results

529 | 3.1 Forecast error at NCRFC's gages

530 | In general, the NCRFC's forecasts are well calibrated across the entire dataset. The average
531 | error, defined as observation minus the forecast, is zero for most gages. For lead times longer
532 | than three days, a slight underestimation by the forecast is noticeable. By a lead time of 6 days
533 | this underestimation averages 0.41 feet only (~~a, a~~Figure 5a, Figure 6). Extremely low water
534 | levels, defined as below the 10th percentile of observed water levels, are also well calibrated
535 | (Figure 5b, Figure 6). (~~b, b~~). However, when considering higher water levels the picture
536 | changes.¹⁵ The underestimation becomes more pronounced, averaging 0.29 feet for three days of
537 | lead time and 1.14 feet for six days of lead time, when only observations exceeding the 90th
538 | percentile of all observations are considered (Figure 5c, Figure 6). (~~e, e~~). When only looking at

¹² ~~Illinois Environmental Protection Agency: "Illinois River and Lakes Fact Sheets", URL [accessed 04/24/2014]: <http://dnr.state.il.us/education/aquatic/aquaticillinoisrivlakefactshts.pdf>~~

¹³ ~~Source: http://waterdata.usgs.gov/nwis/nwisman/?site_no=05558300&agency_cd=USGS~~

¹⁴ ~~Source: http://waterdata.usgs.gov/nwis/nwisman/?site_no=05587060&agency_cd=USGS~~

¹⁵ ~~The gages MOR12 and MMO12 are upstream of a dam. It is likely that the forecasts performed so poorly there, because the dam operators deviated from the schedules that they provide the river forecast centers to base their calculations on.~~

560 observations that exceeded the minor flood stages corresponding to each gage,¹⁶ the
 561 underestimation averages 0.45 feet for three days of lead time and 1.51 feet for 6 days of lead
 562 time (Figure 5d, Figure 6). (Figure 6d, Table 2d). However, some gages, such as Morris
 563 (MORI2), Marseilles Lock/Dam (MMOI2) – both on the Illinois River – and Marshall Town on
 564 the Iowa River (MIWI4) experience *average* errors of 5 to 12 feet for water levels higher than
 565 minor flood stage. The gages MORI2 and MMOI2 are upstream of a dam. It is likely that the
 566 forecasts performed so poorly there, because the dam operators deviated from the schedules that
 567 they provide the river forecast centers to base their calculations on.

568 **Figure 6~~5~~:** Forecast error for 82 river gages that the NCRFC publishes daily forecasts for. In anti-
 569 clockwise direction starting at the top left: (a) Average error; (b) error on days that the water level
 570 did not exceed the 10th percentile of observations; (c) error on days that the water level exceeded the
 571 90th percentile of observations; (d) error on days that the water level exceeded minor flood stage.

572 Figure 6: Empirical cumulative distribution function (ecdf) of forecast error at 82 river gages for
 573 six lead times. Vertical lines show the median forecast error of the corresponding subset.

574 ~~Table 2: Error statistics for the forecast error a) of the whole dataset; b) on days that the water~~
 575 ~~level did not exceed the 10th percentile of observations; c) on days that the water level exceeded the~~
 576 ~~90th percentile of observations; d) on days that the water level exceeded minor flood stage.~~

577 3.2 Including more variables Identifying the best-performing sets of independent variables

578

579 In total, the Brier Skill Score (BSS) for 31 ~~variable combination~~ joint predictors (Table 1) across
 580 various lead times and event threshold have been compared. Across 82 river gages, it has been

¹⁶ ~~Flood stages are based on the damage done by previous floods. It depends on the context, e.g., the shape of the river bed and the development of the river shores, which water levels cause damage. Therefore, it depends on the river gage which percentiles of observed water levels the flood stages correspond with.~~

603 analyzed (a) which combinations perform best and worst most often, and (b) which ~~sets of~~
604 ~~variables~~joint predictor delivers the best BSSs on average.

605 3.2.1 Frequency Analysis

606 For ~~each the four~~ lead time (i.e., one to four days) and ~~various the eight~~ event thresholds (i.e.,
607 10th, 25th, 75th, 90th percentiles as well as the four flood stages), we counted ~~how often~~at how
608 many river gages each ~~variable combination~~joint predictor resulted in the highest and the lowest
609 BSS ~~across the 82 river gages~~. Figure 7 shows that for water levels below the 50th percentile
610 ~~variable combination~~joint predictors with four or more independent variables return the best
611 BSSs most often, while those with one and two ~~variables~~ predictors perform worst most often.
612 For thresholds higher than the 50th percentile the distributions gradually become ~~more flat~~flatter.
613 For the 90th percentile, a clear trend is no longer detectable. Given that the frequency
614 distributions for the extreme events in Figure 7 are relatively uniform, it seems as if extreme
615 events are characterized by different processes at different gages. The same set of histograms for
616 the four flood stages (i.e., action, minor, moderate, and major) confirms this (Figure 8). Across
617 lead times, there is a slight trend noticeable that single ~~variables~~ predictors tend to be the worst
618 combination more often for longer lead times. This suggests that~~us~~, the further out one is
619 forecasting, the more important it becomes to include more data in the ~~model~~configuration.

620 **Figure 7: Histograms of ~~variable combination~~joint predictors returning the best and worst Brier**
621 **Skill Scores across 82 river gages. Each row of histograms refers to an event threshold defined as a**
622 **percentile of the observed water levels, and each column to a lead time. The dotted vertical lines in**
623 **the histograms distinguish ~~variable combination~~joint predictors with different numbers of**
624 **independent variables.**

625 **Figure 8: Histograms of ~~variable-combination~~joint predictors** returning the best and worst Brier
626 Skill Scores across 82 river gages. Each row of histograms refers to a flood stage, and each column
627 to a lead time. The dotted vertical lines in the histograms distinguish ~~variable-combination~~joint
628 ~~predictors~~ with different numbers of ~~independent~~ variables.

629 3.2.2 Best performing combinations on average

630 For each river gage, the combinations have been ranked by BSSs. It was found that the more
631 ~~independent~~ variables are included in a ~~set~~joint predictor, the higher that set of ~~variables~~
632 ~~predictors~~ will rank on average (Figure 9). However, for extremely high water levels, this trend
633 gradually reverses (Figure 10). For action stage¹⁷ and minor flood stage,¹⁸ a slightly increasing
634 trend is still visible. For moderate¹⁹-and major flood stage,²⁰ combinations with fewer
635 ~~independent~~ variables rank higher on average. The most likely explanation is that extreme events
636 like major and moderate flood stage are infrequent. After all, major flood stage equals 90th to
637 100th percentiles at the various gages. This data scarcity can lead to overfitting when using more
638 predictors.

639 Considering these findings and those of the frequency analysis earlier, the
640 ~~model~~configuration for the various river gages can generally be based on the same ~~variable~~
641 ~~combination~~joint predictors of four or more ~~independent~~ variables. But for extremely high water
642 levels, a ~~model~~configuration specific to each river gage has to be built in order to achieve high
643 BSSs.

¹⁷~~-Across the 82 stations, action stage corresponds with water levels between the 60th and 100th percentile.~~

¹⁸~~-Across the 82 stations, minor flood stage corresponds with water levels between the 70th and 100th percentile.~~

¹⁹~~-Across the 82 stations, moderate flood stage corresponds with water levels between the 80th and 100th percentile.~~

²⁰~~-Across the 82 stations, major flood stage corresponds with water levels between the 90th and 100th percentile.~~

644 The combinations including the forecast (indicated by gray vertical lines in Figure 9 and
645 Figure 10) perform less well than those that exclude it. Plotting the independent variables against
646 the forecast error as the dependent variable makes the reason visible (Figure 11, Figure 12).

647 Without a transformation into the normal domain, the ~~forecast does not provide a lot of~~
648 ~~information for the QR model~~ scatterplot of forecast and forecast error does not show a trend.
649 After NQT, the percentiles show trends laid out like a fan. -In contrast, ~~the other four variables~~
650 ~~do not lend themselves for linear quantile regression after performing NQT~~ the other four
651 predictors become uniform distributions after NQT transformation. There is no trend detectable
652 anymore. Further research is necessary to reconcile these two types of ~~variables~~ predictors. A
653 possible solution could be to ~~build~~ define QR ~~model~~ configurations for subsets of the transformed
654 dependent and independent variable.

655 **Figure 9: Average rank for each ~~variable combination~~ joint predictor for one to four days of lead**
656 **time and four percentiles of observed water levels. Vertical gray lines indicate ~~variable~~**
657 **~~combination~~ joint predictors including the forecast.**

658 **Figure 10: Average rank for each ~~variable combination~~ joint predictor for one to four days of lead**
659 **time and four flood stages. Vertical gray lines indicate ~~variable combination~~ joint predictors**
660 **including the forecast.**

661 **Figure 11: Independent variables plotted against the forecast error for Hardin IL with 3 days of**
662 **lead time. First row: Forecast; second row: past forecast errors; third row: ~~rise rates~~ rates of rise.**

663 **Figure 12: Independent variables after transforming into the Gaussian domain plotted against the**
664 **forecast error for Hardin IL with 3 days of lead time. First row: Forecast; second row: past forecast**
665 **errors; third row: ~~rise rates~~ rates of rise.**

666 **3.2.3 Brier Skill Score**

667 Figure 13 illustrates the BSS when using the forecast as the only predictor as studied by Weerts
668 et al. (2011). Confirming Wood et al.'s findings (2009), additionally including the rise rate
669 of rise and forecasts errors as independent variables into the QR model configuration improves
670 the Brier Skill Score (BSS) significantly. ~~illustrates the BSS when using the model as~~
671 ~~originally introduced by Weerts et al.~~ Using the best performing ~~variable combination~~ joint
672 predictors instead, gives an upper bound of the BSSs that can be achieved at best. This
673 configuration increases the mean and decreases the standard deviation (~~→~~) (Figure 14, Figure 16).
674 The performance improves most where all ~~model~~ configurations perform worst: at the 10th
675 percentile. Possibly, the configurations do not perform well for low percentiles, because the
676 dependent variable – the forecast error – exhibits very little variance at those water levels, i.e.,
677 the average error is very small (Figure 16).²⁴ The decrease of the BSSs with lead time also
678 becomes considerably less with this configuration. Additionally, ~~an~~ one-size-fits-all approach
679 was tested to investigate, whether customizing the QR ~~model~~ configuration to each river gage
680 would be worth it. In this configuration, the ~~rise rates~~ rates of rise in the past 24 and 48 hours and
681 the forecast errors 24 and 48 hours ago serve as the independent variables (combination 30). It
682 was found that this approach returns only slightly worse results than working with the best

²⁴ ~~Possibly, the model configurations do not perform well for low percentiles, because the dependent variable – the forecast error – exhibits very little variance at those water levels, i.e., the average error is very small (Figure 6: Empirical cumulative distribution function (ecdf) of forecast error at 82 river gages for six lead times. Vertical lines show the median forecast error of the corresponding subset.~~

Table 2).

683 performing configuration for each river gage deviation (Figure 15, Figure 16). ~~(;)-~~ Accordingly,
684 the same ~~variable combination~~joint predictor can be used for all river gages.

685 As ~~shown in, already discussed earlier,~~ this last conclusion is not true for extremely high
686 water levels. Including more independent variables does improve the BSSs considerably
687 deviation (Figure 17,18, and 19). ~~(and ;)-~~ However, for each river gage the best ~~combination of~~
688 ~~variables~~joint predictor needs to be identified separately. Because data to ~~build models~~define
689 configurations is scarce for extreme levels, the QR ~~model~~configurations all perform less well for
690 each increase in flood stage.

691

692 **Table 3: Mean and standard deviation three QR configurations: the original using the transformed**
693 **forecast only as independent variable; the best performing combination for each river gage (upper**
694 **performance limit); rise rates in the past 24 and 48 hours and the forecast errors 24 and 48 hours**
695 **ago as independent variable (one size fits all solution).**

696 **Figure 13: Brier Skill Scores of the original forecast only QR model configuration (i.e., using the**
697 **transformed forecast as the only independent variable) for four lead times and percentiles of**
698 **observed water levels.**

699 **Figure 14: Brier Skill Scores for four lead times and percentiles of observed water levels using the**
700 **best ~~variable combination~~joint predictor for each river gage as independent variables in the QR**
701 **model configuration.**

702 **Figure 15: Brier Skill Scores for four lead times and percentiles of observed water levels using a**
703 **one-size-fits-all approach (i.e., rr24, rr48, err24, err48) for the independent variables in the QR**
704 **model configuration.**

705 **Figure 16: Empirical cumulative density functions of three QR configurations predicting**
706 **exceedance probabilities of the 10th, 25th, 75th, and 90th percentile: the configuration using the**
707 **transformed forecast as the only independent variable [NQT fcst]; the best performing combination**
708 **for each river gage (upper performance limit) [Best combis]; rates of rise in the past 24 and 48**

709 hours and the forecast errors 24 and 48 hours ago as independent variable (one-size-fits-all
710 solution) [rr+err24/48].

711

712 **Figure 17: Brier Skill Scores of the original-forecast-only QR model-configuration (i.e., using the**
713 **transformed forecast as the only independent variable) for four lead times and flood stages.**

714 **Figure 18: Brier Skill Scores for four lead times and flood stages of observed water levels using the**
715 **best variable-combination-joint predictor for each river gage as independent variables in the QR**
716 **model-configuration.**

717 Figure 19: Empirical cumulative density functions of three QR configurations predicting
718 exceedance probabilities of the Action, Minor, Moderate, and Major Flood Stage: the configuration
719 using the transformed forecast as the only independent variable [NOT fcst]; the best performing
720 combination for each river gage (upper performance limit) [Best combis]

721

722 The fact that the Brier Score can be de-composed into reliability, resolution and
723 uncertainty allows a closer look at which improvements are being achieved by including more
724 variables-predictors than just the forecast. Figure-18Figure 20 shows that the original-forecast-
725 only QR model-configuration as studied by Weerts et al. (2011)(2011) has high reliability (i.e.,
726 the reliability is close to zero). The Brier Score and the Brier Skill Score mainly improve when
727 using rise-rates-rates of rise and forecast errors as independent variables, because the resolution
728 increases. This confirms the finding by Wood et al. (2009) that QR error models should be based
729 on rate of rise (as well as lead time). The forecast quality improves along other dimensions
730 metrics as well, i.e., the areas under the ROC curves and the ranked probability skill score
731 (RPSS) increase. The first weighs missed alarms against false alarms and has a perfect score
732 equal to one. The latter is a version of the Brier Skill Score. While the Brier Skill Score pertains

733 to a binary event, the RPSS can take into account various event categories. Its perfect score
734 equals one (e.g., WWRP/WGNE, 2009)(~~e.g., WWRP/WGNE, 2009~~).

735 **Figure 1820:** Comparison of the ~~original forecast-only~~ QR model configuration (i.e., only
736 transformed forecast as independent variables) and the one-size-fits-all approach (i.e., ~~rise~~
737 ~~rates~~rates of rise and forecast errors as independent variables) using various measures of forecast
738 quality: Brier Score (BS), Brier Skill Score (BSS), Reliability (Rel), Resolution (Res), Uncertainty
739 (Unc), Area under the ROC curve (ROCA), ranked probability score (RPS), ranked probability
740 skill score (RPSS). Lead time: 3 days; 75th percentile of observation levels as threshold. The left
741 figure zooms in on the right figure to make changes in reliability and resolution better visible.

742 3.3 Robustness

743 The impact of the length of the training dataset on the model configuration's performance
744 measured by the Brier Skill Score (BSS) was assessed for the one-size-fits-all QR
745 model configuration (i.e., ~~rise~~ratesrates of rise and forecast errors as independent variables for all
746 gages) for Hardin and Henry on the Illinois River. We were particularly interested in testing how
747 many years of training data are necessary to achieve satisfactory forecasting results. Each year
748 between 2003 and 2013 was forecast by QR model configurations trained ~~on on one year up to~~
749 however many years of archived forecasts were ~~available~~available in that year, i.e., the forecasts
750 for 2005 is produced by a model trained on less data than those for 2013. Then, the BSS for that
751 year (e.g., 2005 or 2013) was computed.

752 Figure 21 and Figure 22 show that training datasets shorter than three years result in very
753 low BSSs for low event thresholds (Q10) at Henry and Hardin, show that for those gages, For the
754 other event thresholds, it ~~does not~~barely matters for the BSS how many years are included in the
755 training dataset. That is good new news, if stationarity cannot be assumed (Milly et al.,
756 2008)(~~Milly et al., 2008~~), a step-change in river regime has occurred, or forecast data have not

757 | been archived in the past. In those cases, only short training datasets are available. Only needing
758 | short time series to define a skillful QR configuration implies that the configuration parameters
759 | can be updated regularly. This way, changing relationships between predictors etc. can be taken
760 | into account.

761 | -However, the BSS varies considerably for what year is being forecast. The forecast
762 | performance varies greatly, especially for the 10th and 25th percentile of observed water levels. It
763 | is likely, that a very large dataset, including more infrequent events, would improve these results.
764 | However, most river forecast centers only recently started archiving forecasts in a text-format, so
765 | that even having ten years' worth of data is an exception. To illustrate that point, the National
766 | Climatic Data Center has archived data from 2001 onwards available in their HDSS Access
767 | System. ²²

768 | To generalize the result, the same analysis as just described for Hardin and Henry was
769 | repeated for all 82 gages. Following that, a regression analysis was executed with the BSS score
770 | as the dependent variable and the river gages and forecast years as factorial independent
771 | variables and the lead time, event thresholds, and number of training years as numerical
772 | independent variables (Table 2). The forecast performance was found to vary statistically
773 | significantly across all those dimensions except the number of training years. This results in a
774 | very wide range of Brier Skill Scores (Figure 22). Accordingly, for the user, it is particularly
775 | difficult to know how much to trust a forecast, if the performance depends so much on context.
776 | Likewise, this is case for the QR configuration based on the forecast only (not shown).

²² ~~To illustrate that point, the National Climatic Data Center has archived data from 2001 onwards available in their HDSS Access System.~~

777 A closer look at the regression coefficients (Table 2) provides interesting insights. For
778 low event thresholds, the BSSs are much worse than for high thresholds. The QR configurations
779 might be performing less well for low event thresholds, because the variance in the dependent
780 variable – the forecast error – is smaller. After all, river forecasts have much smaller errors for
781 lower water levels. The illustrative cases of Henry and Hardin, described above, indicate that
782 using longer time series to predict exceedance probabilities of low event thresholds improves
783 forecast performance.

784 As expected, the BSSs slightly decrease with lead time. Regarding the forecast quality for
785 each forecast year, the regression is slightly biased. The earlier years are included less often in
786 the dataset with on average less years' worth of data in their training dataset, because, for
787 example, unlike for the year 2013, ten years of training data were not available for the year 2006.
788 Nonetheless, the regression indicates that 2008 was particularly difficult to forecast and 2012
789 relatively easy, i.e., they are associated with relatively low and high coefficients respectively
790 (Table 2).

791 The performance of the forecast additionally depends on the river gage. The coefficients
792 of the river gages, included as factors in the regression, have been excluded from Table 2 for the
793 sake of brevity. Instead, Figure 23 maps the geographic position of the river gages with the color
794 code indicating each gage's regression coefficient. The coefficients are lower, and therefore the
795 Brier Skill Scores are lower, for gages far upstream a river and those close to confluences. At
796 least for the gages at confluences, the QR model could probably be improved by including the
797 rise rates at the river gages on the other joining river into the regression.

798

799 **Figure-2119:** Brier Skill Score for various forecast years and various sizes of training dataset across
800 different lead times (colors) and event thresholds (plots) for Hardin, IL (HARI2). The filled-in end
801 point of each line indicates the BSS for the forecast year on the x-axis with one year in the training
802 dataset. Each point further to the left stands for one additional training year for that same forecast
803 year.

804 **Figure -2220:** Brier Skill Score for various forecast years and various sizes of training dataset
805 across different lead times (colors) and event thresholds (plots) for Henry, IL (HNYI2). The filled-
806 in end point of each line indicates the BSS for the forecast year on the x-axis with one year in the
807 training dataset. Each point further to the left stands for one additional training year for that same
808 forecast year.

809 **Figure 2123:** Geographical position of rivers. Colors indicate the regression coefficient of each
810 station with the Brier Skill Score as dependent variable.

811 **Figure 2224:** Minimum (black) and maximum (red) Brier Skill Scores for various lead times and
812 event thresholds across locations, size of training dataset and forecast years.

