

First of all, please let me offer my sincere apologies for being over two months late in submitting this review. In hindsight, I should have maybe declined the review request, but as I found the topic of this paper interesting and relevant to my own work also, I didn't want to pass up this opportunity. At the expense of timeliness however.

The Hoss and Fischbeck paper is an interesting and very worthwhile addition to the existing literature on predictive hydrological uncertainty insofar as it explores the optimal selection of predictors to configure a Quantile Regression based statistical post-processor for estimating predictive uncertainty. However, the paper requires some work prior to its publication. I don't think the the computations underlying the manuscript need to be modified in great extent, but the description and the analysis thereof can be improved substantially. I hope to give some suggestions on how that can be done.

#### General comments

- The manuscript could benefit from a more substantial “hydrological analysis” of the forecasts made. Post-processors can be used to find statistical relations between predictors and predictands. There needs to be correlation and causality. The paper could benefit from a more in-depth analysis of the latter: what does the ‘forecast error’ depend on? Here, the authors choose rate of rise and past forecast error: these appear to be more or less randomly chosen, and are subsequently applied to all forecasting locations considered. However, I think that an analysis of the hydrology of the basins considered, in conjunction with the forecasting models for those basins, could reveal important information on how those models are expected to perform. How are the models calibrated? What does this mean for extreme events? Is the relation between predictors and predictand stationary across ‘normal flow regimes’ and ‘extremes’? This likely varies with basin, and therefore one should consider varying post-processing configurations with basin also.
- There is one important assumption underlying the use of statistical post-processors: stationarity of the joint predictor, predictand distributions. The paper would benefit from a discussion thereof, particularly in relation to the results section, and the ‘robustness’ section contained therein.
- “First US application” is irrelevant to the science and also incorrect, as Wood et al (see reference in Weerts et al, 2011) applied QR previously. This comes back a couple of times in the paper. Also, QR was originally devised by Roger Koenker; not by Weerts et al (I wish!).

- Different users have different needs for uncertainty information; it is not universally true that users benefit most from probabilities of exceedence or non-exceedence. Likewise, not all users are interested in extreme events per sé. This comes back a couple of times in the paper.
- I would recommend to streamline use of terms:
  - ‘predictor’ or ‘independent variable’
  - ‘predictand’ or ‘dependent variable’
  - preferably omit use of ‘variable’ in context of statistical post-processors, as its interpretation can be ambiguous
  - ‘configuration’ rather than ‘model’ (to avoid confusion with underlying hydrological models)
- Please consider removing the footnotes. If the text contained therein is important, include it in the main body of the paper. If not, you may want to consider omitting it altogether.
- Practicalities of data access are not too relevant to the science and I would suggest omitting descriptions of why certain data sources could (not) be accessed and how much effort that would require. Instead, you could turn the argument around and say: “this and this is available and we’re trying to assess if there is any signal that can contribute to better probabilistic forecasts.”

### Specific comments

#### Introduction:

- Some elements can be safely omitted from the introduction:
  - Discussion on QPF forecasts
  - Discussion of RFC produced “outlooks”
- Verifying by means of BSS only is somewhat limited I think, but it does fit with the authors’ wish to verify exceedence probabilities only. Why not, however, use a range of verification metrics? See, for example, some of the recent Brown and Seo papers as well as some of my own work (where the verification approach was inspired on the Brown/Seo papers).
- “Rate of rise” is more commonly used than “rise rate” I think.

#### 2.2 Brier Skill Score:

- The ‘method’ section would benefit from a subsection on verification metrics. That section would then include the current sub-section on BSS, but also some discussion of other metrics now included in the ‘results’ section.
- A decomposition of Brier’s probability score is included; what’s

missing, is a note on how these decompositions are computed in terms of *skill*. See one of the Brown and Seo papers for how that's done. Also, no quantified decompositions are shown in the results/analysis section?

### 2.3 Proposed addition

- The current title "Proposed addition: more than one independent variable" suggests that it is the *number* of predictors that's important. This is not necessarily so - it's content, not just quantity that's relevant. Please consider retitling this section.
- This section could really benefit from some 'hydrological intelligence': what are the factors determining level of accuracy of model predictions? Are these already included in the model itself somehow? If so, how? If not, why not? To me, it is still an open question: what to include in a model, and what to include in a post-processor? Where is the boundary between statistical modeling and modeling of physical processes? This point is one that the authors should also re-visit in the discussion/conclusions section.
- Table 1: "forecast error 24 hours ago". I understand this to be 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 - correct? Maybe good to state this.

### 2.5 Data:

- First sentence may be omitted, or moved to the introduction.
- The manuscript would benefit from a custom made map showing the forecasting locations and basin delineations.

### 3.2.2 Best performing combinations

- The forecasts for extreme conditions perform worse when using multiple predictors. Why - overfitting? Some in-depth analysis would be good.

### 3.3 Robustness

- I think the 'robustness' analysis could, and should, be simplified by using a leave-one-year-out analysis. Length of training set is less relevant than stationarity of joint predictand, predictor distributions. Why not simply use all of the available data most efficiently and then discuss any drops in forecast quality? Also, the current analysis results in a difference in sample size and this would require an analysis of the uncertainty in resulting BSS - which

is likely bigger for smaller samples. With a leave-one-year-out analysis, sample size would be equal and the authors would be more easily forgiven for not analysing uncertainty.

- Some hydrologic analysis could contribute to explaining why forecast quality is different between locations.

#### “Future work”

- Yes, more analysis on which predictors to use could work. Please refer to my earlier comments also on statistical modeling versus numerical modeling of physical processes, and on using knowledge of the hydrology of basins to determine meaningful predictors.

#### Figures:

- The multi-plot figures contain a lot of white space between plots. As some horizontal and vertical axes are identical across plots within the figure, I would suggest eliminating the in-between space altogether. In figures 10 and 11, this can be done for the vertical axes also. In R: `par(mar = c(.5,0,0,0))` and then `plot(..., xaxt="n")` for plots where you can omit horizontal axis.

Additional specific comments are included in attached, annotated PDF.























For low event thresholds, the BSSs are much worse than for high thresholds, and the BSSs slightly decrease with lead time (Table 4). The regression is slightly biased regarding the forecast quality for each forecast year. The earlier years are included less often in the dataset with on average less years' worth of data in their training dataset, because, for example, unlike for the year 2013, ten years of training data were not available for the year 2008. Nonetheless, the regression indicates that 2008 was particularly difficult to forecast and 2012 relatively easy, i.e. they are associated with relatively low and high coefficients respectively (Table 4). The performance of the forecast additionally depends on the river gage. The coefficients of the river gages, included as factors in the regression, have been excluded from Table 4 for the sake of brevity. Instead, Fig. 21 maps the geographic position of the river gages with the color code indicating each gage's regression coefficient. The coefficient is lower, and therefore the Brier Skill Scores are lower, for gages far upstream a river and those close to confluences. The latter is particularly visible where the Illinois River and the Mississippi River join. At least for the gages at confluences, the QR model could probably be improved by including the rise rates at the river gages on the other joining river into the regression.