813 4 Conclusion

814 In this study, quantile regression (QR) has been applied to estimate the probability of the river
815 water level exceeding various event thresholds (i.e., 10th, 25th, 75th, 90th percentiles of observed
816 water levels as well as the four flood stages of each river gage). ~~This is the first study applying
817 this method to the U.S. American context. Additionally, it~~It further develops the ~~method
818 application of QR to estimating river forecast uncertainty by (a) including more comparing
819 different sets of independent variables, (b) and testing the method~~technique's robustness across
820 locations, lead times, event thresholds, forecast years and sizes of training dataset.

821 _____
822 ~~Most importantly~~When compared to the configuration using only the forecast, it was found that
823 including ~~rise rates~~rates of rise in the past 24 and 48 hours and the forecast errors of 24 and 48
824 hours ago as independent variables improves the performance of the QR ~~model~~configuration, as

825 measured by the Brier Skill Score. This confirms Wood et al.'s (2009) finding that QR error
826 models should be a function of rate of rise and lead time. Since the reliability was already high
827 with the original QR method as proposed by Weerts et al., The configuration with the forecast as
828 the only independent variable, as studied by Weerts et al. (2011), produced estimates with high
829 reliability. Including the other four predictors mentioned above mainly~~the new configuration~~
830 ~~mainly~~ increases the resolution.

831 For extremely high water levels, the combinations of independent variables that perform best
832 vary across stations. On those days, combinations of fewer independent variables perform better
833 than those that include more. The most likely explanation is that QR configurations based on
834 large joint predictors result in overfitting the data. In contrast to these extremely high event
835 thresholds, larger sets of variables-predictors work better than smaller ones for non-extreme and
836 low event thresholds. Additionally, customizing the set of predictors to the event thresholds does
837 not improve the BSS much. a one-size-fits-all approach (i.e. the rise rates and forecasts errors as
838 independent variables) performs satisfactorily for those cases.

839 When forming a joint predictor, the independent variables rates of rise and forecast errors do
840 not combine well with the forecast itself. To account for heteroscedasticity, the forecast was
841 transformed into the Gaussian domain. However, no trend is detectable anymore between
842 forecast error and the rates of rise or the previous forecast errors after applying NQT to those
843 variables. Therefore, it is difficult to combine these two predictors. A possible solution could be
844 to define QR configurations for subsets of the transformed data. However, such an approach
845 drastically decreases the amount of data available for each configuration.

846

847 ~~The new independent variables—rise rates and forecast errors—do not combine well with~~
848 ~~forecast itself. The latter was the only variable included in the original QR configuration as~~
849 ~~studied by Weerts et al. and López-López et al.. To account for heteroscedasticity, the forecast~~
850 ~~was transformed into the Gaussian domain. However, the rise rates and the forecast errors do not~~
851 ~~lend themselves for linear quantile regression after such a transformation. Therefore, it is~~
852 ~~difficult to combine these two variables. A possible solution could be to build regression models~~
853 ~~for subsets of the transformed data. However, such an approach drastically decreases the amount~~
854 ~~of data available for each model.~~

855 The ~~proposed-studied QR method-configurations~~ areis relatively robust to the size of training
856 dataset, which is convenient if stationarity cannot be assumed (Milly et al., 2008)(~~Milly et al.,~~
857 ~~2008~~), a step-change in the river regime has occurred, or – as is the case for most river forecast
858 centers – only recent forecast data have been archived. However, the performance of the
859 ~~methodtechnique~~ does dependdepends heavily on the river gage, the lead time, event threshold
860 and year that are being forecast. This results in a very wide range of Brier Skill Scores. This
861 means that the danger remains that forecast users make good experiences with a forecast one
862 year or at one location and assume it is equally reliable in other locations and every year. As is
863 the case with most other forecasts, an indication of forecast uncertainty needs to be
864 communicated alongside the exceedance probabilities generated by our approach.

865 The ~~proposed-studied QR approach-configurations~~ performs less well for longer lead times,
866 for gages far upstream a river or close to confluences, for low event thresholds and extremely
867 high ones. The ~~QR model-configurations~~ might be performing less well for low event thresholds,
868 because the variance in the dependent variable – the forecast error – is smaller. After all, river

869 forecasts have much smaller errors for lower water levels. In turn, for extremely high water
870 levels, the scarcity of data decreases the ~~model configuration's~~ performance.

871 *Future Work*

872 ~~This method techniques~~ can be further developed in several ways to achieve higher Brier Skill
873 Scores and more robustness. First, more independent variables can be added. Trials with a
874 different ~~method technique~~, classification trees, showed that the observed precipitation, the
875 precipitation forecast (i.e., POP – probability of precipitation) and the upstream water levels
876 significantly improve ~~models forecasting performance~~. Presumably, this is the case, because the
877 ~~QPF-forecast used in this study~~ includes the precipitation forecast ~~only~~ for only the next 12
878 hours. However, currently, the precipitation data and forecasts can only be requested in chunks
879 of a month, three chunks per day, from the NCDC's HDSS Access System.²³ For a period of 12
880 years, requesting such data for several weather stations²⁴ is obviously time-consuming; n-ot
881 least, because the geographical units of the weather forecasts bulletins do not correspond with
882 those of the river forecast bulletins. Upstream water levels can easily be included after manually
883 determining the upstream gage(s) for each of the 82 NCRFC gages. To improve ~~model~~
884 performance at gages close to river confluences, the upstream water level of the gages on the
885 joining river should be included as well.

886 Different approaches of sub-setting the data to improve ~~models results performance~~ also
887 warrant consideration. Particularly, clustering the data by variability seems promising. However,
888 early trials indicated that this ~~method technique~~ is very sensitive to the training dataset.

²³ ~~URL [accessed July 2014]:~~

~~<http://edo.ncdc.noaa.gov/pls/plhas/HAS.FileAppSelect?datasetname=9957ANX>~~

²⁴ ~~The geographical units of the weather forecasts bulletins do not correspond with those of the river forecast bulletins.~~

889 As mentioned above, the QR ~~method~~approach works less well for low than for high event
890 thresholds. Further study should investigate, why that is the case, and identify possible solutions.
891 The current study focused on extremely high event thresholds, i.e., flood stages, but not on lower
892 ones, i.e., below the 50th percentile of observed water levels.

893 ~~Last~~Additionally, the ~~proposed-studied method~~technique would need to be verified for gages
894 for which the NCRFC does not publish daily forecasts. Ignorance of the uncertainty inherent in
895 river forecasts ~~have~~has had some of the most unfortunate impacts on decision-making in Grand
896 Forks, ND and Fargo, ND (~~Pielke, 1999; Morss, 2010~~)(~~Pielke, 1999; Morss, 2010~~). Both of those
897 stages are discontinuously forecast NCRFC gages.

898 Finally, this paper uses a brute force approach by simply calculating and comparing all
899 possible combinations of independent variables. Mathematically more challenging stepwise
900 quantile regression would not only be more elegant, but also provide better safeguards against
901 overfitting the data.

902 *Acknowledgements:*

903 Many thanks to Grant Weller who suggested looking into quantile regression to predict forecast
904 errors. We would like to thank the two reviewers for their insightful comments. The paper
905 greatly benefitted from their comments. As to funding, Frauke Hoss is supported by an ERP
906 fellowship of the German National Academic Foundation and by the Center of Climate and
907 Energy Decision Making (SES-0949710), through a cooperative agreement between the National
908 Science Foundation and Carnegie Mellon University (CMU).~~To ensure anonymity, this section~~
909 ~~will be added after the review process.~~

References

[Alexander, M., Harding, M. and Lamarche, C.: Quantile Regression for Time-Series-Cross-Section-Data, Int. J. Stat. Manag. Syst., 4\(1-2\), 47–72, 2011.](#)

[Bogner, K., Pappenberger, F. and Cloke, H. L.: Technical Note: The normal quantile transformation and its application in a flood forecasting system, Hydrol. Earth Syst. Sci., 16\(4\), 1085–1094, doi:10.5194/hess-16-1085-2012, 2012.](#)

[Brier, G. W.: Verification of Forecasts Expressed in Terms of Probability, Mon. Weather Rev., 78\(1\), 1–3, doi:10.1175/1520-0493\(1950\)078<0001:VOFEIT>2.0.CO;2, 1950.](#)

[Brown, J. D. and Seo, D.-J.: Evaluation of a nonparametric post-processor for bias correction and uncertainty estimation of hydrologic predictions, Hydrol. Process., 27\(1\), 83–105, doi:10.1002/hyp.9263, 2013.](#)

[Demargne, J., Wu, L., Regonda, S. K., Brown, J. D., Lee, H., He, M., Seo, D.-J., Hartman, R., Herr, H. D., Fresch, M., Schaake, J. and Zhu, Y.: The Science of NOAA’s Operational Hydrologic Ensemble Forecast Service, Bull. Am. Meteorol. Soc., 95\(1\), 79–98, doi:10.1175/BAMS-D-12-00081.1, 2013.](#)

[Hsu, W. and Murphy, A. H.: The attributes diagram A geometrical framework for assessing the quality of probability forecasts, Int. J. Forecast., 2\(3\), 285–293, doi:10.1016/0169-2070\(86\)90048-8, 1986.](#)

[Ikeda, M., Ishigaki, T. and Yamauchi, K.: Relationship between Brier score and area under the binormal ROC curve, Comput. Methods Programs Biomed., 67\(3\), 187–194, doi:10.1016/S0169-2607\(01\)00157-2, 2002.](#)

Illinois Department of Natural Resources: Aquatic Illinois - Illinois Rivers and Lakes Fact Sheets, [online] Available from:
<http://dnr.state.il.us/education/aquatic/aquaticillinoisrivlakefactshts.pdf> (Accessed 3 February 2015), 2011.

Jolliffe, I. T. and Stephenson, D. B.: Forecast Verification: A Practitioner's Guide in Atmospheric Science, John Wiley & Sons., 2012.

Kelly, K. S. and Krzysztofowicz, R.: A bivariate meta-Gaussian density for use in hydrology, Stoch. Hydrol. Hydraul., 11(1), 17–31, doi:10.1007/BF02428423, 1997.

Koenker, R.: Quantile Regression, Cambridge University Press., 2005.

Koenker, R.: quantreg: Quantile Regression, R Package Version 505 [online] Available from: <http://CRAN.R-project.org/package=quantreg> (Accessed 27 August 2014), 2013.

Koenker, R. and Bassett, G.: Regression Quantiles, Econometrica, 46(1), 33, doi:10.2307/1913643, 1978.

Koenker, R. and Machado, J. A. F.: Goodness of Fit and Related Inference Processes for Quantile Regression, J. Am. Stat. Assoc., 94(448), 1296–1310, doi:10.1080/01621459.1999.10473882, 1999.

Leahy, C. P.: Objective Assessment and Communication of Uncertainty in Flood Warnings., 2007.

López López, P., Verkade, J. S., Weerts, A. H. and Solomatine, D. P.: Alternative configurations of Quantile Regression for estimating predictive uncertainty in water level forecasts for the Upper Severn River: a comparison, Hydrol. Earth Syst. Sci. Discuss., 11(4), 3811–3855, 2014.

Milly, P. C. D., Betancourt, J., Falkenmark, M., Hirsch, R. M., Kundzewicz, Z. W., Lettenmaier, D. P. and Stouffer, R. J.: Stationarity Is Dead: Whither Water Management?, Science, 319(5863), 573–574, doi:10.1126/science.1151915, 2008.

Montanari, A. and Brath, A.: A stochastic approach for assessing the uncertainty of rainfall-runoff simulations, Water Resour. Res., 40(1), W01106, doi:10.1029/2003WR002540, 2004.

Montanari, A. and Grossi, G.: Estimating the uncertainty of hydrological forecasts: A statistical approach, Water Resour. Res., 44(12), W00B08, doi:10.1029/2008WR006897, 2008.

Morss, R. E.: Interactions among Flood Predictions, Decisions, and Outcomes: Synthesis of Three Cases, Nat. Hazards Rev., 11(3), 83–96, doi:10.1061/(ASCE)NH.1527-6996.0000011, 2010.

National Climatic Data Center: HDSS Access System, [online] Available from: <http://cdo.ncdc.noaa.gov/pls/plhas/HAS.FileAppSelect?datasetname=9957ANX>; (Accessed 15 July 2014), 2014.

National Research Council: Completing the Forecast: Characterizing and Communicating Uncertainty for Better Decisions Using Weather and Climate Forecasts, National Academies Press, Washington, DC. [online] Available from: http://www.nap.edu/catalog.php?record_id=11699 (Accessed 18 September 2014), 2006.

Pielke, R. A.: Who Decides? Forecasts and Responsibilities in the 1997 Red River Flood, Appl. Behav. Sci. Rev., 7(2), 83–101, 1999.

Regonda, S. K., Seo, D.-J., Lawrence, B., Brown, J. D. and Demargne, J.: Short-term ensemble streamflow forecasting using operationally-produced single-valued streamflow forecasts – A Hydrologic Model Output Statistics (HMOS) approach, J. Hydrol., 497, 80–96, doi:10.1016/j.jhydrol.2013.05.028, 2013.

Seo, D. J.: Hydrologic Ensemble Processing Overview, [online] Available from: http://www.nws.noaa.gov/oh/hrl/hymb/docs/hep/events_announce/Hydro_Ens_Overview_DJ.pdf (Accessed 29 January 2015), 2008.

Seo, D.-J., Herr, H. D. and Schaake, J. C.: A statistical post-processor for accounting of hydrologic uncertainty in short-range ensemble streamflow prediction, Hydrol Earth Syst Sci Discuss, 3(4), 1987–2035, doi:10.5194/hessd-3-1987-2006, 2006.

Solomatine, D. P. and Shrestha, D. L.: A novel method to estimate model uncertainty using machine learning techniques, Water Resour. Res., 45, doi:10.1029/2008WR006839, 2009.

USGS: Stream Site - USGS 05558300 Illinois River at Henry, IL, [online] Available from: http://waterdata.usgs.gov/nwis/inventory/?site_no=05558300&agency_cd=USGS (Accessed 2 February 2015a), 2015.

USGS: Stream Site - USGS 05587060 Illinois River at Hardin, IL, [online] Available from: http://waterdata.usgs.gov/il/nwis/inventory/?site_no=05587060& (Accessed 3 February 2015b), 2015.

Weerts, A. H., Winsemius, H. C. and Verkade, J. S.: Estimation of predictive hydrological uncertainty using quantile regression: examples from the National Flood Forecasting System (England and Wales), Hydrol Earth Syst Sci, 15(1), 255–265, doi:10.5194/hess-15-255-2011, 2011.

Welles, E., Sorooshian, S., Carter, G. and Olsen, B.: Hydrologic Verification: A Call for Action and Collaboration, Bull. Am. Meteorol. Soc., 88(4), 503–511, doi:10.1175/BAMS-88-4-503, 2007.

Wikipedia: Brier score, [online] Available from:

http://en.wikipedia.org/w/index.php?title=Brier_score&oldid=619686224 (Accessed 27 August 2014), 2014.

Wilson, L. J.: Verification of probability and ensemble forecasts, [online] Available from:

http://www.swpc.noaa.gov/forecast_verification/Assets/Tutorials/Ensemble%20Forecast%20Verification.pdf (Accessed 27 August 2014), n.d.

Wood, A. W., Wiley, M. and Nijssen, B.: Use of quantile regression for calibration of hydrologic forecasts, [online] Available from:

<http://ams.confex.com/ams/89annual/wrfredirect.cgi?id=10049>, 2009.

WWRP/WGNE: Methods for probabilistic forecasts. Forecast Verification – Issues, Methods and FAQ, [online] Available from:

http://www.cawcr.gov.au/projects/verification/verif_web_page.html#BSS (Accessed 27 August 2014), 2009.

Alexander, M., Harding, M. and Lamarche, C.: Quantile Regression for Time-Series Cross-Section Data, Int. J. Stat. Manag. Syst., 4(1-2), 47-72, 2011.

Anon: Brier score, Wikipedia Free Encycl. [online] Available from:

http://en.wikipedia.org/w/index.php?title=Brier_score&oldid=619686224 (Accessed 27 August 2014), 2014.

Bogner, K., Pappenberger, F. and Cloke, H. L.: Technical Note: The normal quantile transformation and its application in a flood forecasting system, Hydrol. Earth Syst. Sci., 16(4), 1085-1094, doi:10.5194/hess-16-1085-2012, 2012.

Brier, G. W.: Verification of Forecasts Expressed in Terms of Probability, Mon. Weather Rev., 78(1), 1-3, doi:10.1175/1520-0493(1950)078<0001:VOFEIT>2.0.CO;2, 1950.

~~Bröcker, J.: Estimating reliability and resolution of probability forecasts through decomposition of the empirical score, *Clim. Dyn.*, 39(3–4), 655–667, doi:10.1007/s00382-011-1191-1, 2012.~~

~~Demargne, J., Wu, L., Regonda, S. K., Brown, J. D., Lee, H., He, M., Seo, D. J., Hartman, R., Herr, H. D., Fresch, M., Schaake, J. and Zhu, Y.: The Science of NOAA's Operational Hydrologic Ensemble Forecast Service, *Bull. Am. Meteorol. Soc.*, 95(1), 79–98, doi:10.1175/BAMS-D-12-00081.1, 2013.~~

~~Ferro, C. a. T. and Fricker, T. E.: A bias-corrected decomposition of the Brier score, *Q. J. R. Meteorol. Soc.*, 138(668), 1954–1960, doi:10.1002/qj.1924, 2012.~~

~~Hsu, W. and Murphy, A. H.: The attributes diagram A geometrical framework for assessing the quality of probability forecasts, *Int. J. Forecast.*, 2(3), 285–293, doi:10.1016/0169-2070(86)90048-8, 1986.~~

~~Ikeda, M., Ishigaki, T. and Yamauchi, K.: Relationship between Brier score and area under the binormal ROC curve, *Comput. Methods Programs Biomed.*, 67(3), 187–194, doi:10.1016/S0169-2607(01)00157-2, 2002.~~

~~Jolliffe, I. T. and Stephenson, D. B.: *Forecast Verification: A Practitioner's Guide in Atmospheric Science*, John Wiley & Sons., 2012.~~

~~Kelly, K. S. and Krzysztofowicz, R.: A bivariate meta-Gaussian density for use in hydrology, *Stoch. Hydrol. Hydraul.*, 11(1), 17–31, doi:10.1007/BF02428423, 1997.~~

~~Koenker, R.: *Quantile Regression*, Cambridge University Press., 2005.~~

~~Koenker, R.: *quantreg: Quantile Regression*, R Package Version 505 [online] Available from: <http://CRAN.R-project.org/package=quantreg> (Accessed 27 August 2014), 2013.~~

~~Koenker, R. and Bassett, G.: Regression Quantiles, *Econometrica*, 46(1), 33, doi:10.2307/1913643, 1978.~~

Koenker, R. and Machado, J. A. F.: Goodness of Fit and Related Inference Processes for Quantile Regression, *J. Am. Stat. Assoc.*, 94(448), 1296–1310, doi:10.1080/01621459.1999.10473882, 1999.

Leahy, C. P.: Objective Assessment and Communication of Uncertainty in Flood Warnings., 2007.

López-López, P., Verkade, J. S., Weerts, A. H. and Solomatine, D. P.: Alternative configurations of Quantile Regression for estimating predictive uncertainty in water level forecasts for the Upper Severn River: a comparison, *Hydrol. Earth Syst. Sci. Discuss.*, 11(4), 3811–3855, 2014.

Milly, P. C. D., Betancourt, J., Falkenmark, M., Hirsch, R. M., Kundzewicz, Z. W., Lettenmaier, D. P. and Stouffer, R. J.: Stationarity Is Dead: Whither Water Management?, *Science*, 319(5863), 573–574, doi:10.1126/science.1151915, 2008.

Montanari, A. and Brath, A.: A stochastic approach for assessing the uncertainty of rainfall-runoff simulations, *Water Resour. Res.*, 40(1), W01106, doi:10.1029/2003WR002540, 2004.

Montanari, A. and Grossi, G.: Estimating the uncertainty of hydrological forecasts: A statistical approach, *Water Resour. Res.*, 44(12), W00B08, doi:10.1029/2008WR006897, 2008.

Morss, R. E.: Interactions among Flood Predictions, Decisions, and Outcomes: Synthesis of Three Cases, *Nat. Hazards Rev.*, 11(3), 83–96, doi:10.1061/(ASCE)NH.1527-6996.0000011, 2010.

Morss, R. E., Lazo, J. K. and Demuth, J. L.: Examining the use of weather forecasts in decision scenarios: results from a US survey with implications for uncertainty communication, *Meteorol. Appl.*, 17(2), 149–162, doi:10.1002/met.196, 2010.

National Research Council: Completing the Forecast: Characterizing and Communicating Uncertainty for Better Decisions Using Weather and Climate Forecasts, National Academies

Press, Washington, DC. [online] Available from:
http://www.nap.edu/catalog.php?record_id=11699 (Accessed 18 September 2014), 2006.

NCAR Research Applications Laboratory, N. R. A.: verification: Weather Forecast Verification Utilities. [online] Available from: <http://cran.r-project.org/web/packages/verification/index.html> (Accessed 27 August 2014), 2014.

Pielke, R. A.: Who Decides? Forecasts and Responsibilities in the 1997 Red River Flood, *Appl. Behav. Sci. Rev.*, 7(2), 83–101, 1999.

R Core Team: R: A language and environment for statistical computing., [online] Available from: [http://www.R-project.org/.](http://www.R-project.org/), 2014.

Regonda, S. K., Seo, D. J., Lawrence, B., Brown, J. D. and Demargne, J.: Short-term ensemble streamflow forecasting using operationally-produced single-valued streamflow forecasts—A Hydrologic Model Output Statistics (HMOS) approach, *J. Hydrol.*, 497, 80–96, doi:10.1016/j.jhydrol.2013.05.028, 2013.

Seo, D. J., Herr, H. D. and Schaake, J. C.: A statistical post-processor for accounting of hydrologic uncertainty in short-range ensemble streamflow prediction, *Hydrol Earth Syst Sci Discuss*, 3(4), 1987–2035, doi:10.5194/hessd-3-1987-2006, 2006.

Solomatine, D. P. and Shrestha, D. L.: A novel method to estimate model uncertainty using machine learning techniques, *Water Resour. Res.*, 45, doi:10.1029/2008WR006839, 2009.

Stephenson, D. B., Coelho, C. A. S. and Jolliffe, I. T.: Two Extra Components in the Brier Score Decomposition, *Weather Forecast.*, 23(4), 752–757, doi:10.1175/2007WAF2006116.1, 2008.

Weerts, A. H., Winsemius, H. C. and Verkade, J. S.: Estimation of predictive hydrological uncertainty using quantile regression: examples from the National Flood Forecasting System

(England and Wales), *Hydrol Earth Syst Sci*, 15(1), 255–265, doi:10.5194/hess-15-255-2011, 2011.

Welles, E., Sorooshian, S., Carter, G. and Olsen, B.: Hydrologic Verification: A Call for Action and Collaboration, *Bull. Am. Meteorol. Soc.*, 88(4), 503–511, doi:10.1175/BAMS-88-4-503, 2007.

Wilks, D. S.: *Statistical Methods in the Atmospheric Sciences*, Academic Press, San Diego., 1995.

Wilson, L. J.: Verification of probability and ensemble forecasts, [online] Available from: http://www.swpc.noaa.gov/forecast_verification/Assets/Tutorials/Ensemble%20Forecast%20Verification.pdf (Accessed 27 August 2014), n.d.

WWRP/WGNE: Methods for probabilistic forecasts. Forecast Verification—Issues, Methods and FAQ, [online] Available from: http://www.cawcr.gov.au/projects/verification/verif_web_page.html#BSS (Accessed 27 August 2014), 2009.

Tables

Table 1: ~~Variable Combination~~ Joint predictors

Combi	fcst	err24	err48	rr24	rr48	Combi	fcst	err24	err48	rr24	rr48
1	●					16	●	●	●		
2		●				17	●	●		●	
3			●			18	●	●			●
4				●		19	●		●	●	
5					●	20	●		●		●
6	●	●				21	●			●	●
7	●		●			22		●	●	●	
8	●			●		23		●	●		●
9	●				●	24		●		●	●
10		●	●			25			●	●	●
11		●		●		26	●	●	●	●	
12		●			●	27	●	●	●		●
13			●	●		28	●	●		●	●
14			●		●	29	●		●	●	●
15				●	●	30		●	●	●	●
						31	●	●	●	●	●

fcst = forecast; rr24, rr48 = rise rate of rise in the past 24 and 48 hours;

err24, err 48 = forecast error 24 and 48 hours ago

The forecast error equals the difference between the current (i.e., at issue time of the forecast) water level and the forecast that was produced 24/48 hours ago.

Table 2: Error statistics for the forecast error a) of the whole dataset; b) on days that the water level did not exceed the 10th percentile of observations; c) on days that the water level exceeded the 90th percentile of observations; d) on days that the water level exceeded minor flood stage.

Average errors of 82 gages	Lead-Time					
	Day-1	Day-2	Day-3	Day-4	Day-5	Day-6
a) ALL OBSERVATIONS						
Minimum	-0.21	-0.08	-0.09	-0.07	-0.04	0.02
Median	0.01	0.02	0.06	0.13	0.22	0.30
Mean	0.01	0.04	0.10	0.18	0.30	0.41
Maximum	0.19	0.21	0.76	1.65	2.62	3.47
b) OBSERVATIONS < 10th PERCENTILE						
Minimum	-1.2	-0.35	-0.38	-0.41	-0.38	-0.39
Median	-0.03	-0.04	-0.05	-0.05	-0.04	-0.04
Mean	-0.06	-0.06	-0.06	-0.06	-0.05	-0.04
Maximum	0.03	0.04	0.05	0.12	0.17	0.25
c) OBSERVATIONS > 90th PERCENTILE						
Minimum	-0.11	-0.23	-0.31	-0.38	-0.38	-0.27
Median	-0.01	0.02	0.15	0.32	0.55	0.81
Mean	0.01	0.09	0.29	0.55	0.82	1.14
Maximum	0.34	1.01	3.12	5.13	6.81	8.56
d) OBSERVATIONS > FLOOD STAGE						
Minimum	-0.20	-0.30	-0.44	-0.63	-0.78	-0.80
Median	-0.02	-0.03	0.22	0.45	0.78	1.10
Mean	0.01	0.17	0.45	0.80	1.14	1.51
Maximum	0.65	2.44	5.70	8.37	10.40	11.74

Table 3: Mean and standard deviation three QR configurations: the original using the transformed forecast only as independent variable; the best performing combination for each river gage (upper performance limit); rise rates in the past 24 and 48 hours and the forecast errors 24 and 48 hours ago as independent variable (one-size-fits-all solution).

	Q10	Q25	Q75	Q90	Q10	Q25	Q90
	Day 1				Day 2		
NQT-fest	0.34 (0.52)	0.65 (0.36)	0.90 (0.07)	0.88 (0.08)	0.24 (0.57)	0.59 (0.35)	0.88 (0.08)
Best combi.s	0.54 (0.34)	0.78 (0.18)	0.93 (0.05)	0.91 (0.06)	0.49 (0.36)	0.74 (0.19)	0.91 (0.06)
Rise rate 24/48 +error 24/48*	0.49 (0.41)	0.77 (0.18)	0.92 (0.05)	0.93 (0.06)	0.42 (0.44)	0.73 (0.19)	0.93 (0.06)
	Day 3				Day 4		
NQT-fest	0.20 (0.61)	0.56 (0.33)	0.81 (0.10)	0.75 (0.15)	0.19 (0.55)	0.55 (0.31)	0.75 (0.15)
Best combi.s	0.47 (0.37)	0.74 (0.17)	0.89 (0.05)	0.85 (0.09)	0.46 (0.37)	0.73 (0.18)	0.85 (0.09)
Rise rate 24/48 +error 24/48*	0.40 (0.44)	0.72 (0.19)	0.88 (0.06)	0.84 (0.11)	0.39 (0.43)	0.71 (0.20)	0.84 (0.11)
	Action	Minor	Moderate	Major	Action	Minor	Major
	Day 1				Day 2		
NQT-fest	0.81 (0.27)	0.42 (1.12)	0.38 (1.02)	-0.80 (2.07)	0.68 (0.59)	0.41 (0.90)	0.38 (1.02)
Best combi.s	0.86 (0.26)	0.78 (0.27)	0.73 (0.24)	0.36 (0.66)	0.82 (0.29)	0.73 (0.28)	0.36 (0.66)
	Day 3				Day 4		
NQT-fest	0.67 (0.37)	0.37 (0.87)	-0.09 (1.42)	-1.69 (2.24)	0.62 (0.35)	0.22 (1.00)	-0.09 (1.42)
Best combi.s	0.81 (0.26)	0.71 (0.31)	-0.64 (0.23)	-0.19 (0.76)	0.79 (0.26)	0.69 (0.30)	-0.19 (0.76)

* Combination 30

Table 2: Regression results

	<u>Coef.</u>	<u>St.Dev.</u>	
<u>Intercept</u>	<u>-0.206</u>	<u>0.031</u>	<u>***</u>
<u>Event thresholds</u>	<u>0.265</u>	<u>0.003</u>	<u>***</u>
<u>Lead Times</u>	<u>-0.021</u>	<u>0.003</u>	<u>***</u>
<u>Forecast Years</u>			
<u>2004</u>	<u>-0.266</u>	<u>0.020</u>	<u>***</u>
<u>2005</u>	<u>-0.081</u>	<u>0.018</u>	<u>***</u>
<u>2006</u>	<u>-0.125</u>	<u>0.017</u>	<u>***</u>
<u>2007</u>	<u>-0.129</u>	<u>0.017</u>	<u>***</u>
<u>2008</u>	<u>-0.203</u>	<u>0.017</u>	<u>***</u>
<u>2009</u>	<u>-0.125</u>	<u>0.016</u>	<u>***</u>
<u>2010</u>	<u>-0.140</u>	<u>0.017</u>	<u>***</u>
<u>2011</u>	<u>-0.128</u>	<u>0.016</u>	<u>***</u>

<u>2012</u>	<u>0.056</u>	<u>0.017</u>	<u>***</u>
<u>2013</u>	<u>-0.054</u>	<u>0.016</u>	<u>***</u>
<u>Number of Years in Training Dataset</u>	<u>0.001</u>	<u>0.001</u>	
<u>River Gages</u>			<u>***</u>
<i><u>For the sake of brevity, the 82 river gages included in the regression as factors are omitted here.</u></i>			
<u>R²</u>		<u>0.26</u>	
<u>Adjusted R²</u>		<u>0.25</u>	
<u>P-Values: *** – <0.001; ** – 0.01; * – 0.05; . – 0.1</u>			

Figures

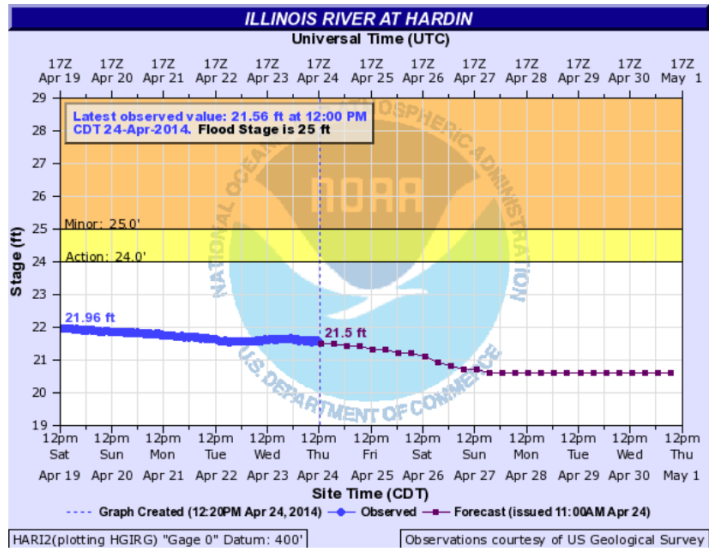


Figure 1: Deterministic short-term weather forecast in six hour intervals as published by the NWS for Hardin, IL on 24 April 2014.

Source:<http://water.weather.gov/ahps2/hydrograph.php?wfo=lsx&gage=hari2>.

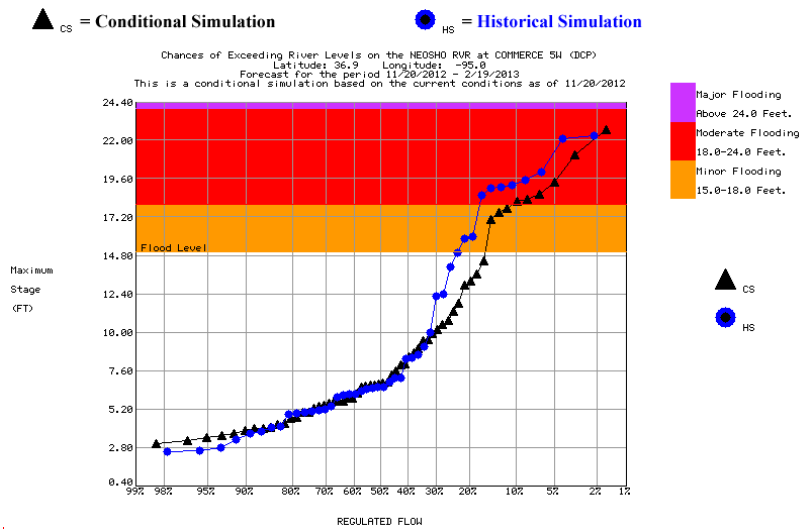


Figure 2: Probabilistic long-term forecast as published by the NWS for Commerce, OK on December 14th, 2012: Exceedance curve for three months period. (Not available for Hardin, IL). Source: <http://water.weather.gov/ahps2/hydrograph.php?wfo=tsa&gage=como2>

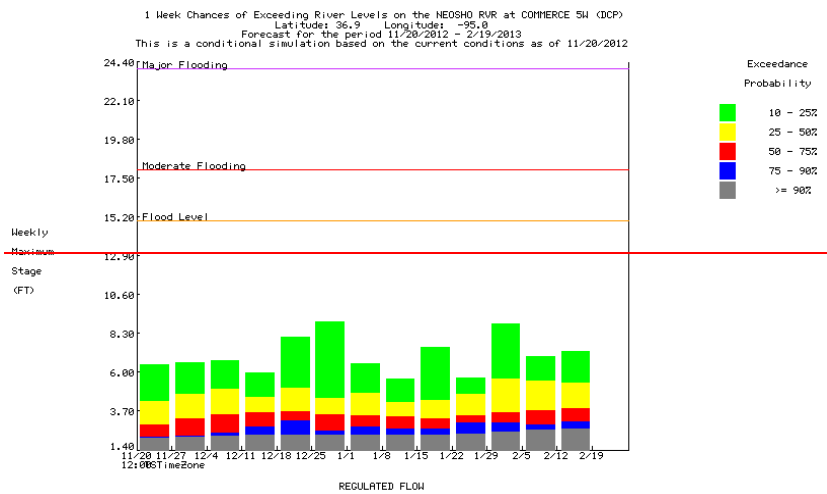


Figure 3: Probabilistic long-term forecast as published by the NWS for Commerce, OK on December 14th, 2012: Bar plot for each week of a three-month period. (Not available for Hardin, IL). Source: <http://water.weather.gov/ahps2/hydrograph.php?wfo=tsa&gage=como2>

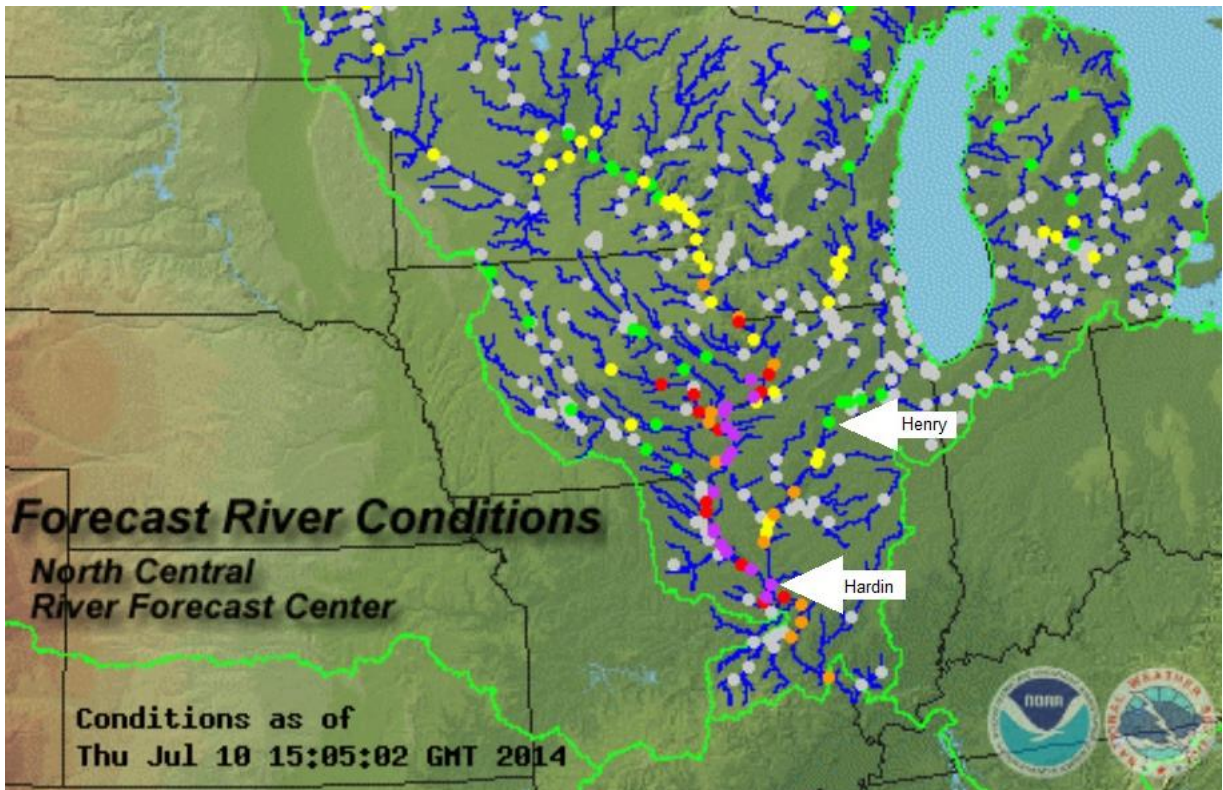
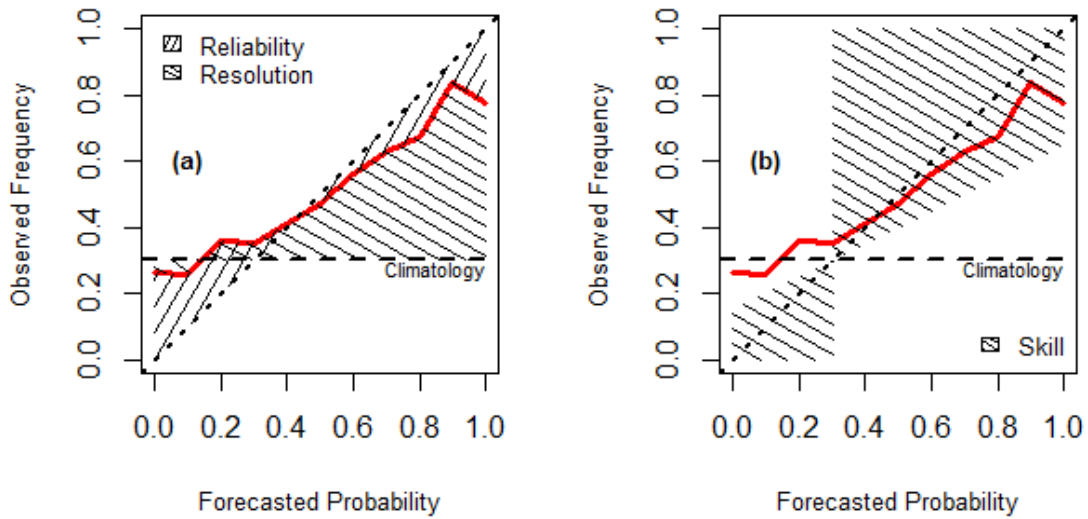


Figure 24: Theory behind Brier Skill Score illustrated for an imaginary forecast (red line): (a) reliability and resolution; (b) skill. In figure a, the area representing reliability should be as small, and for resolution as large as possible. The forecast has skill ($BSS > 0$), i.e., performs better than ~~random-guessing~~ the reference forecast, if it is inside the shaded area in the figure b. Ideally, the forecast would follow the diagonal ($BSS=1$). (Adapted from Hsu and Murphy, 1986; Wilson,

n.d.(Adapted from Hsu and Murphy, 1986; Wilson, n.d.).