#### 4 Conclusions

In this study, quantile regression (QR) has been applied to estimate the probability of the river water level exceeding various event thresholds (i.e., 10th, 25th, 75th, 90th percentiles of observed water levels as well as the four flood stages of each river gage). This is the first study applying this method to the US American context. Additionally, it further develops the method by including more independent variables and testing the method's robustness across locations, lead times, event thresholds, forecast years and sizes of training dataset.

Most importantly, it was found that including rise rates in the past 24 and 48 h and the forecast errors of 24 and 48 h ago as independent variables improves the performance

11301

of the QR model, as measured by the Brier Skill Score. Since the reliability was already high with the original QR method as proposed by Weerts et al. (2011), the new configuration mainly increases the resolution.

For extremely high water levels, the combinations of independent variables that perform best vary across stations. On those days, combinations of fewer variables perform better than those that include more. In contrast to these extremely high event thresholds, larger sets of variables work better than smaller ones for non-extreme and low event thresholds. Additionally, a one-size-fits-all approach (i.e. the rise rates and forecast errors as independent variables) performs satisfactorily for those cases.

The new independent variables – rise rates and forecast errors – do not combine well with forecast itself. The latter was the only variable included in the original QR configuration as studied by Weerts et al. (2011) and López López et al. (2014). To account for heteroscedasticity, the forecast was transformed into the Gaussian domain. However, the rise rates and the forecast errors do not lend themselves for linear quantile regression after such a transformation. Therefore, it is difficult to combine these two variables. A possible solution could be to build regression models for subsets of the transformed data. However, such an approach drastically decreases the amount of data available for each model.

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

11302

The proposed approach performs less well for longer lead times, for gages far upstream a river or close to confluences, for low event thresholds and extremely high ones. The model might be performing less well for low event thresholds, because the variance in the dependent variable – the forecast error – is smaller. After all, river forecasts have much smaller errors for lower water levels. In turn, for extremely high water levels, the scarcity of data decreases the model performance.

### Future work

The methods can be further developed in several ways to achieve higher Brier Skill Scores and more robustness. First, more independent variables can be added. Trials with a different method, classification trees, showed that the observed precipitation, the precipitation forecast (i.e., POP – probability of precipitation) and the upstream water levels significantly improve model performance. Presumably, this is the case, because the QPF-forecast includes the precipitation forecast only for the next 12 h. However, currently, the precipitation data and forecasts can only be requested in chunks of a month, three chunks per day, from the NCDC's HDSS Access System.<sup>23</sup> For a period of 12 years, requesting such data for several weather stations<sup>24</sup> is obviously time-consuming. Upstream water levels can easily be included after manually determining the upstream gage(s) for each of the 82 NCRFC gages. To improve model performance at gages close to river confluences, the upstream water level of the gages on the joining river should be included as well.

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

<sup>23</sup>URL: <http://cdo.ncdc.noaa.gov/pls/plhas/HAS.FileAppSelectfidata.setname=9957ANX>, last access: July 2014.

<sup>24</sup>The geographical units of the weather forecasts bulletins do not correspond with those of the river forecast bulletins.

11303

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

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

*Acknowledgements.* To ensure anonymity, this section will be added after the review process.

### References

- Alexander, M., Harding, M., and Lamarche, C.: Quantile regression for time-series-cross-section-data, *International Journal of Statistics and Management System*, 4, 47–72, 2011.
- And Brier Score, Wikipedia Free Encycl., available at: [http://en.wikipedia.org/wiki/Brier\\_score](http://en.wikipedia.org/wiki/Brier_score) (last access: 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, 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, 1–3, 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. Dynam.*, 39, 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, *B. Am. Meteorol. Soc.*, 95, 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. Roy. Meteor. Soc.*, 138, 1954–1960, doi:10.1002/qj.1924, 2012.

11304



- Hsu, W. and Murphy, A. H.: The attributes diagram a geometrical framework for assessing the quality of probability forecasts, *Int. J. Forecasting*, 2, 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. Meth. Prog. Bio.*, 67, 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, 17–31, doi:10.1007/BF02428423, 1997.
- Koenker, R.: *Quantile Regression*, Cambridge University Press, New York, 2005.
- Koenker, R.: `quantreg`: Quantile Regression, R Package Version 505, available at: <http://CRAN.R-project.org/package=quantreg> (last access: 27 August 2014), 2013.
- Koenker, R. and Bassett, G.: Regression quantiles, *Econometrica*, 46, 33–50, 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, 1296–1310, doi:10.1080/01621459.1999.10473882, 1999.
- Leahy, C. P., Sri Srikanthan, G. Amirthanathan, Soori Sooriyakumaran, and Hydrology Unit: Objective Assessment and Communication of Uncertainty in Flood Warnings, in: 5 th Flood Management Conference Warrnamboll, 9–12 October 2007, 1–6, 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.*, 18, 3411–3428, doi:10.5194/hess-18-3411-2014, 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, 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, W01106, doi:10.1029/2003WR002540, 2004.
- Montanari, A. and Grossi, G.: Estimating the uncertainty of hydrological forecasts: a statistical approach, *Water Resour. Res.*, 44, W00B08, doi:10.1029/2008WR006897, 2008.

11305

- Morss, R. E.: Interactions among flood predictions, decisions, and outcomes: synthesis of three cases, *Natural Hazards Review*, 11, 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, 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, available at: [http://www.nap.edu/catalog.php?record\\_id=11699&utm\\_expid=4418042-5.krRTDpXJQISoXLpdo-1Ynw.0](http://www.nap.edu/catalog.php?record_id=11699&utm_expid=4418042-5.krRTDpXJQISoXLpdo-1Ynw.0) (last access: 18 September 2014), 2006.
- NCAR-Research Applications Laboratory, N.-R. A.: *verification: Weather Forecast Verification Utilities*, available at: <http://cran.r-project.org/web/packages/verification/index.html> (last access: 27 August 2014), 2014.
- Pielke, R. A.: Who decides? Forecasts and responsibilities in the 1997 Red River Flood, *Applied Behavioral Science Review*, 7, 83–101, 1999.
- R-Core Team: *R: A language and Environment For Statistical Computing*, available at: <http://www.R-project.org/> (last access: 27 August 2014), 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, 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, W00B11, 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, 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