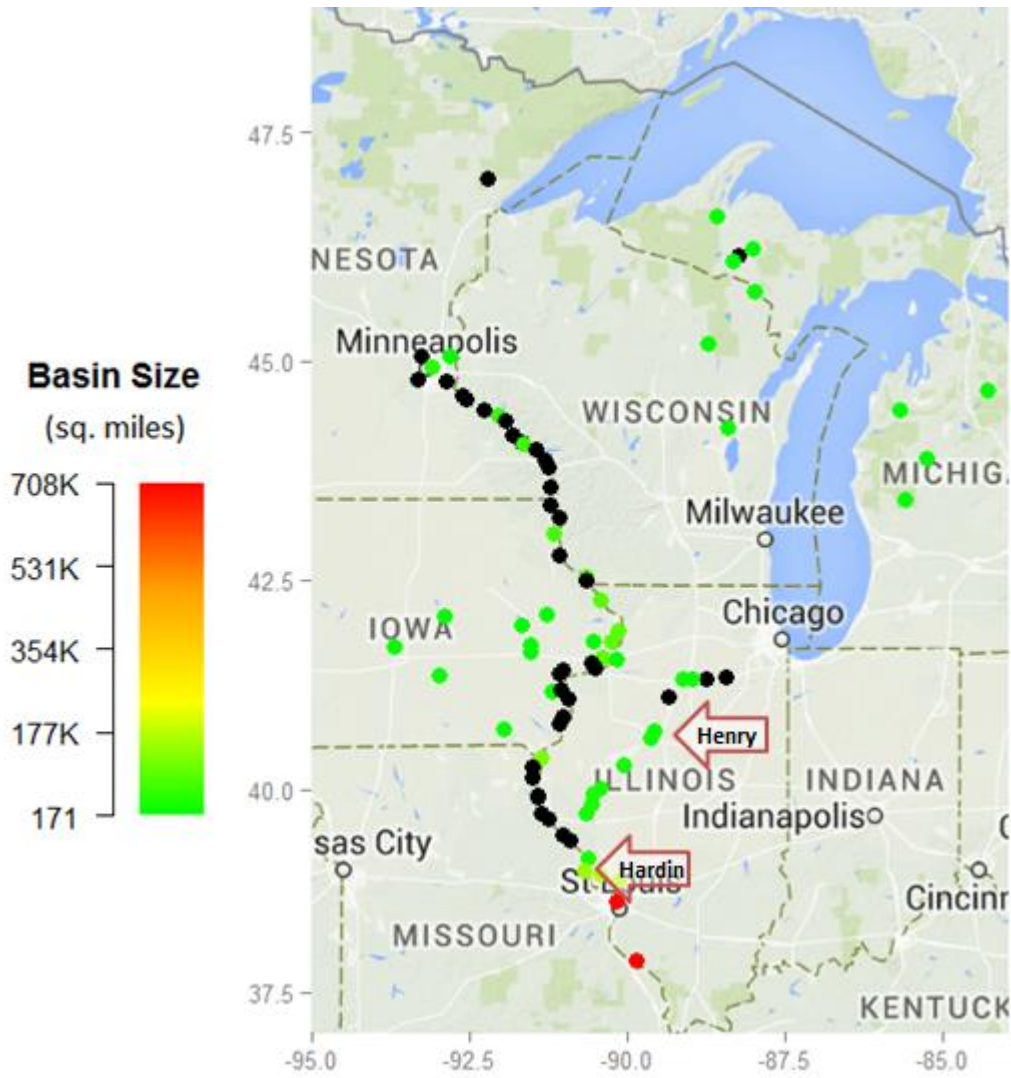


Figure 35: Portion-River gages for which the of-the-North Central River Forecast Centers river gages withpublishes forecasts daily. Henry (HYN12) and Hardin (HARI2) are indicated by the upper and lower red arrow respectively. For gages indicated by black dots the basin size is missing.Source:

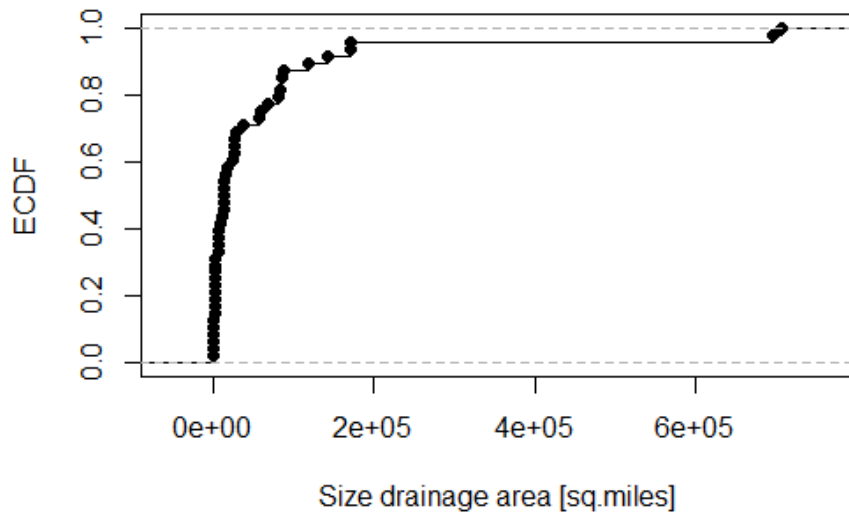


Figure 4: Empirical cumulative density function (ecdf) of sizes of drainage area for the river gages that are being forecasted daily by the NCRFC.

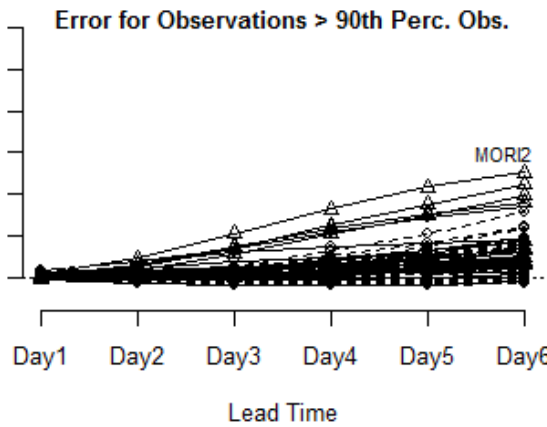
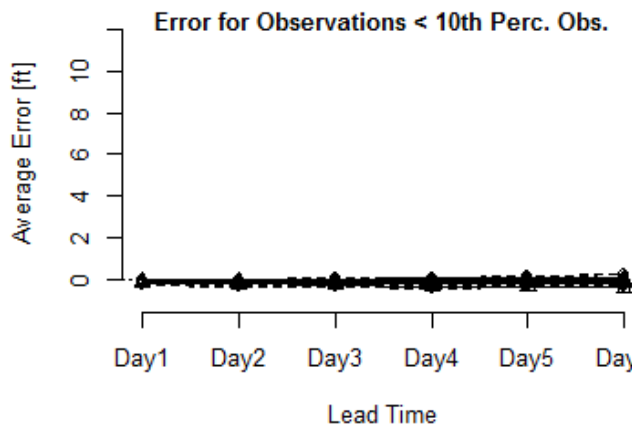
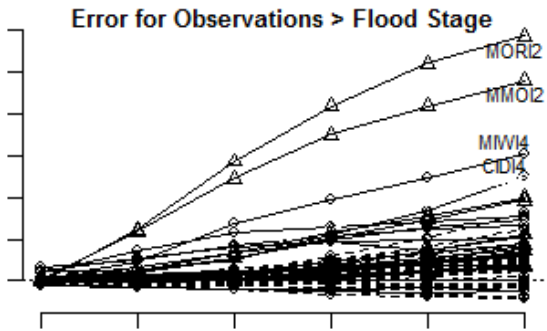
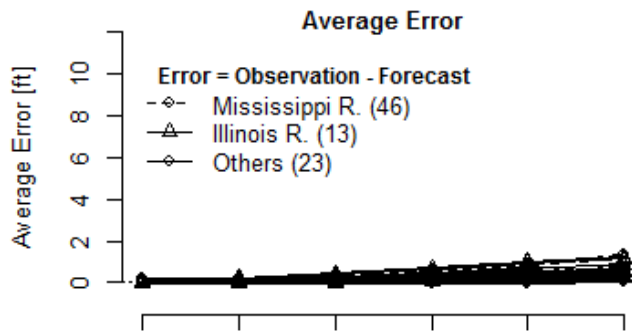
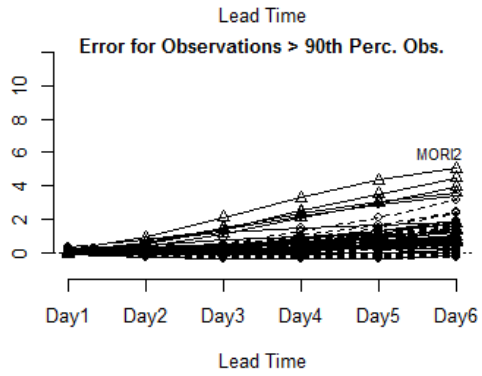
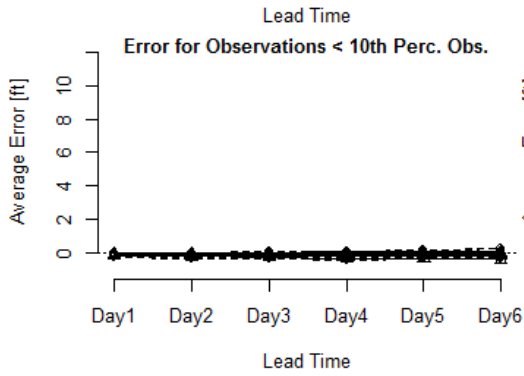
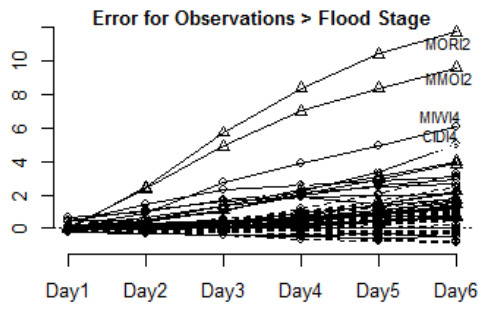
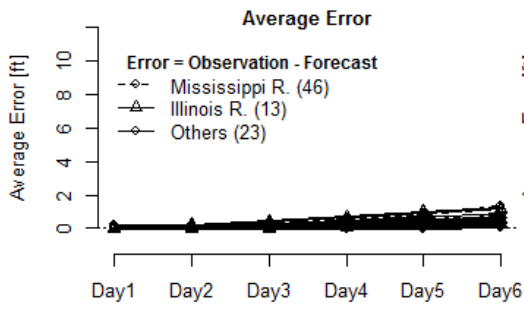


Figure 56: Forecast error for 82 river gages that the NCRFC publishes daily forecasts for. In anti-clockwise direction starting at the top left: (a) Average error; (b) error on days that the water level did not exceed the 10th percentile of observations; (c) error on days that the water level exceeded the 90th percentile of observations; (d) error on days that the water level exceeded minor flood stage.

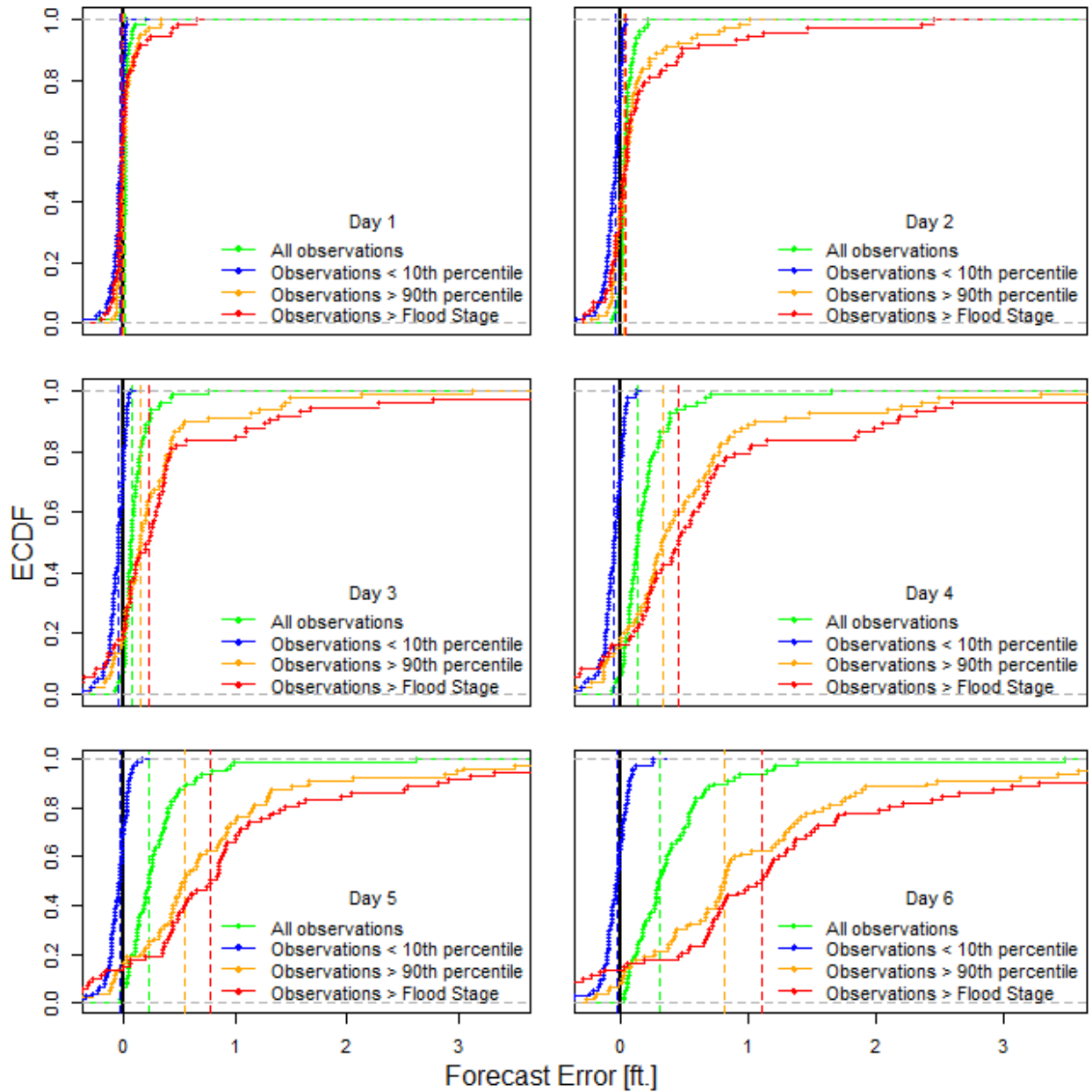
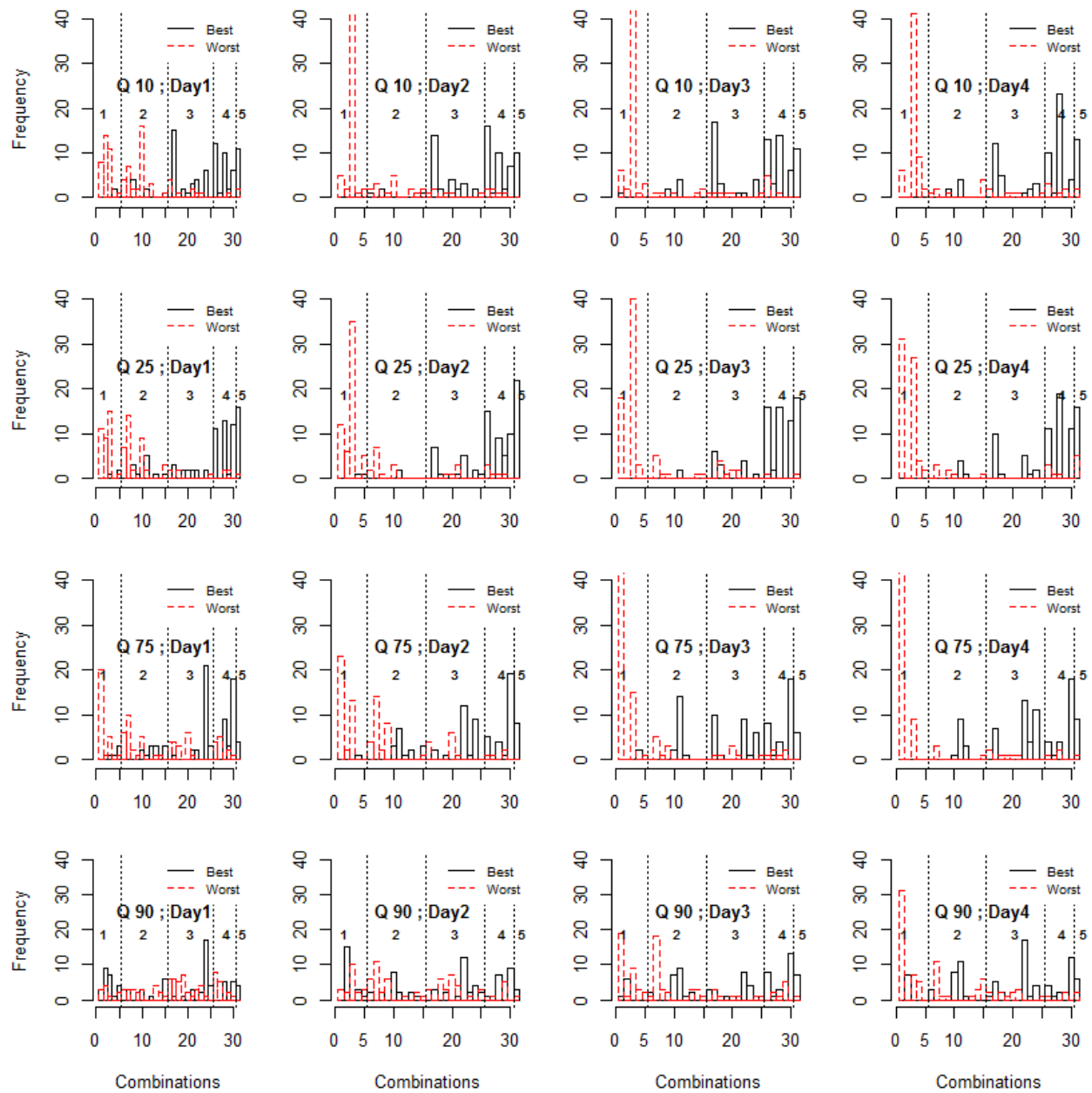


Figure 6: Empirical cumulative distribution function (ecdf) of forecast error at 82 river gages for six lead times. Vertical lines show the median forecast error of the corresponding subset.



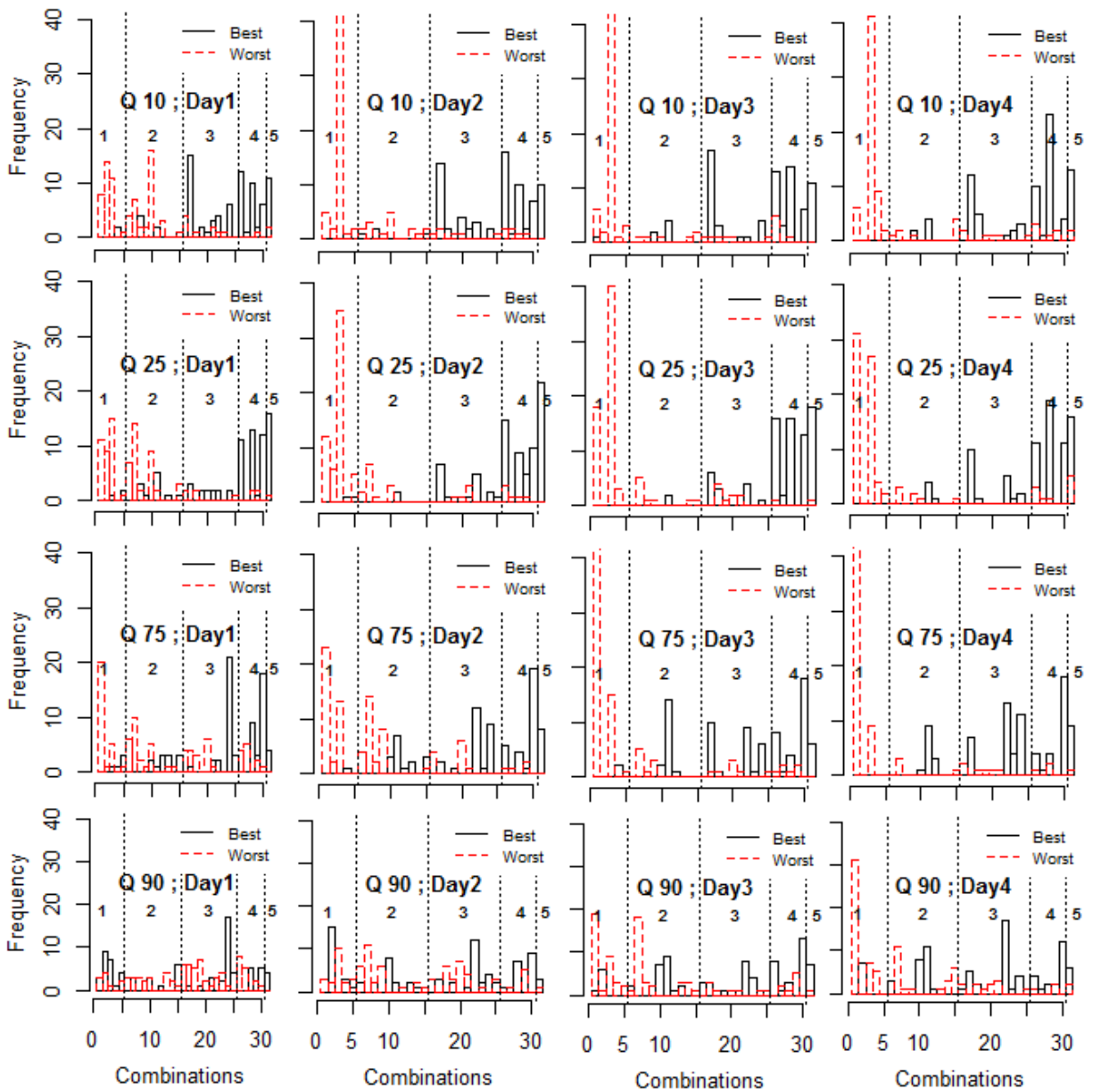
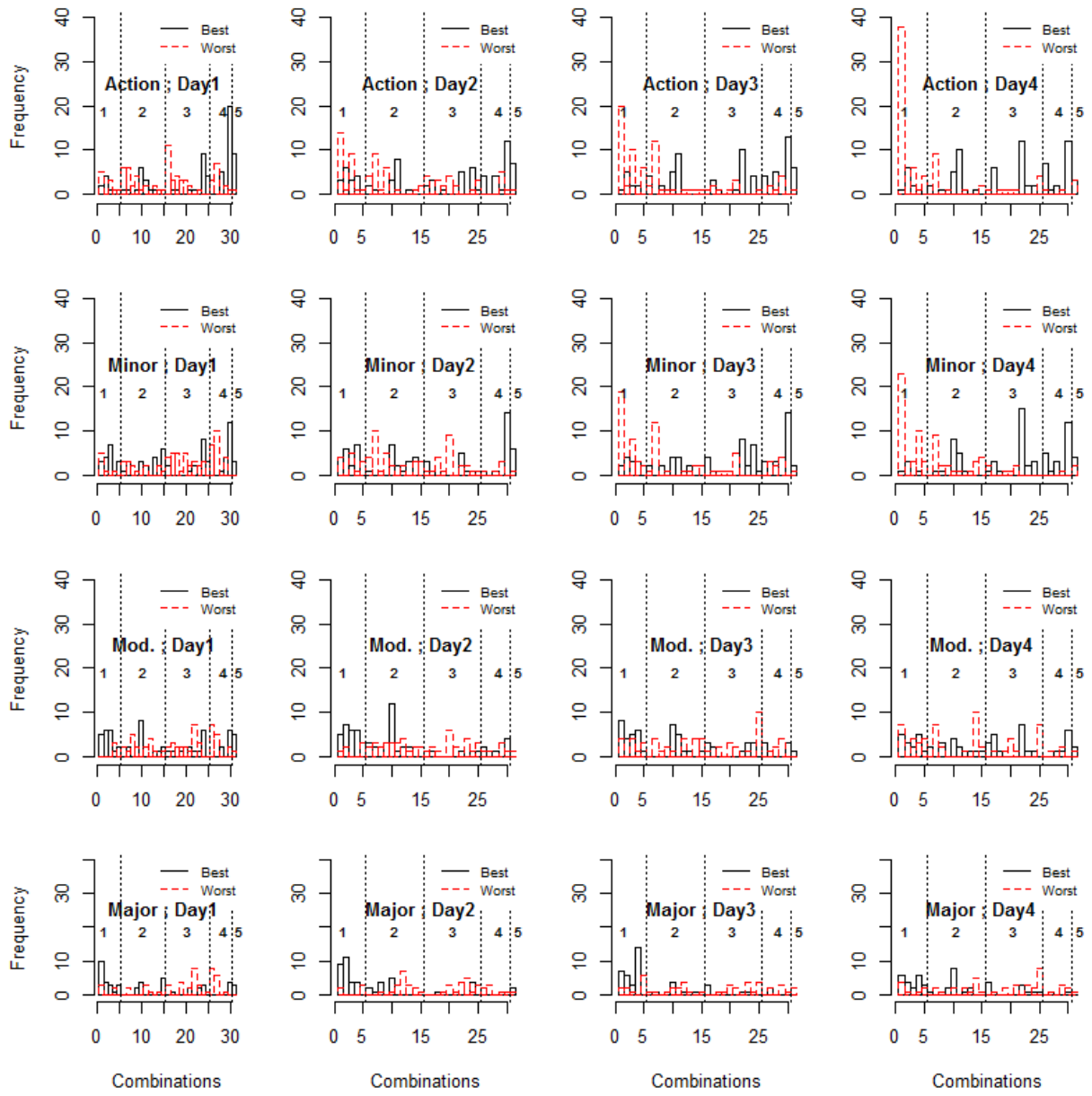


Figure 7: Histograms of variable-combination joint predictors returning the best and worst Brier Skill Scores across 82 river gages. Each row of histograms refers to an event threshold defined as a percentile of the observed water levels, and each column to a lead time. The dotted vertical lines in the histograms distinguish variable-combination joint predictors with different numbers of independent variables.



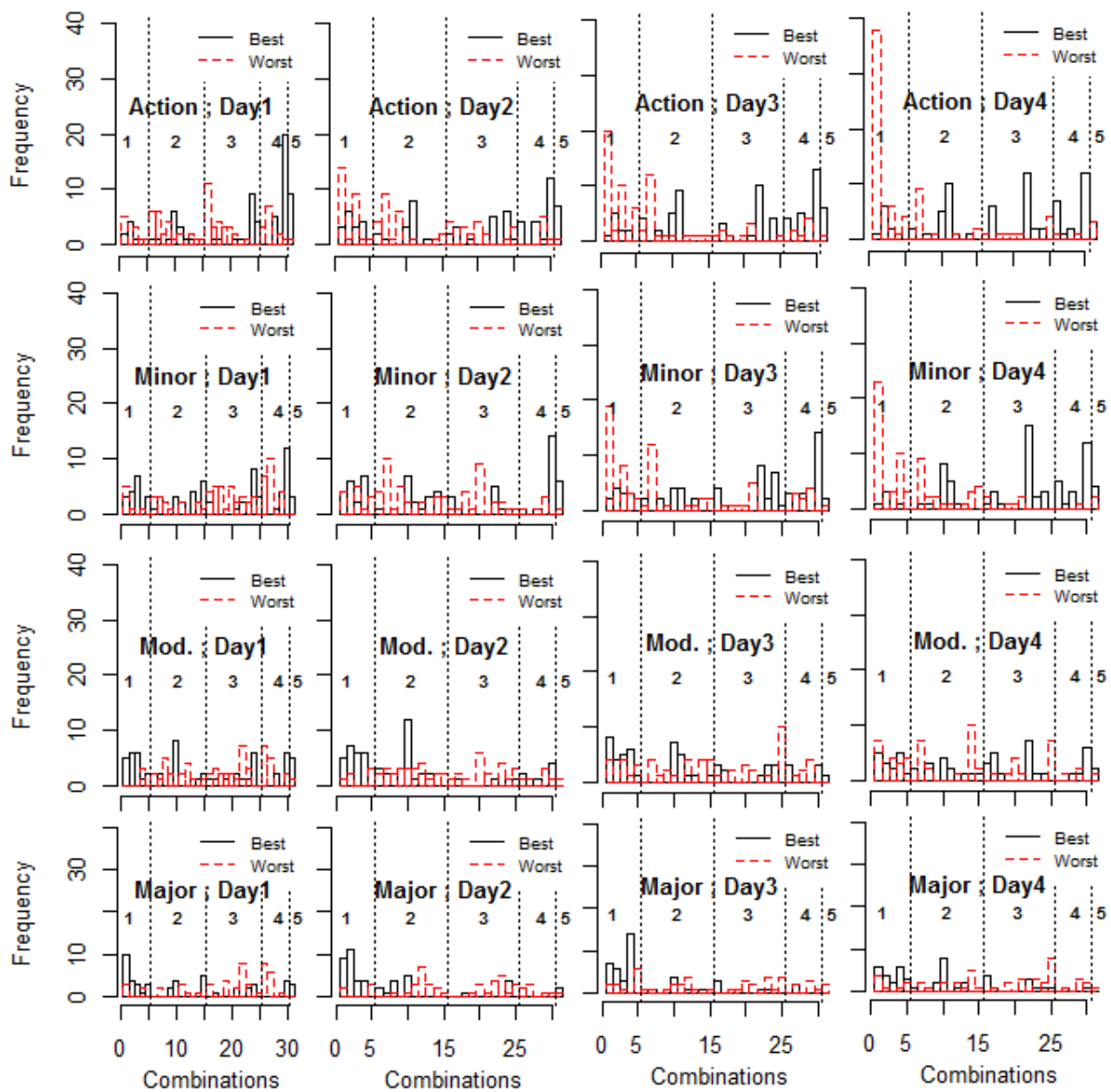
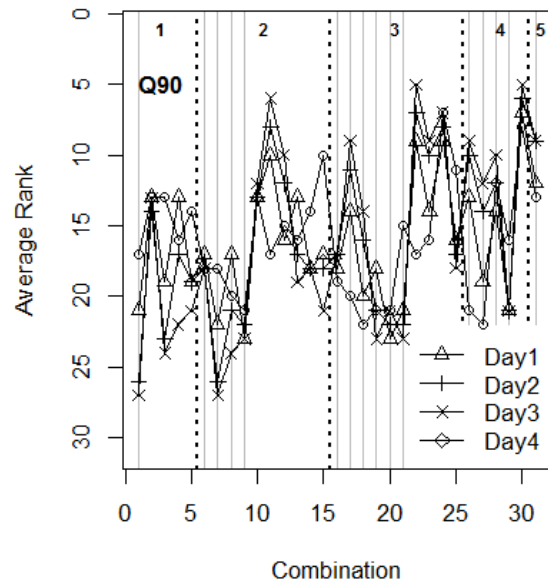
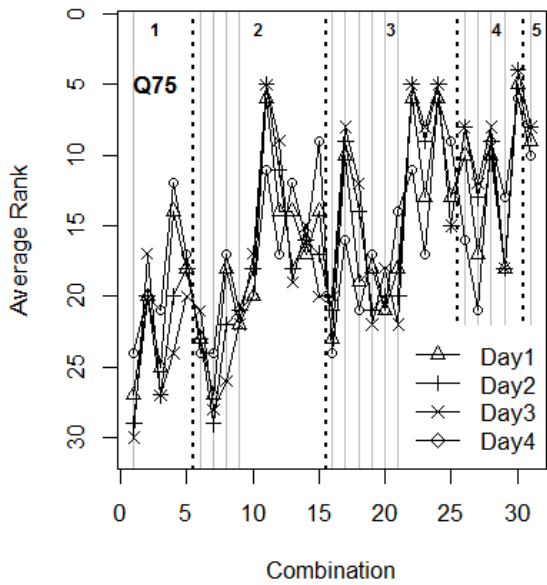
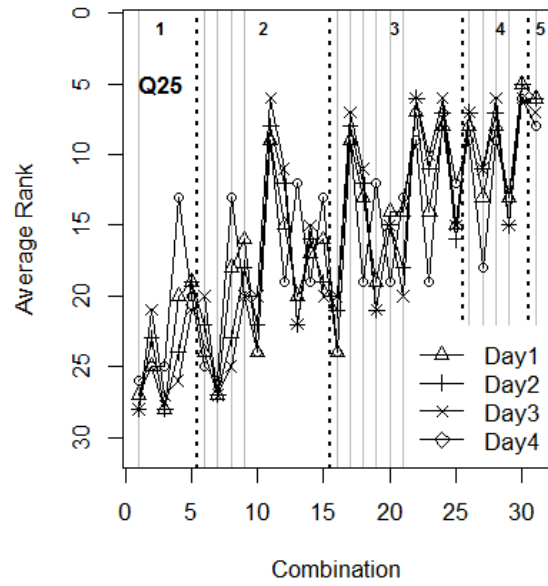
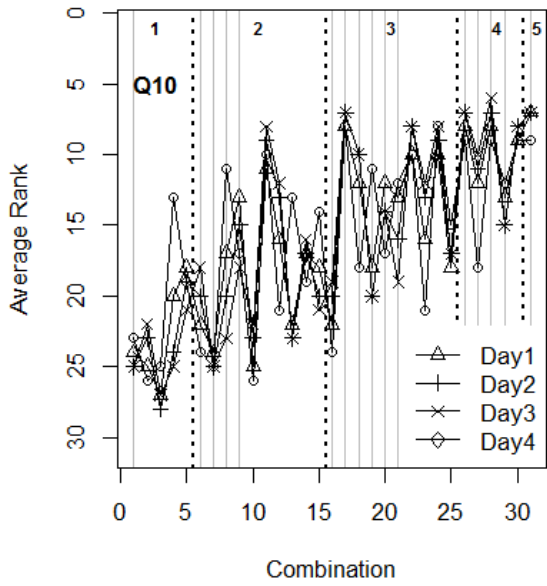


Figure 8: Histograms of **variable-combination joint predictors** returning the best and worst Brier Skill Scores across 82 river gages. Each row of histograms refers to a flood stage, and each column to a lead time. The dotted vertical lines in the histograms distinguish **variable-combination joint predictors** with different numbers of **independent variables**.



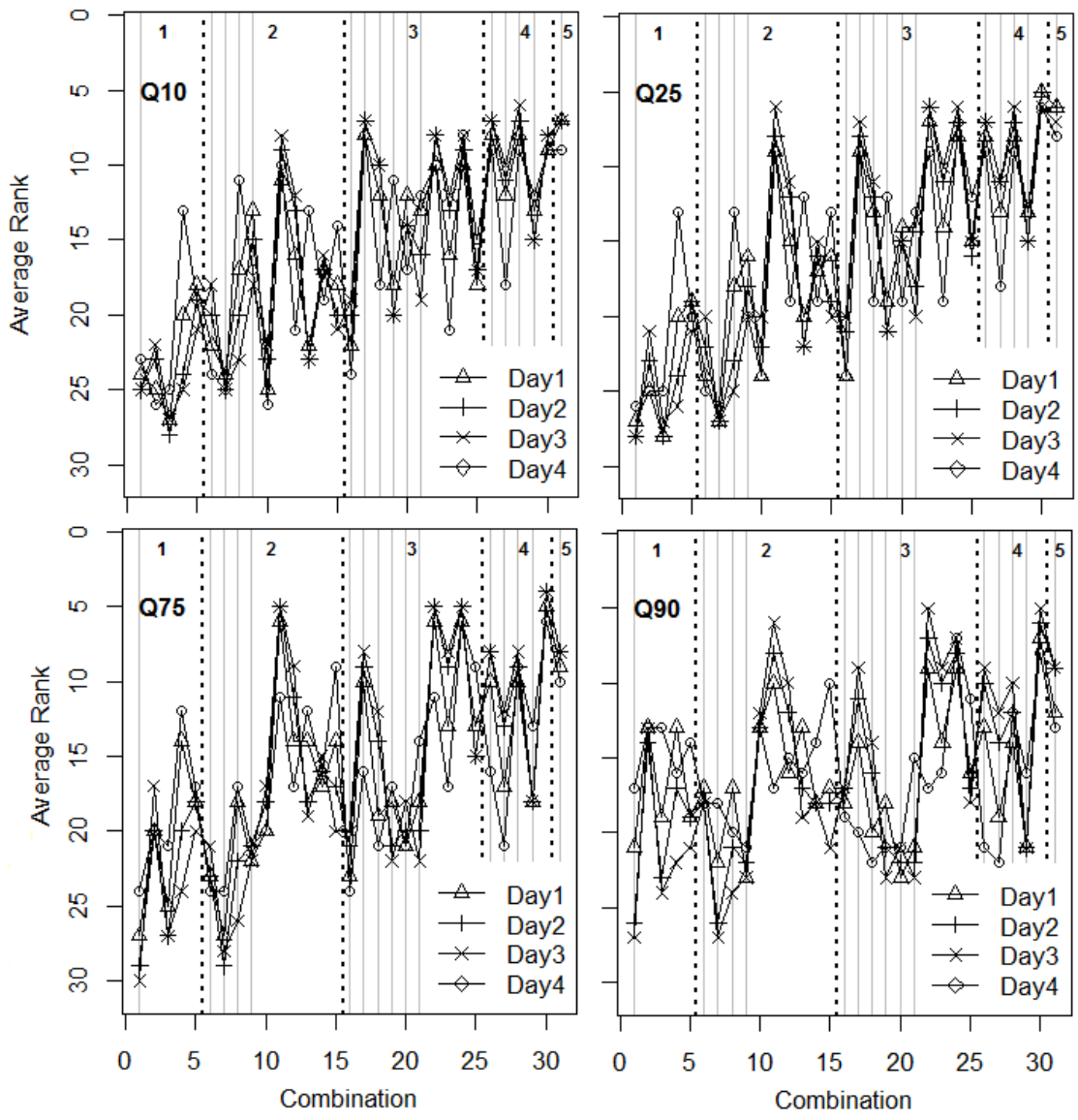
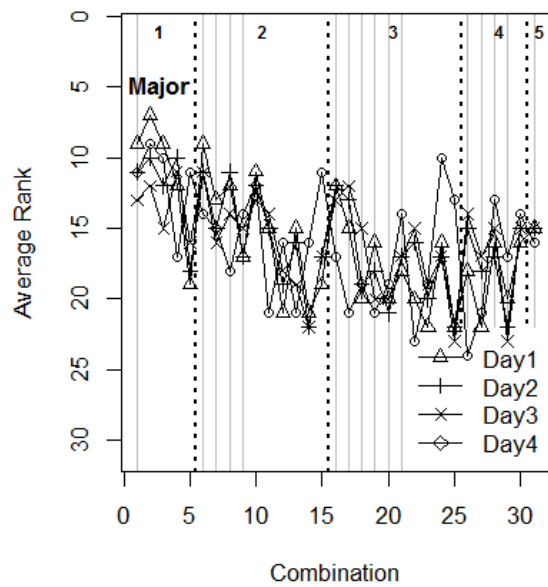
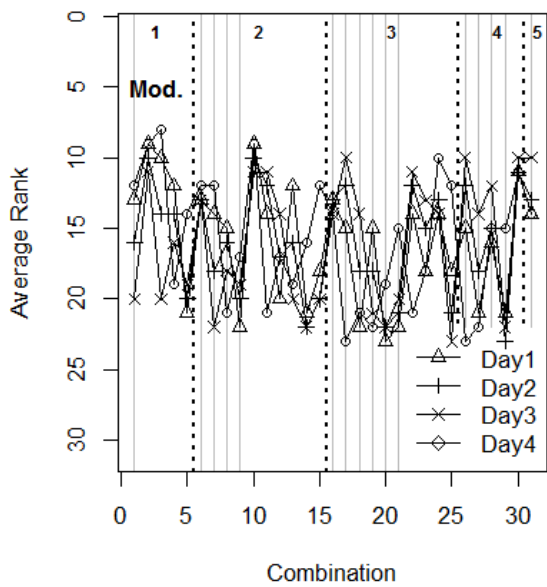
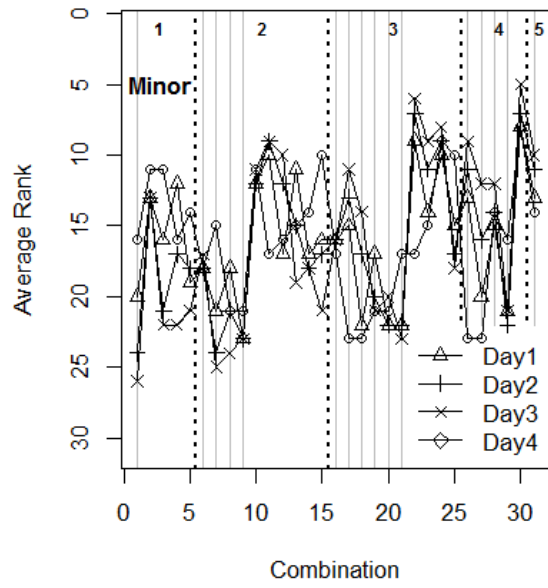
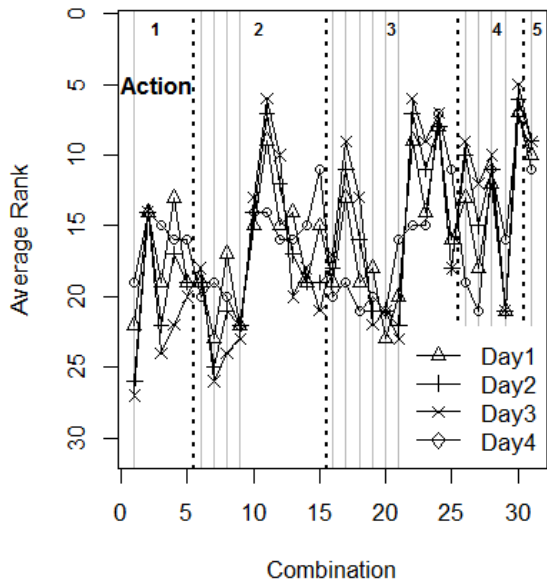


Figure 9: Average rank for each **variable combination joint predictor** for one to four days of lead time and four percentiles of observed water levels. Vertical gray lines indicate **variable combination joint predictors** including the forecast.