11306

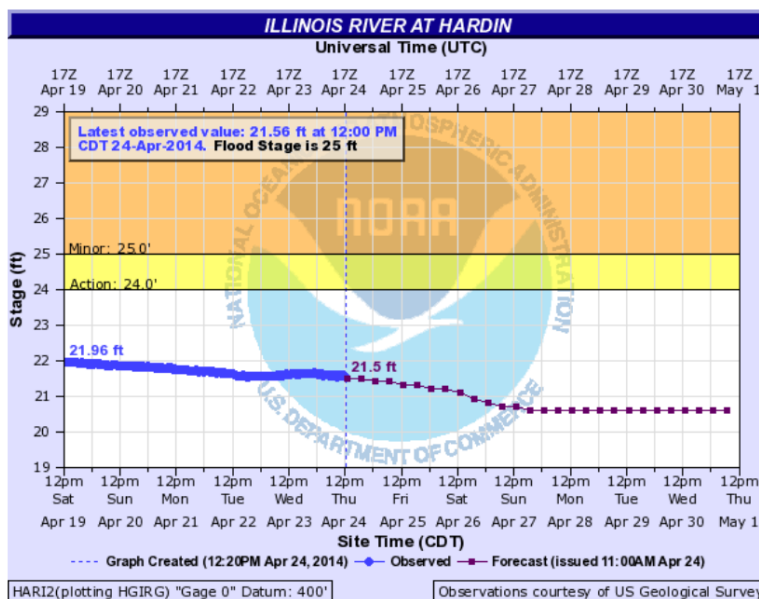




**Table 4.** Regression results.

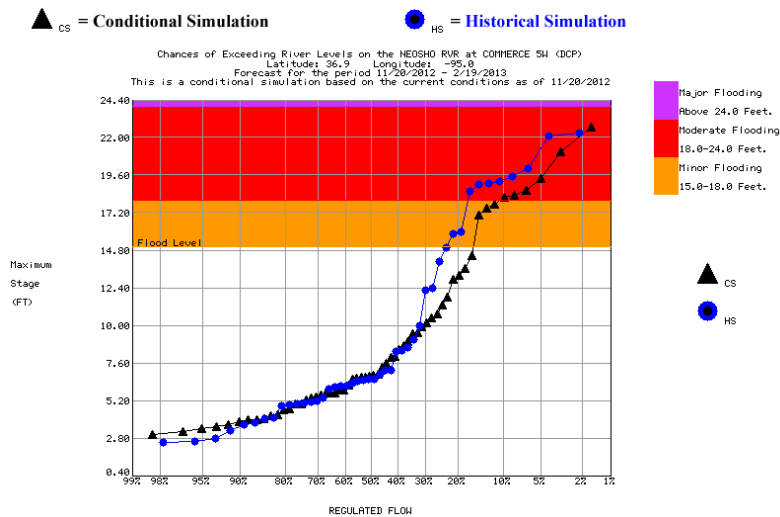
	Coef.	SD	
Intercept	-0.206	0.031	***
Event thresholds	0.265	0.003	***
Lead Times	-0.021	0.003	***
Forecast Years			
2004	-0.266	0.020	***
2005	-0.081	0.018	***
2006	-0.125	0.017	***
2007	-0.129	0.017	***
2008	-0.203	0.017	***
2009	-0.125	0.016	***
2010	-0.140	0.017	***
2011	-0.128	0.016	***
2012	0.056	0.017	***
2013	-0.054	0.016	***
Number of Years in Training Dataset	0.001	0.001	
River Gages			***
<i>For the sake of brevity, the 82 river gages included in the regression as factors are omitted here.</i>			
$R^2$		0.26	
Adjusted $R^2$		0.25	
<i>P values: *** - &lt; 0.001; ** - 0.01; * - 0.05; . - 0.1</i>			

11311



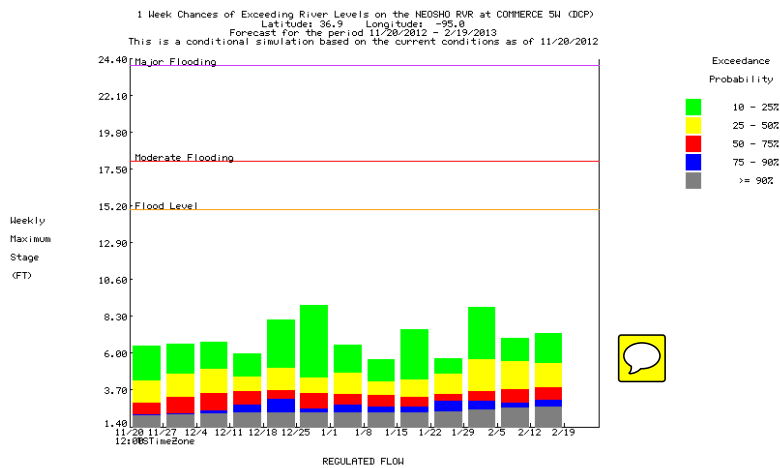
**Figure 1.** Deterministic short-term weather forecast in six hour intervals as published by the NWS for Hardin, IL on 24 April 2014. Source: <http://www.weather.gov/ahps2/hydrograph.php?wfo=lsx&gage=hari2> (last access: 1 October 2014)

11312



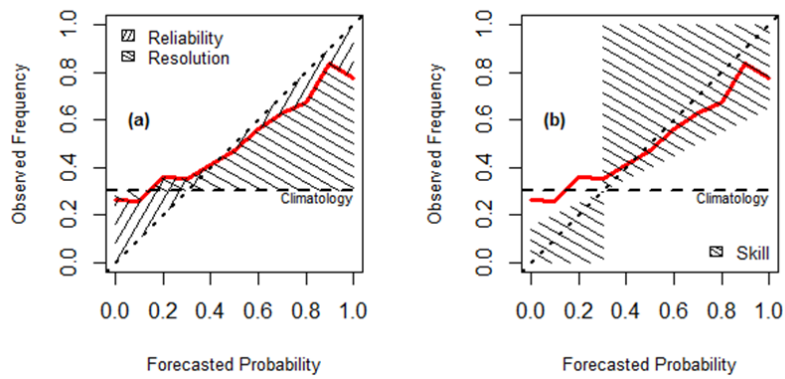
**Figure 2.** Probabilistic long-term forecast as published by the NWS for Commerce, OK on 14 December 2012: exceedance curve for three months period (Not available for Hardin, IL.) Source: [http://water.weather.gov/ahps2/probability\\_information.php?wfo=tsa&gage=COMO2&graph\\_id=2](http://water.weather.gov/ahps2/probability_information.php?wfo=tsa&gage=COMO2&graph_id=2) (last access: 1 October 2014).

11313



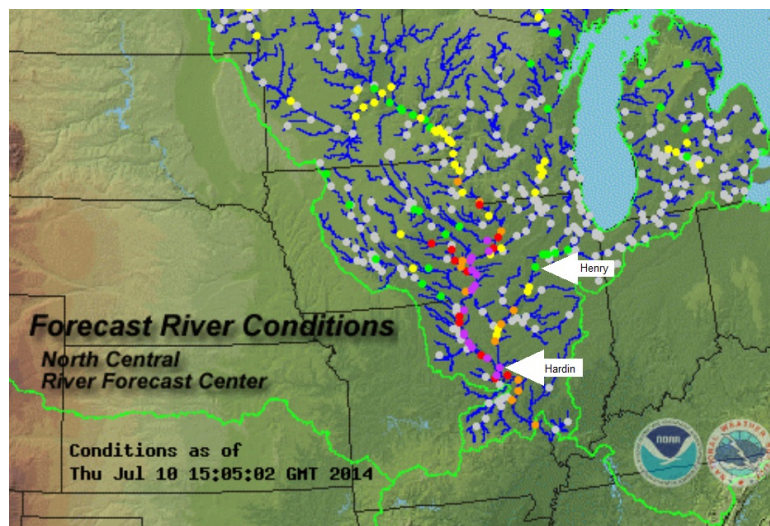
**Figure 3.** Probabilistic long-term forecast as published by the NWS for Commerce, OK on 14 December 2012: bar plot for each week of a three months period. (Not available for Hardin, IL.) Source: [http://water.weather.gov/ahps2/probability\\_information.php?wfo=tsa&gage=COMO2&graph\\_id=0](http://water.weather.gov/ahps2/probability_information.php?wfo=tsa&gage=COMO2&graph_id=0) (last access: 1 October 2014).

11314



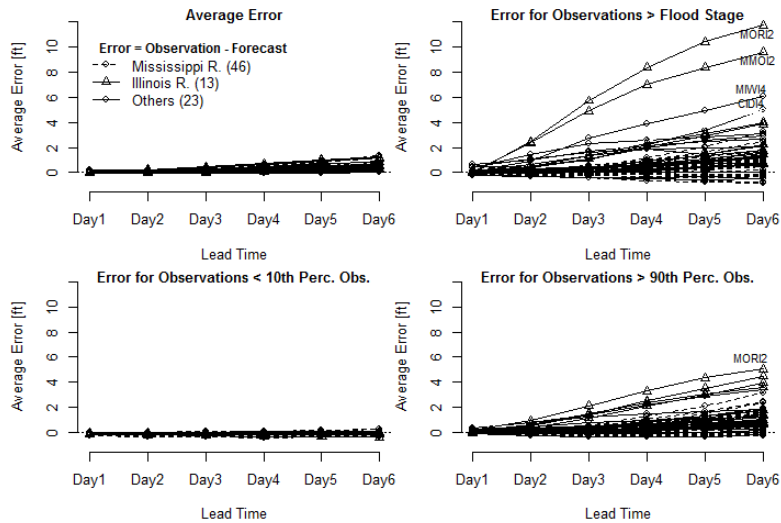
**Figure 4.** Theory behind Brier Skill Score illustrated for an imaginary forecast (red line): **(a)** reliability and resolution; **(b)** skill. In **(a)**, the area representing reliability should be as small, and for resolution as large as possible. The forecast has skill ( $BSS > 0$ ), i.e. performs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

11315



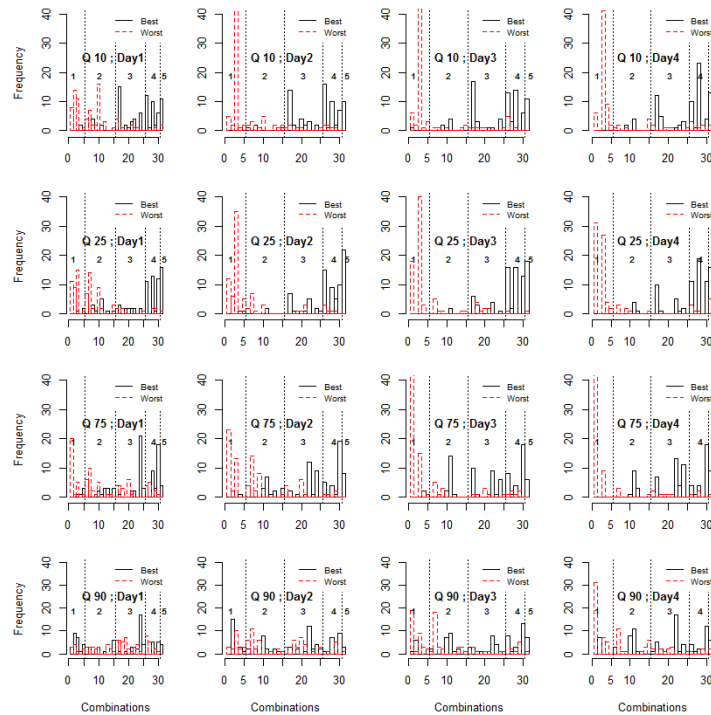
**Figure 5.** Portion of the North Central River Forecast Centers river gages with Henry (HYN12) and Hardin (HARI2) indicated by the upper and lower red arrow respectively. Source: <http://www.crh.noaa.gov/ncrfc/>.

11316



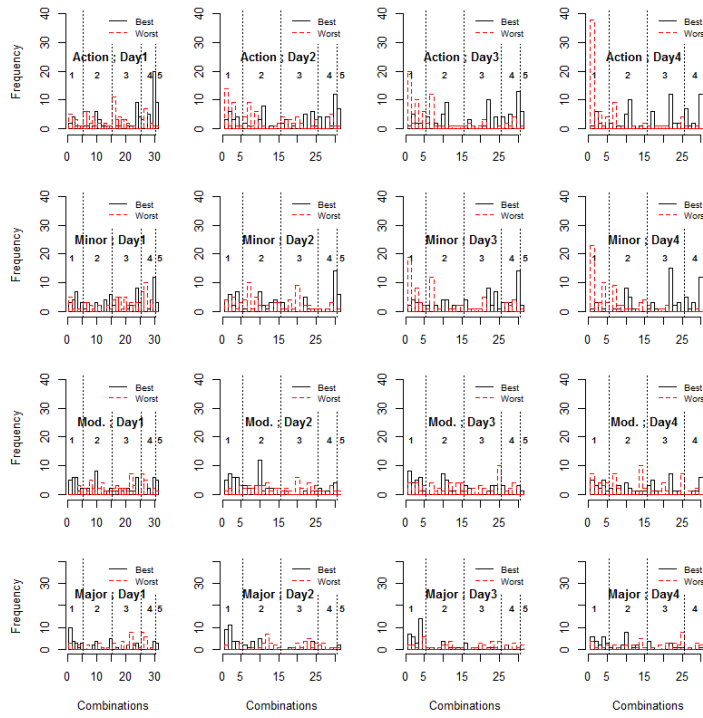
**Figure 6.** Forecast error for 82 river gages that the NCRFC publishes daily forecasts for. In anti-clockwise direction starting at the top left: (a) average error; (b) error on days that the water level did not exceed the 10th percentile of observations; (c) error on days that the water level exceeded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