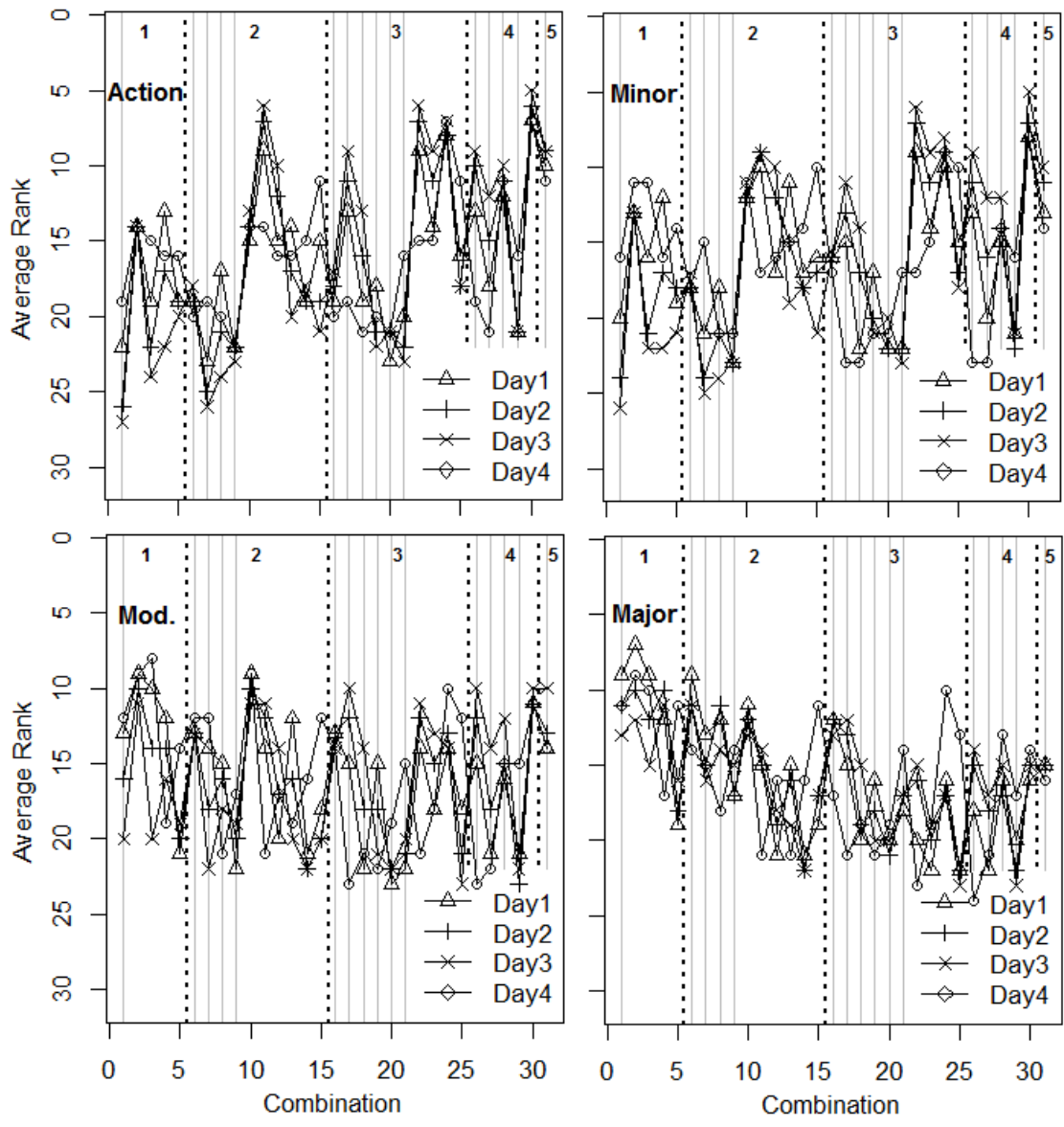
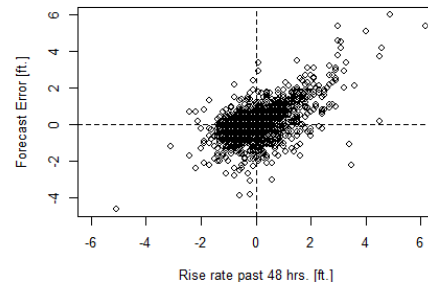
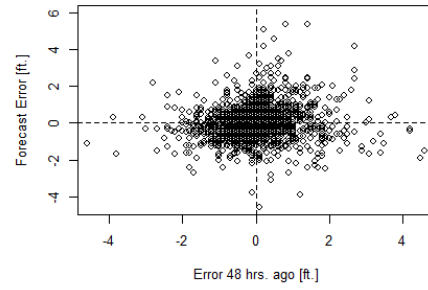
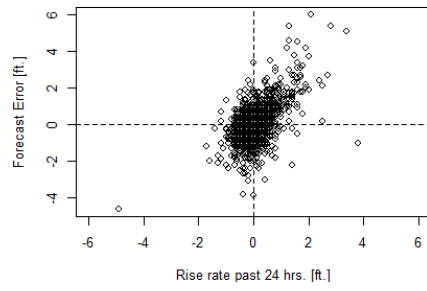
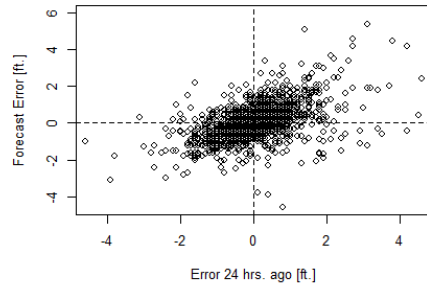
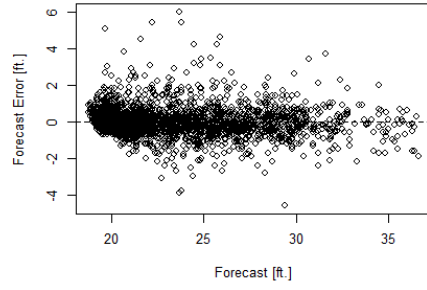


Figure 10: Average rank for each **variable-combination**joint predictor for one to four days of lead time and four flood stages. Vertical gray lines indicate **variable-combination**joint predictors including the forecast.



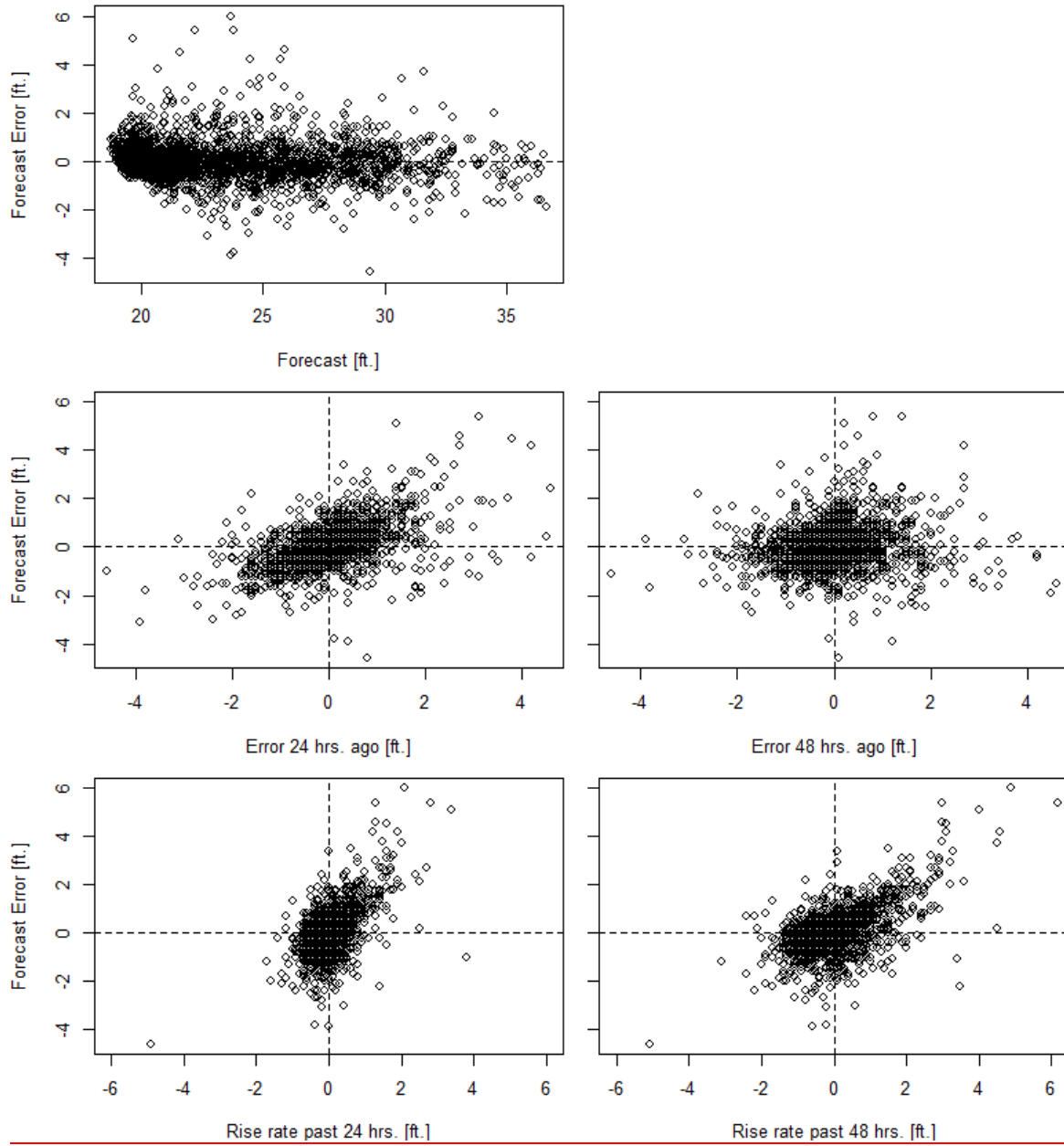
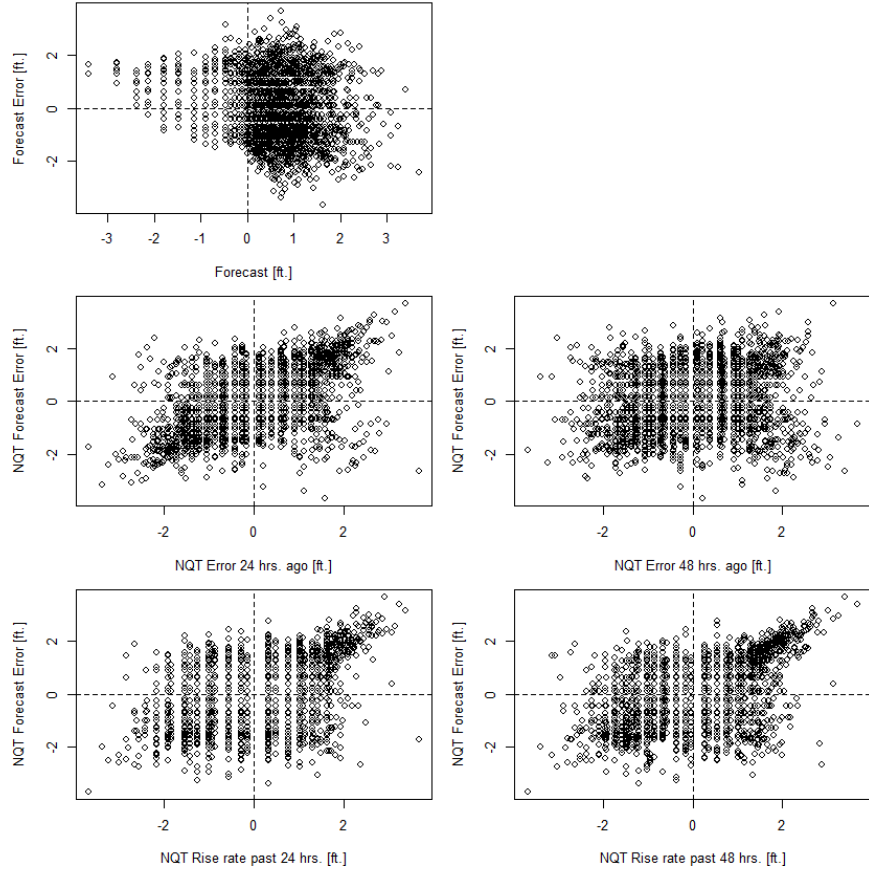


Figure 11: Independent variables plotted against the forecast error for Hardin IL with 3 days of lead time. First row: Forecast; second row: past forecast errors; third row: ~~rise rates~~ rates of rise.



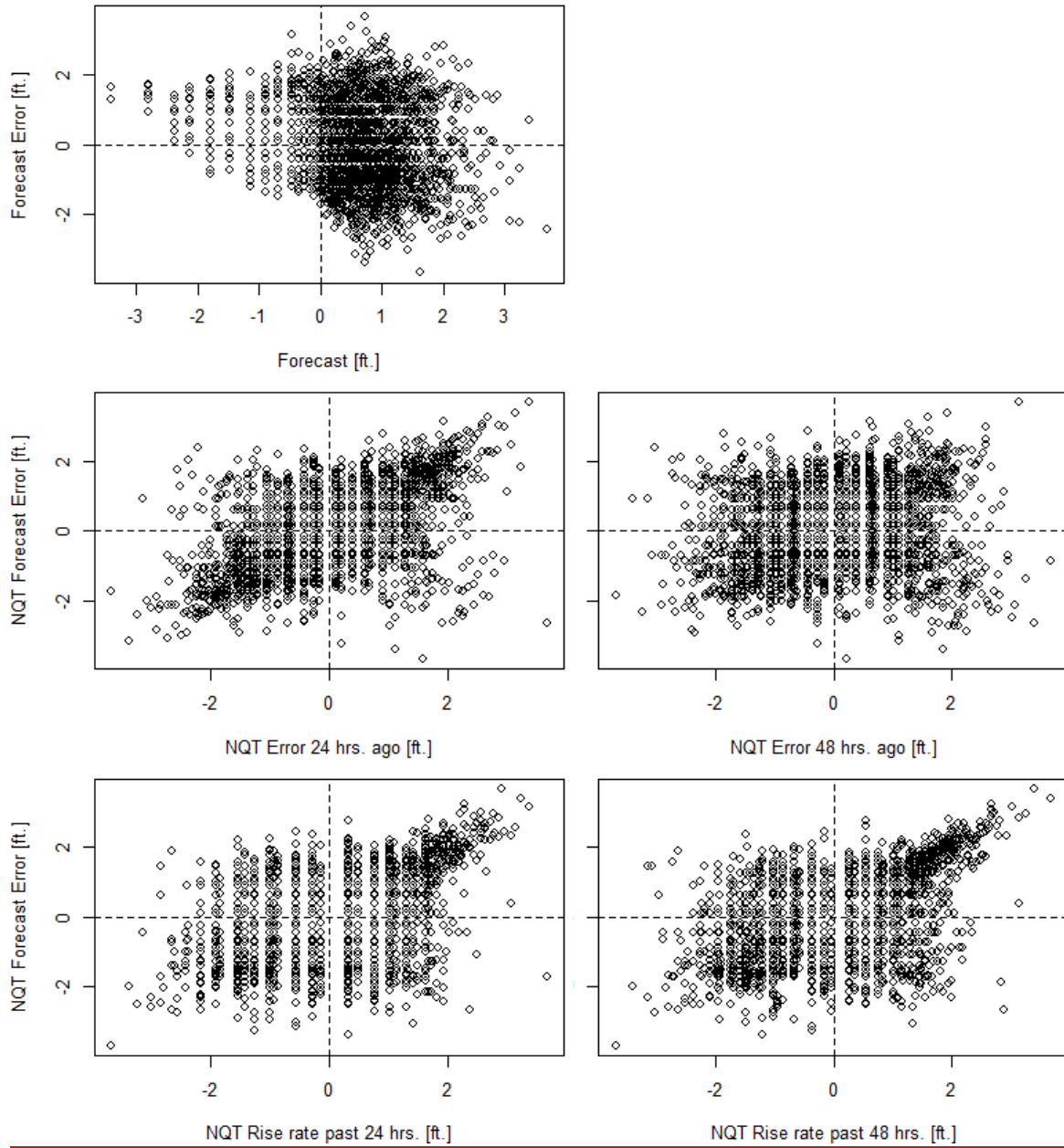


Figure 12: Independent variables after transforming into the Gaussian domain plotted against the forecast error for Hardin IL with 3 days of lead time. First row: Forecast; second row: past forecast errors; third row: rise rates rates of rise.

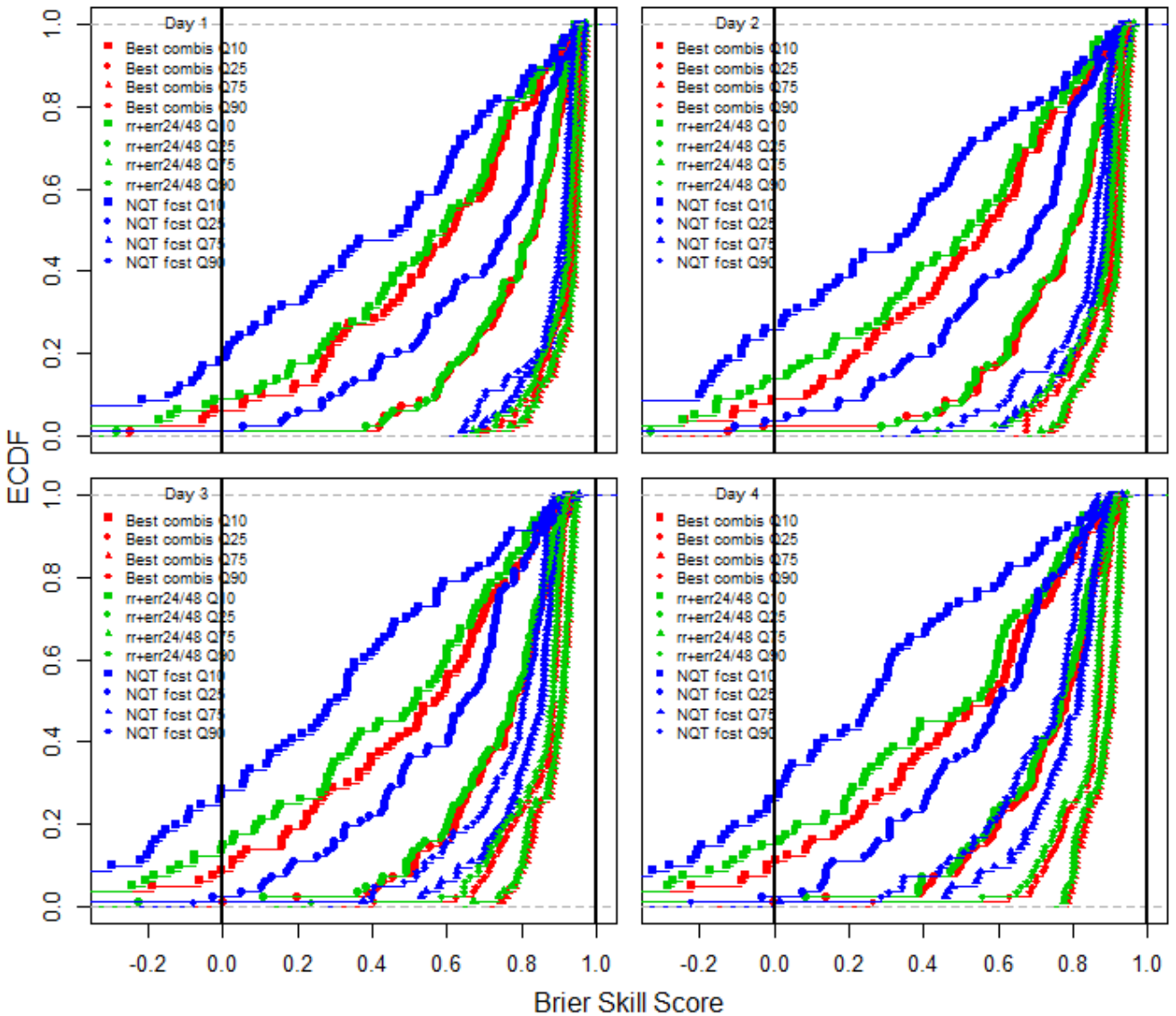


Figure 16: Empirical cumulative density functions of three QR configurations predicting exceedance probabilities of the 10th, 25th, 75th, and 90th percentile: the configuration using the transformed forecast as the only independent variable [NQT fcst]; the best performing combination for each river gage (upper performance limit) [Best combis]; rates of rise in the past 24 and 48 hours and the forecast errors 24 and 48 hours ago as independent variable (one-size-fits-all solution) [rr+err24/48].