11317



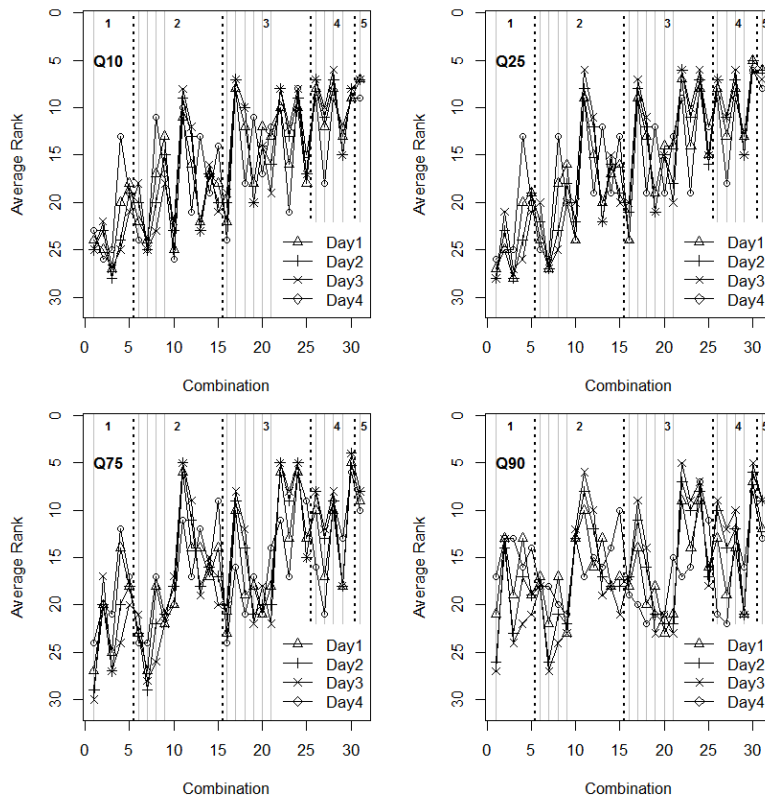
**Figure 7.** Histograms of variable combinations returning the best and worst Brier Skill Scores across 82 river gages. Each row of histograms refers to an event threshold defined as a percentile of the observed water levels, and each column to a lead time. The dotted vertical lines in the histograms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

11318



**Figure 8.** Histograms of variable combinations 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 combinations with different numbers of variables.

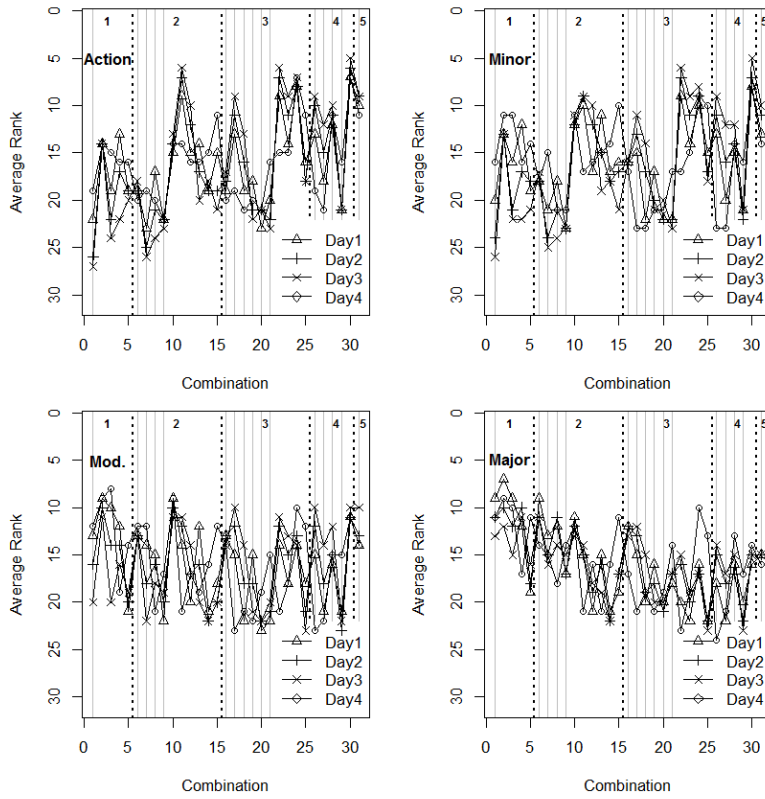
11319



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

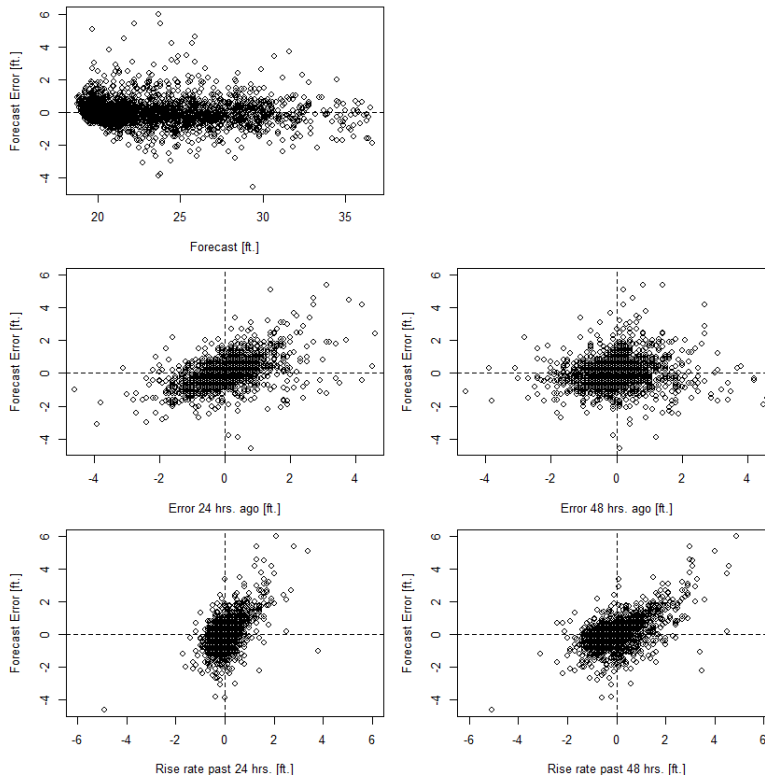
11320





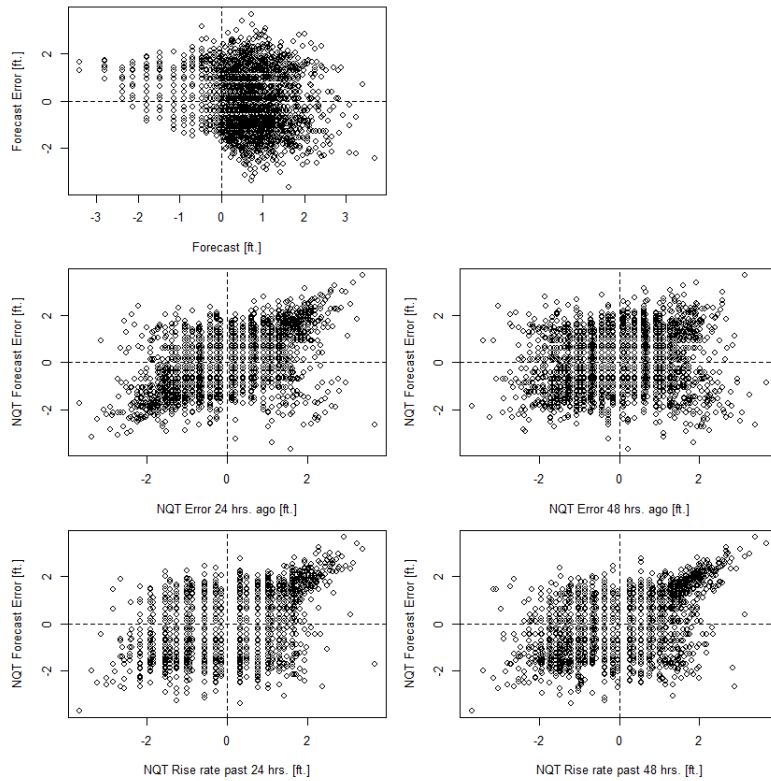
**Figure 10.** Average rank for each variable combination for one to four days of lead time and four flood stages. Vertical gray lines indicate variable combinations including the forecast.

11321



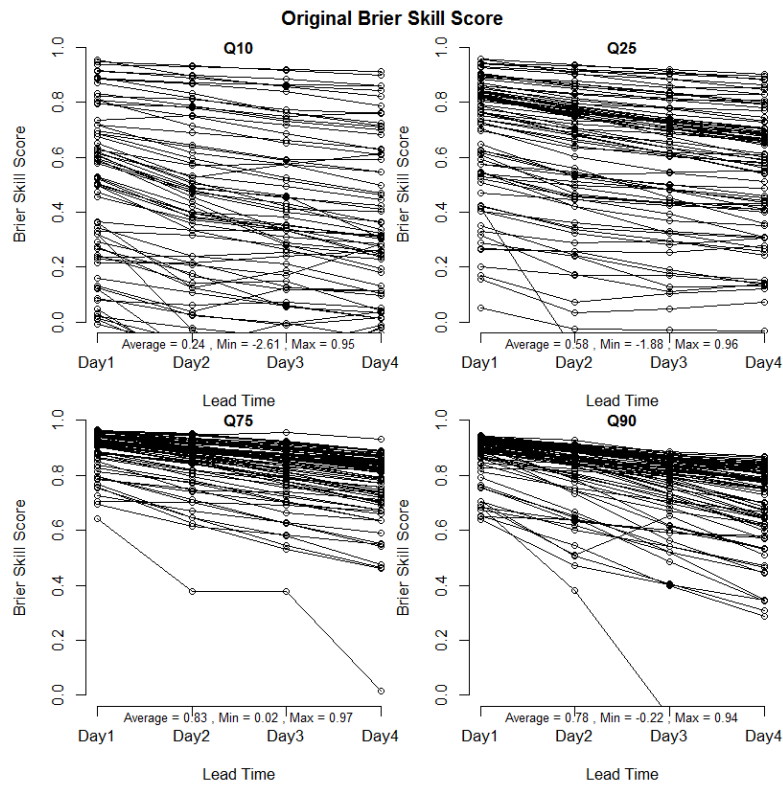
**Figure 11.** Independent variables plotted against the forecast error for Hardin IL with 3 days of lead time. First row: forecast; second row: past forecast errors; third row: rise rates.

11322



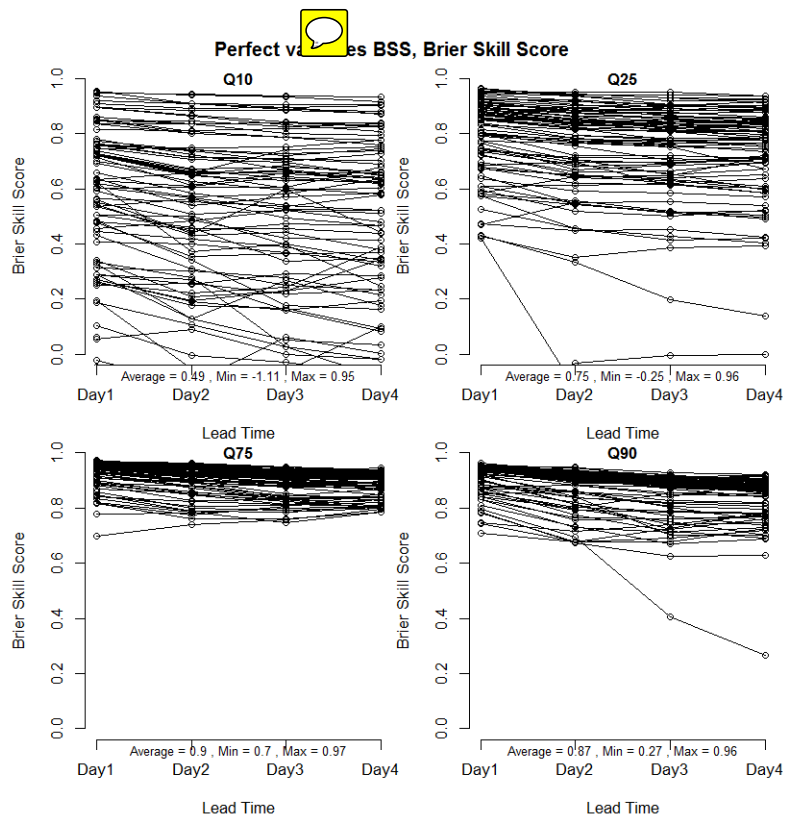
**Figure 12.** Independent variables after transforming into the Gaussian domain plotted against the forecast error for Hardin IL with 3 days of lead time. First row: forecast; second row: past forecast errors; third row: rise rates.

11323



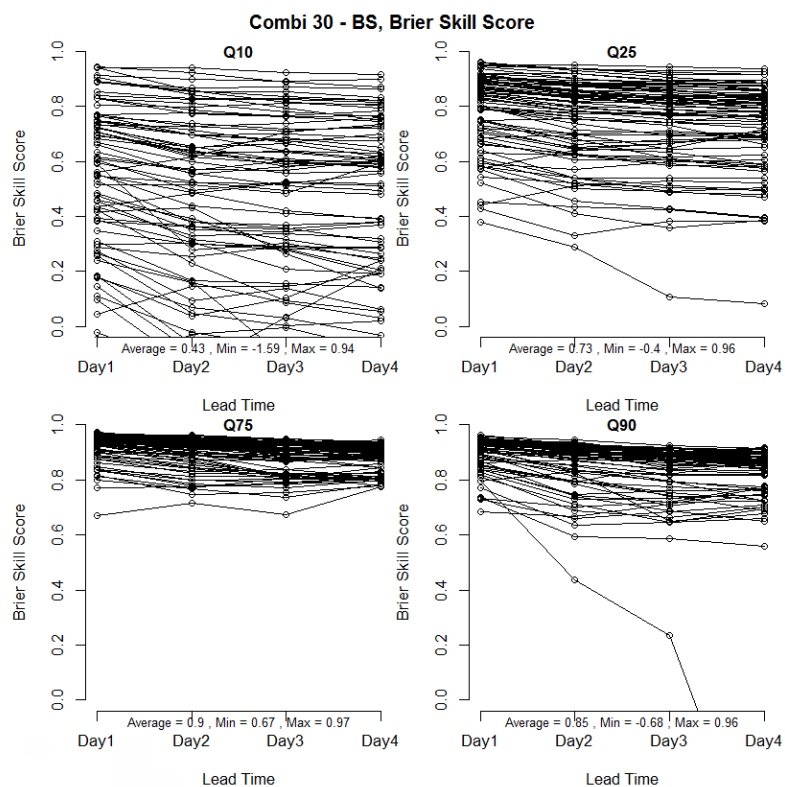
**Figure 13.** Brier Skill Scores of the original QR model (i.e., using the transformed forecast as the only independent variable) for four lead times and percentiles of observed water levels.