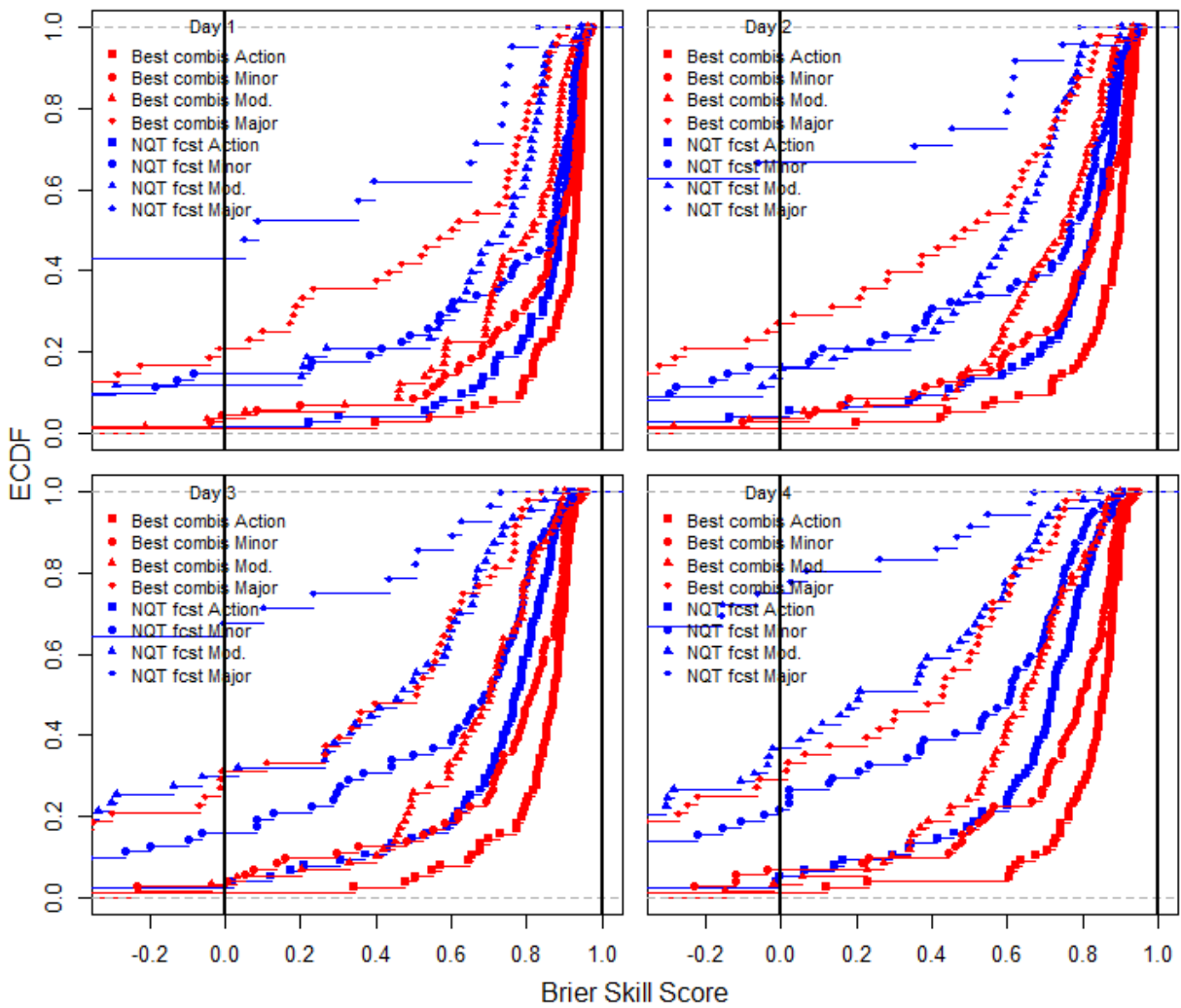
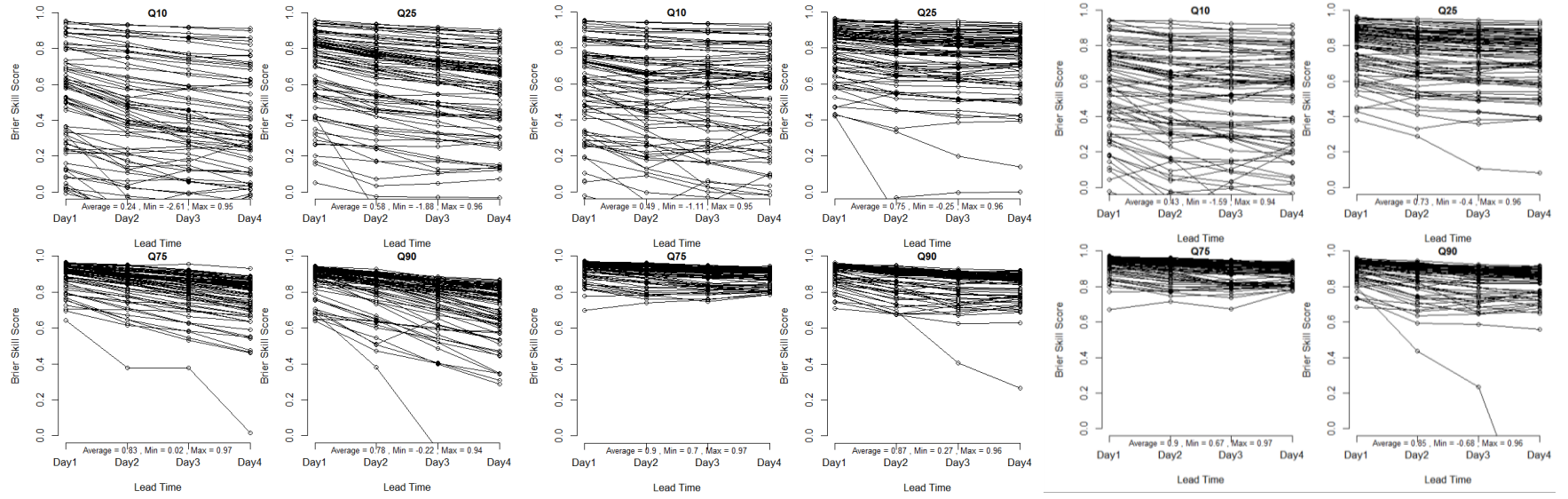


Figure 19: Empirical cumulative density functions of three QR configurations predicting exceedance probabilities of the Action, Minor, Moderate, and Major Flood Stage: the configuration using the transformed forecast as the only independent variable [NQT fcst]; the best performing combination for each river gage (upper performance limit) [Best combis]



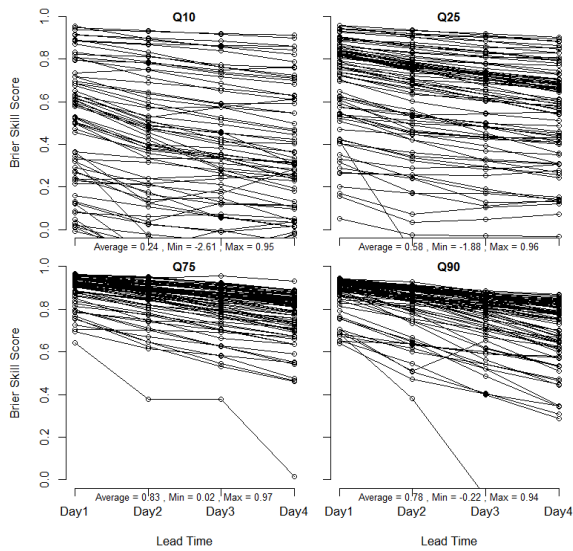


Figure 13: Brier Skill Scores of the **original forecast-only QR model configuration** (i.e., using the transformed forecast as the only independent variable) for four lead times and percentiles of observed water levels.

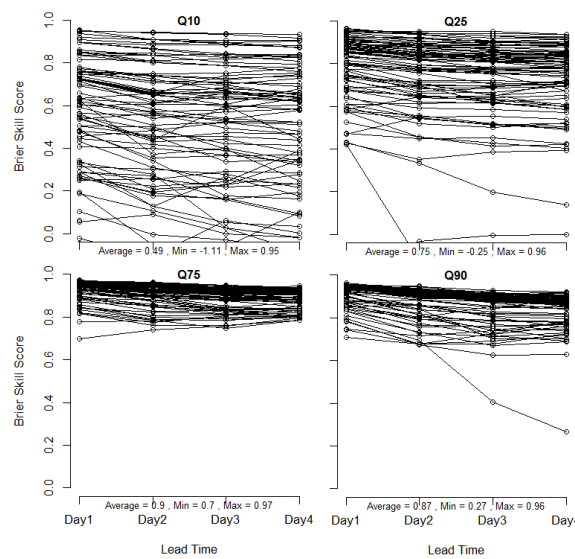


Figure 14: Brier Skill Scores for four lead times and percentiles of observed water levels using the best **variable combination joint predictor** for each river gage as independent variables in the QR **model configuration**.

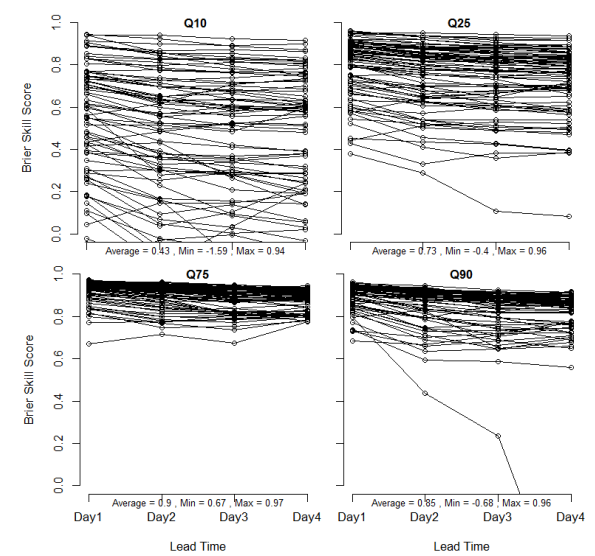
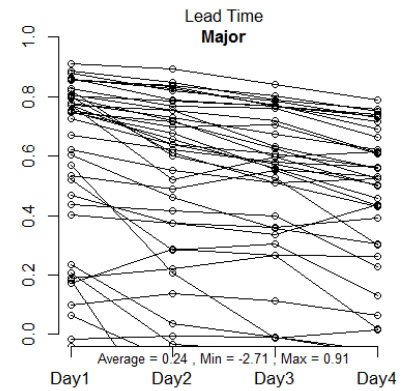
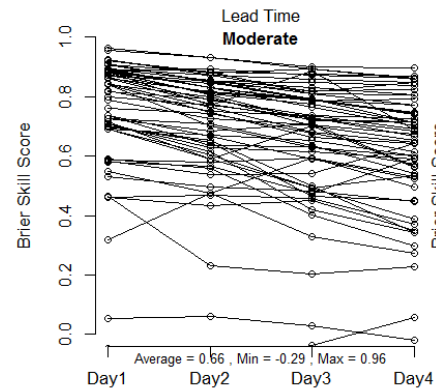
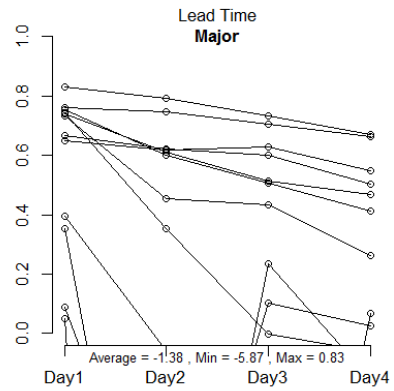
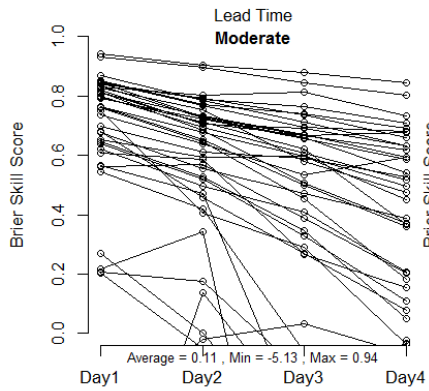
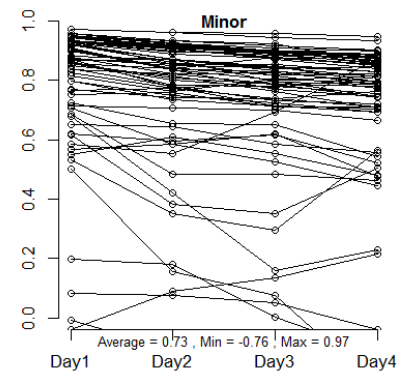
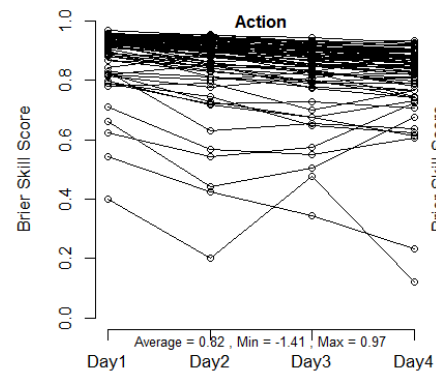
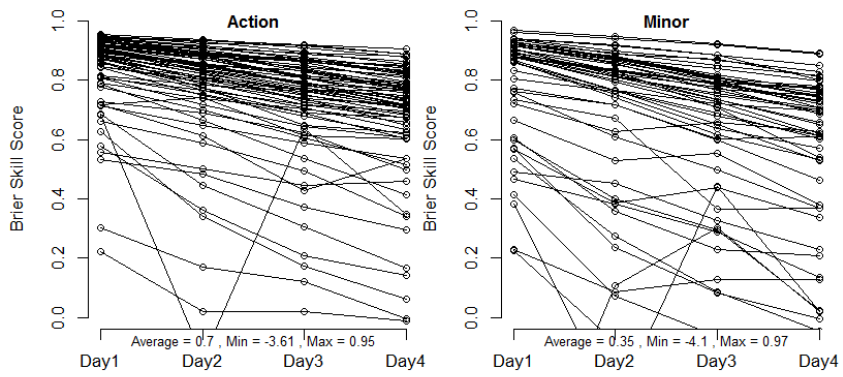


Figure 15: Brier Skill Scores for four lead times and percentiles of observed water levels using a one-size-fits-all approach (i.e., rr24, rr48, err24, err48) for the independent variables in the QR **model configuration**.



Lead Time

Lead Time

Lead Time

Lead Time

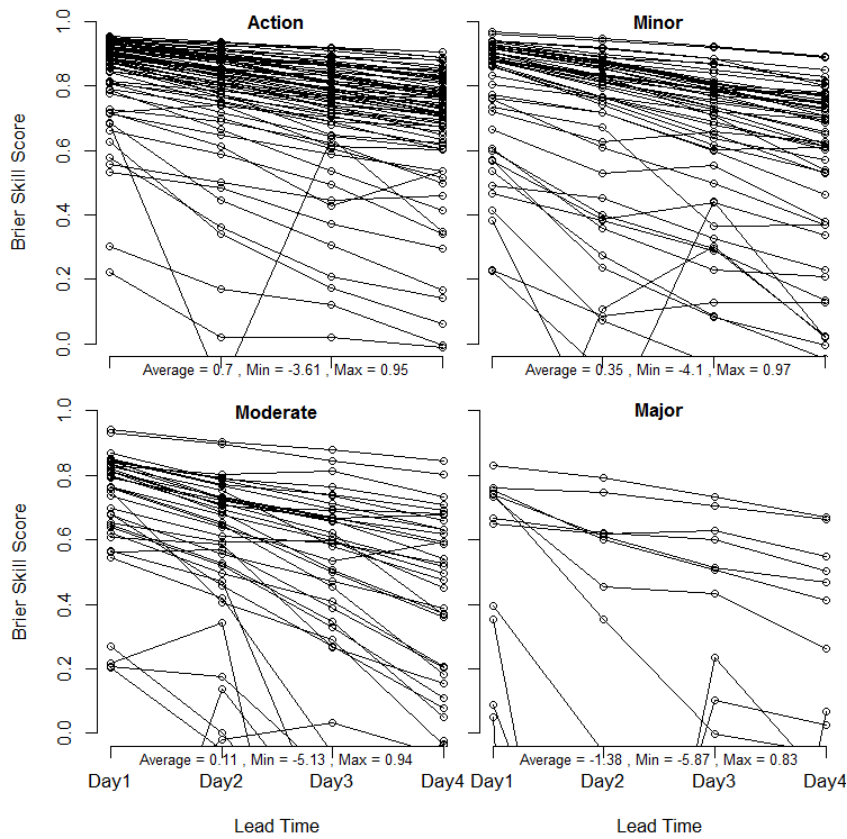


Figure 176: Brier Skill Scores of the forecast-only original-QR model configuration (i.e., using the transformed forecast as the only independent variable) for four lead times and flood stages.

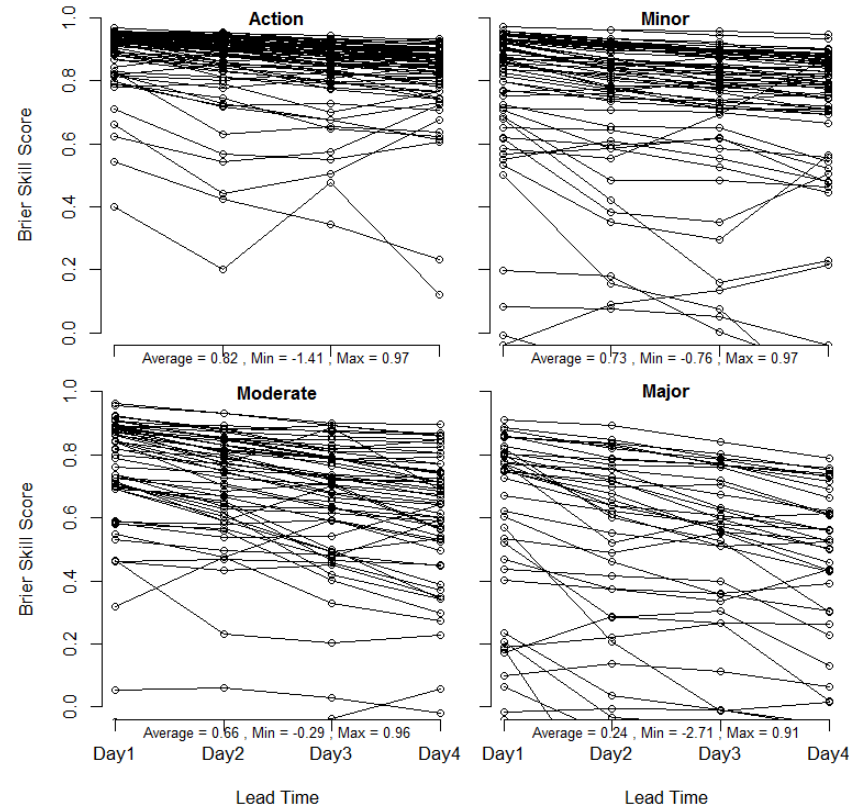


Figure 187: Brier Skill Scores for four lead times and flood stages of observed water levels using the best variable combination joint predictor for each river gage as independent variables in the QR model configuration.