11324



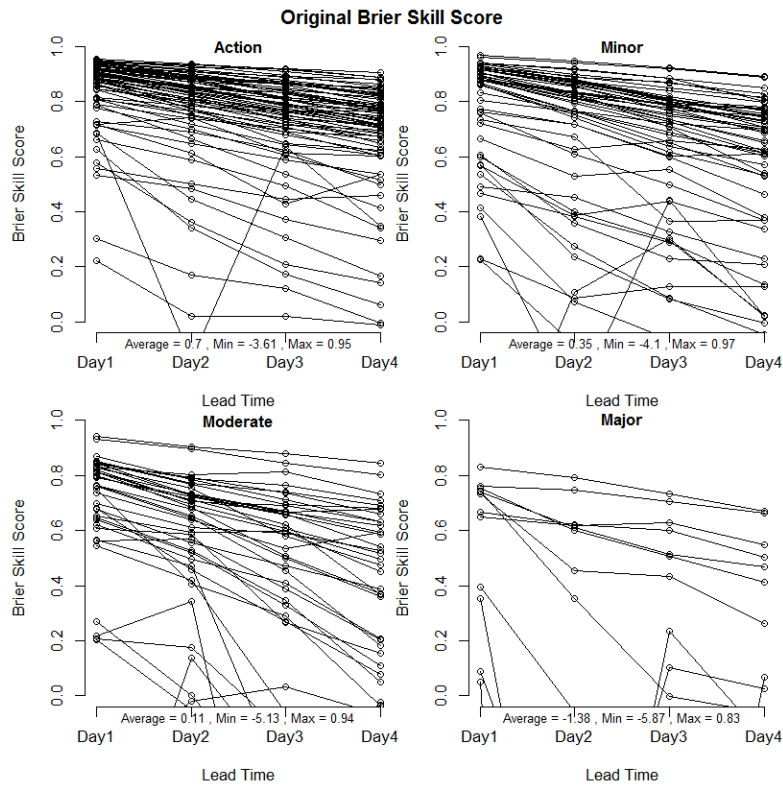
**Figure 14.** Brier Skill Scores for four lead times and percentiles of observed water levels using the best variable combination for each river gage as independent variables in the QR model.

11325



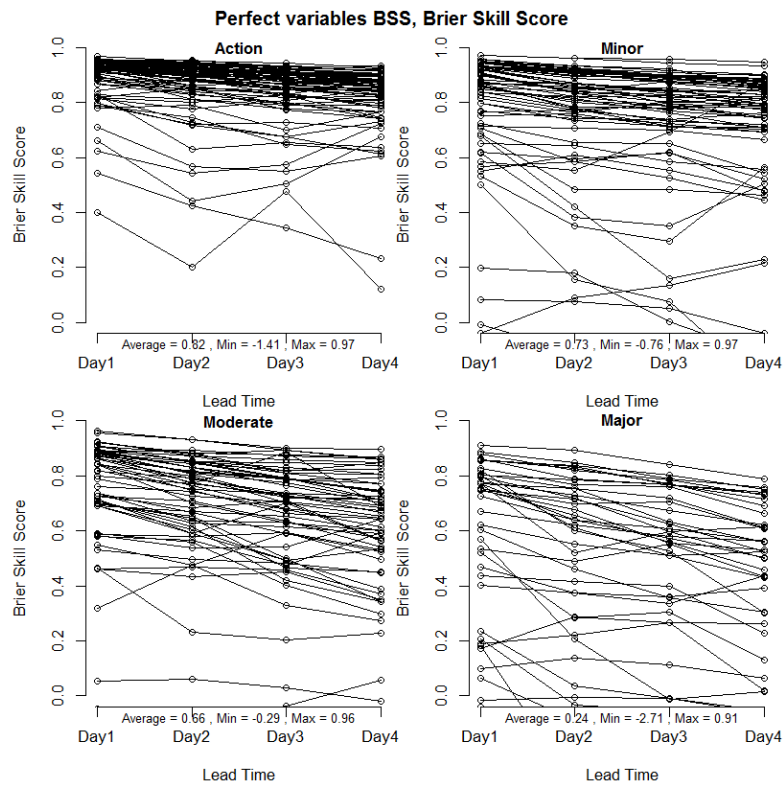
**Figure 15.** Brier Skill Scores for four lead times and percentiles of observed water levels using a one-size-fits-all approach (i.e., rr24, rr48, err24, err48) for the independent variables in the QR model.

11326



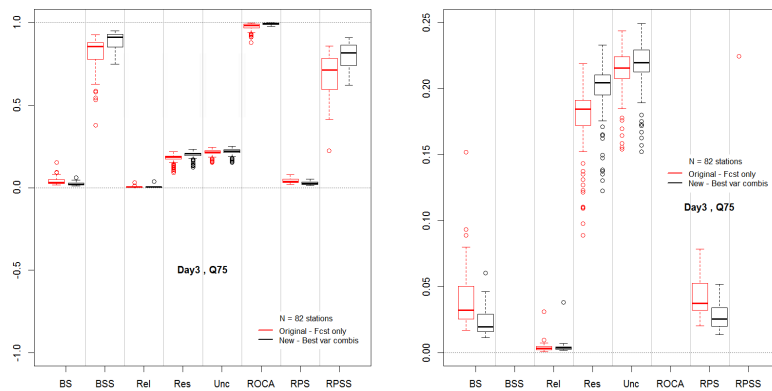
**Figure 16.** Brier Skill Scores of the original QR model (i.e., using the transformed forecast as the only independent variable) for four lead times and flood stages.

11327



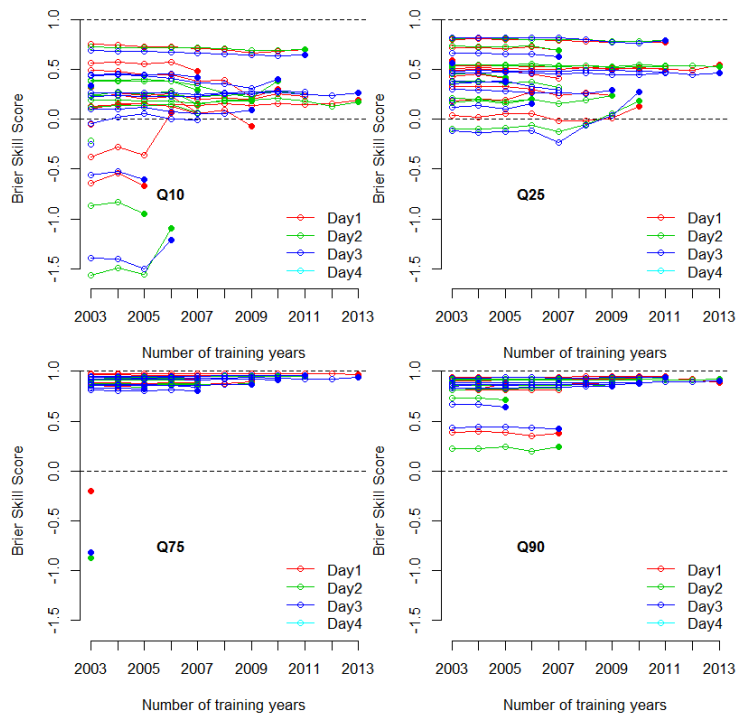
**Figure 17.** Brier Skill Scores for four lead times and flood stages of observed water levels using the best variable combination for each river gage as independent variables in the QR model.

11328



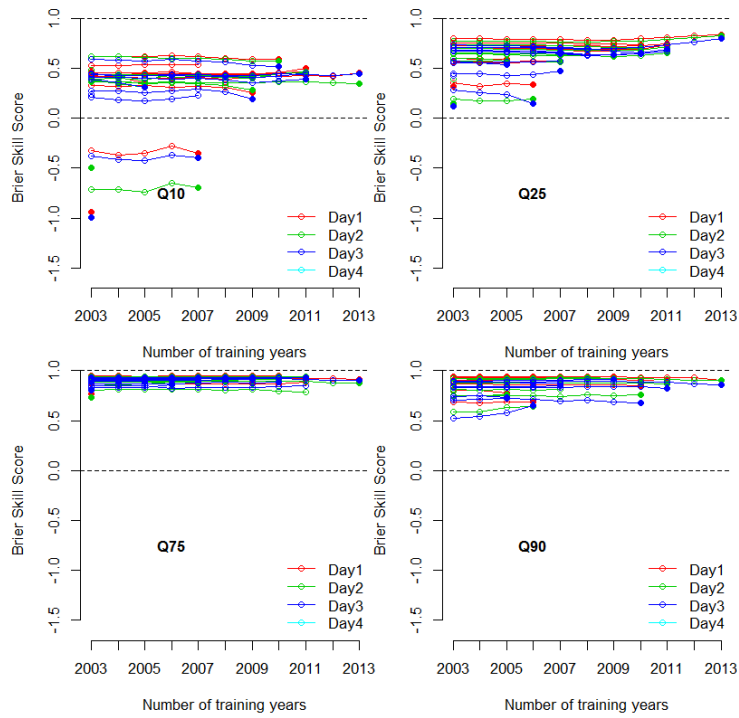
**Figure 18.** Comparison of the original QR model (i.e., only transformed forecast as independent variables) and the one-size-fits-all approach (i.e., rise rates and forecast errors as independent variables) using various measures of forecast quality: Brier Score (BS), Brier Skill Score (BSS), Reliability (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.

11329



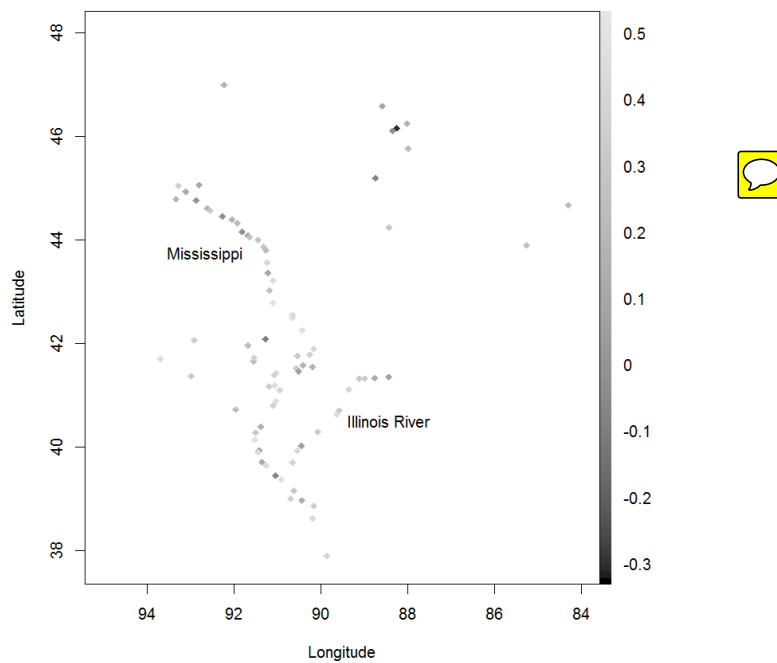
**Figure 19.** Brier Skill Score for various forecast years and various sizes of training dataset across different lead times (colors) and event thresholds (plots) for Hardin, IL (HARI2). The filled-in end point of each line indicates the BSS for the forecast year on the x axis with one year in the training dataset. Each point further to the left stands for one additional training year for that same forecast year.

11330



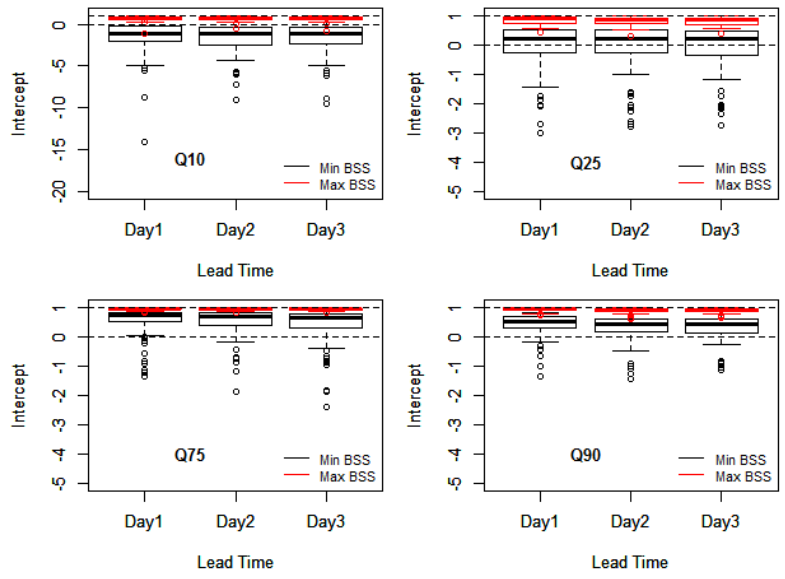
**Figure 20.** Brier Skill Score for various forecast years and various sizes of training dataset across different lead times (colors) and event thresholds (plots) for Henry, IL (HNYI2). The filled-in end point of each line indicates the BSS for the forecast year on the x axis with one year in the training dataset. Each point further to the left stands for one additional training year for that same forecast year.

11331



**Figure 21.** Geographical position of rivers. Colors indicate the regression coefficient of each station with the Brier Skill Score as dependent variable.

11332



**Figure 22.** Minimum (black) and maximum (red) Brier Skill Scores for various lead times and event thresholds across locations, size of training dataset and forecast years.