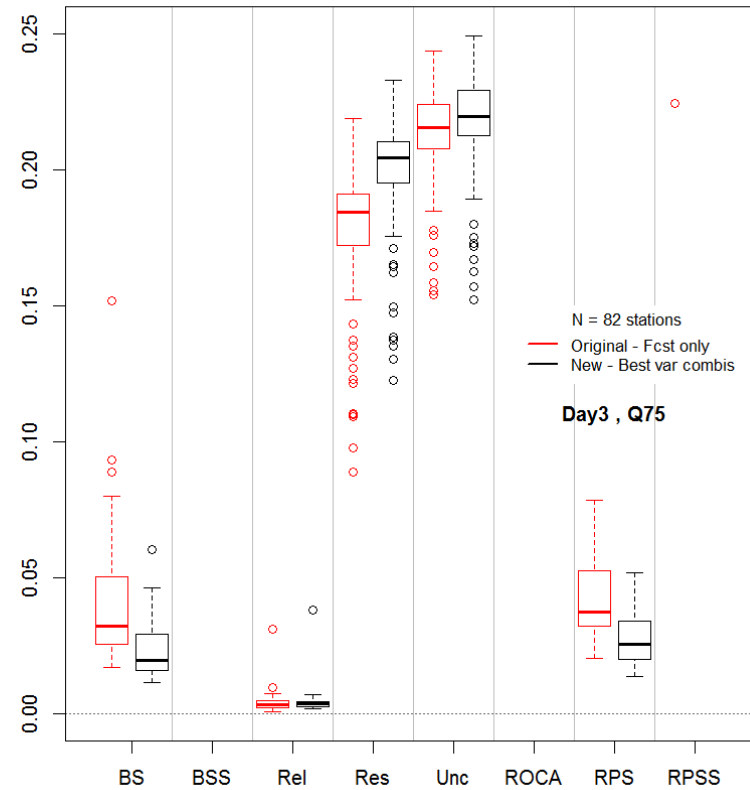
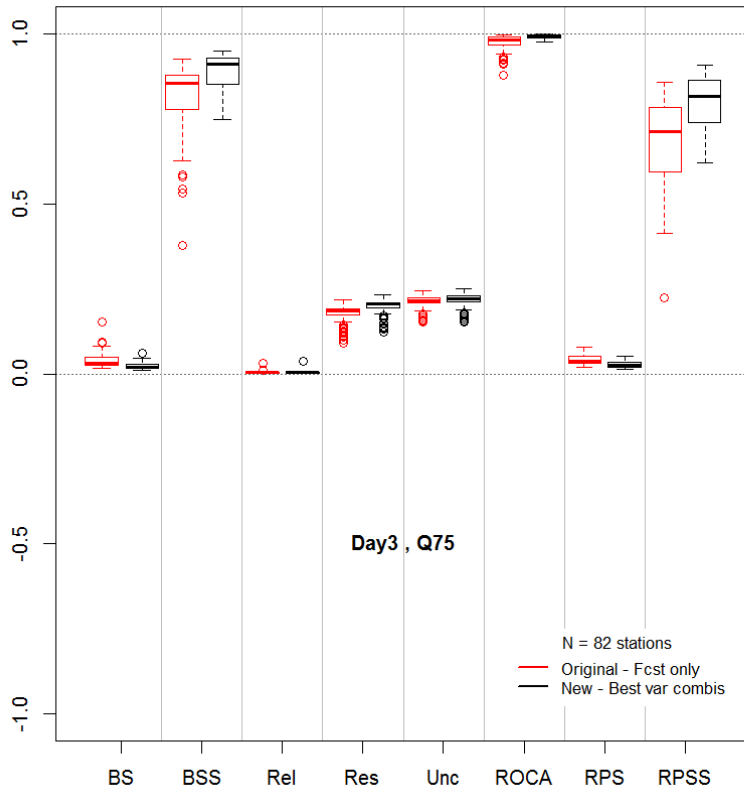
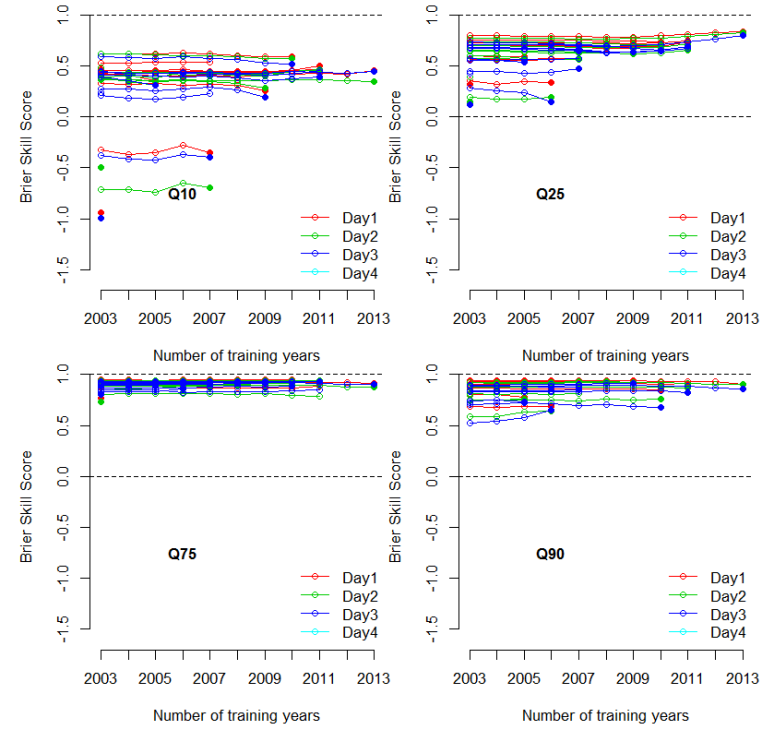
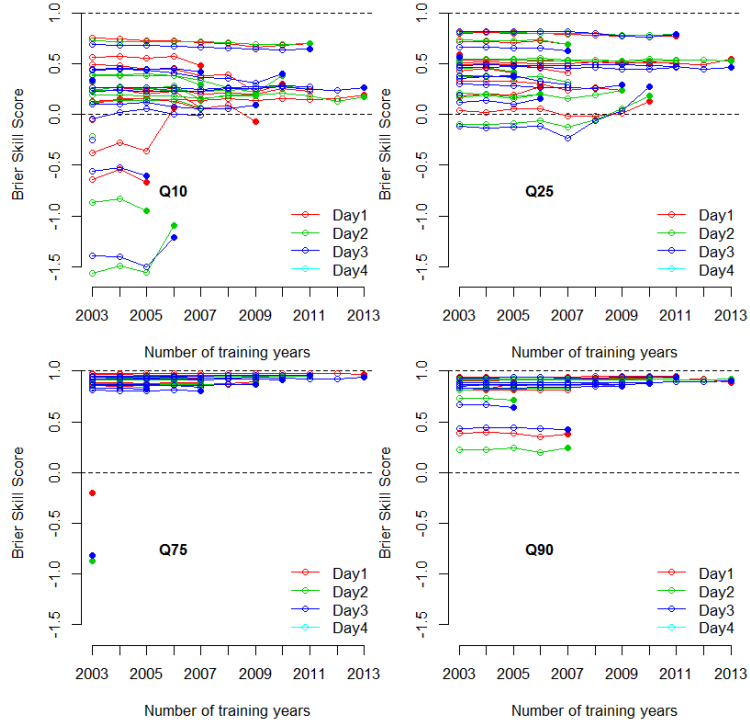


Figure 2018: Comparison of the forecast-only original-QR model configuration (i.e., only transformed forecast as independent variables) and the one-size-fits-all approach (i.e., rise-rates and forecast errors as independent variables) using various measures of forecast quality: Brier Score (BS), Brier Skill Score (BSS), Reliability (Rel), Resolution (Res), Uncertainty (Unc), Area under the ROC curve (ROCA), ranked probability score (RPS), ranked probability skill score (RPSS). Lead time: 3 days; 75th percentile of observation levels as threshold. The left figure zooms in on the right figure to make changes in reliability and resolution better visible.



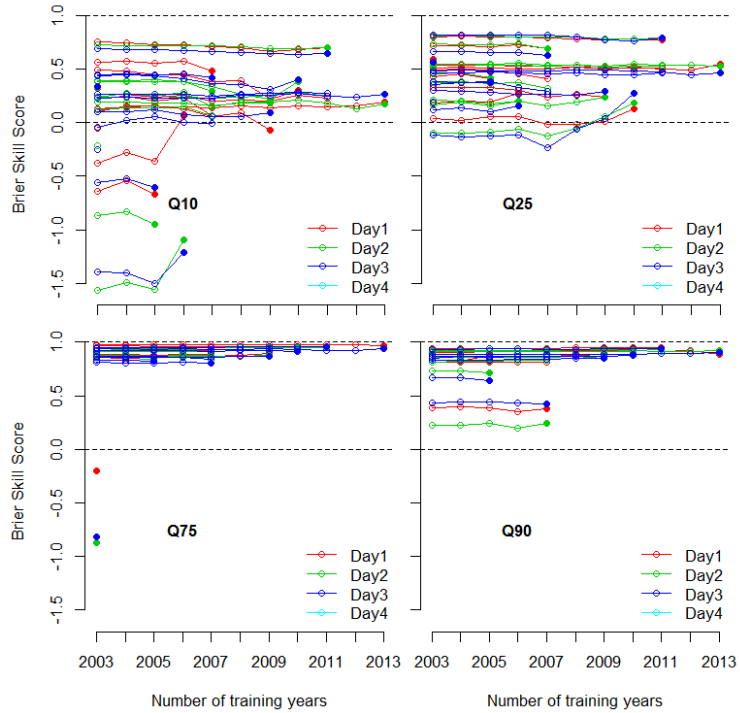


Figure 219: Brier Skill Score for various forecast years and various sizes of training dataset across different lead times (colors) and event thresholds (plots) for Hardin, IL (HARI2). The filled-in end point of each line indicates the BSS for the forecast year on the x-axis with one year in the training dataset. Each point further to the left stands for one additional training year for that same forecast year.

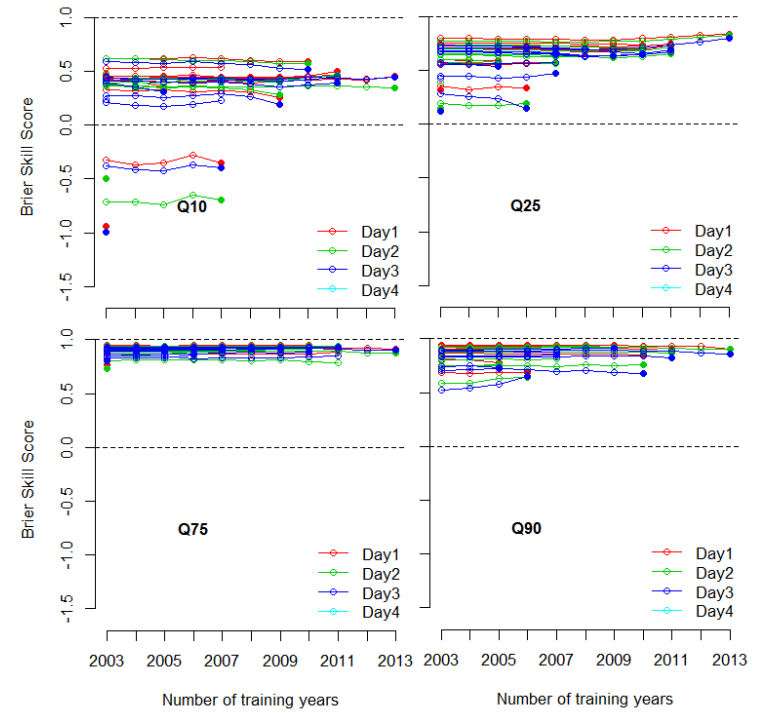


Figure 220: Brier Skill Score for various forecast years and various sizes of training dataset across different lead times (colors) and event thresholds (plots) for Henry, IL (HNYI2). The filled-in end point of each line indicates the BSS for the forecast year on the x-axis with one year in the training dataset. Each point further to the left stands for one additional training year for that same forecast year.

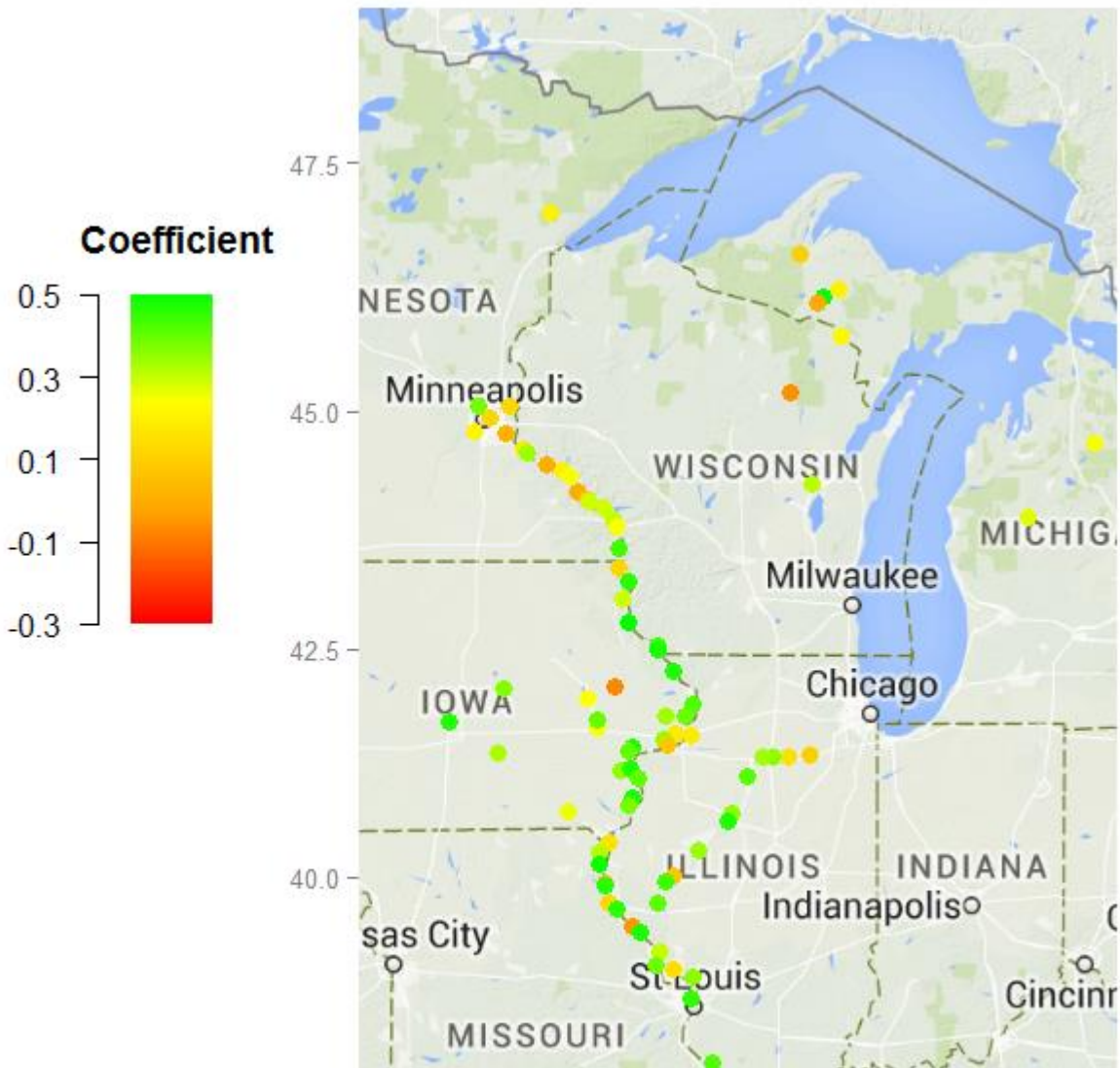
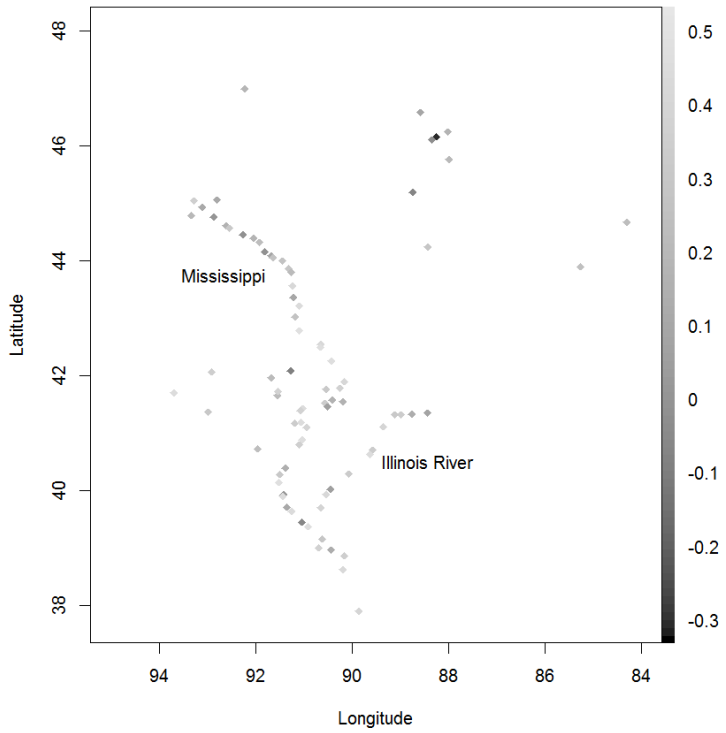


Figure 231: Geographical position of rivers. Colors indicate the regression coefficient of each station with the Brier Skill Score as dependent variable.

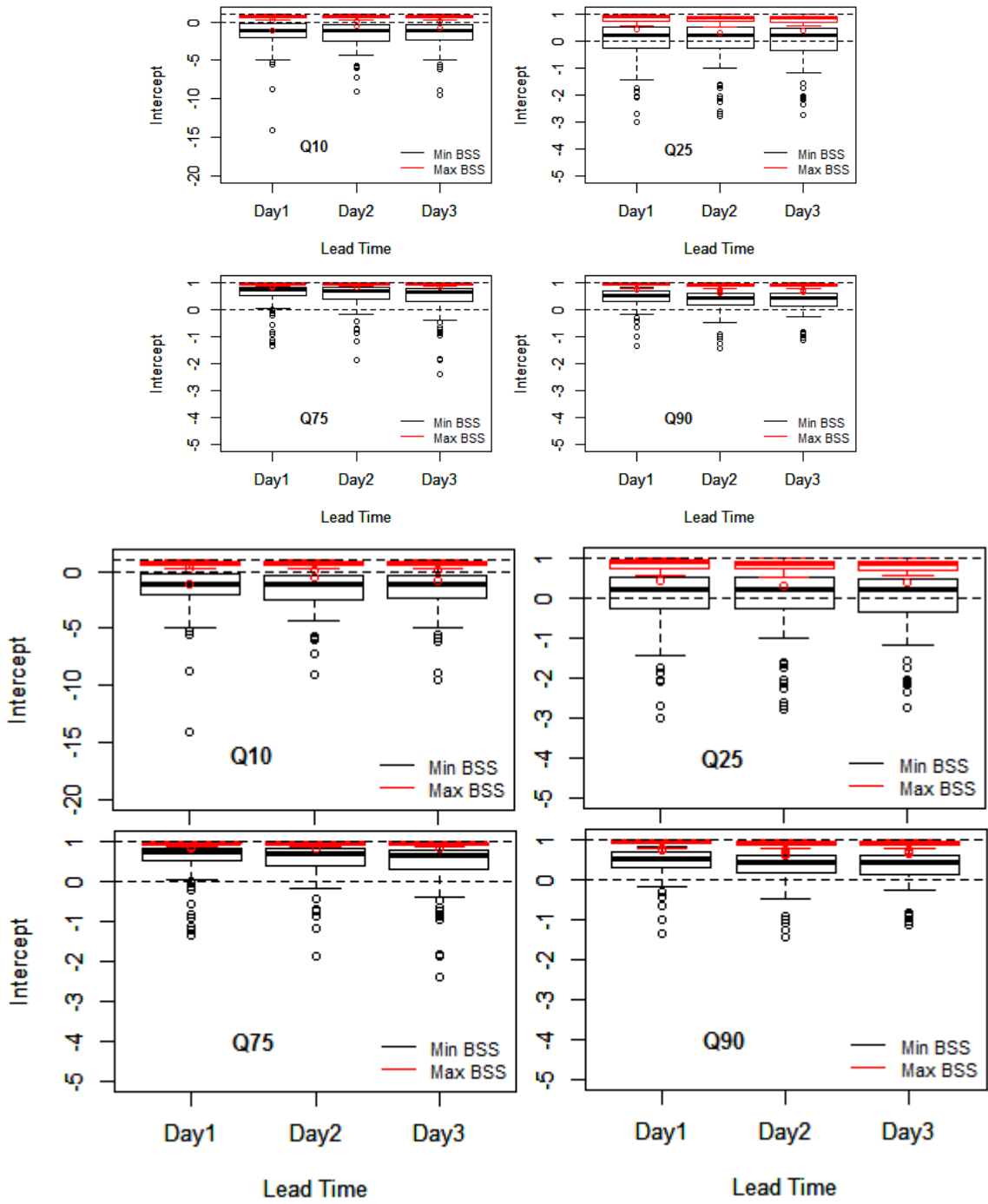


Figure 242: Minimum (black) and maximum (red) Brier Skill Scores for various lead times and event thresholds across locations, size of training dataset and forecast years.