*Research Article*

# Statistical Uncertainty Estimation Using Random Forests and Its Application to Drought Forecast

## Junfei Chen,[1,2] Ming Li,[3] and Weiguang Wang[1]

[1] State Key Laboratory of Hydrology-Water Resources and Hydraulic Engineering, Hohai University, Nanjing 210098, China
[2] Business School, Hohai University, Nanjing 210098, China
[3] CSIRO Mathematics, Informatics and Statistics, Private Bag No. 5, Wembley, WA 6913, Australia

Correspondence should be addressed to Ming Li, li084@csiro.au

Academic Editor: Joao B. R. Do Val

Drought is part of natural climate variability and ranks the first natural disaster in the world. Drought forecasting plays an important role in mitigating impacts on agriculture and water resources. In this study, a drought forecast model based on the random forest method is proposed to predict the time series of monthly standardized precipitation index (SPI). We demonstrate model application by four stations in the Haihe river basin, China. The random-forest- (RF-) based forecast model has consistently shown better predictive skills than the ARIMA model for both long and short drought forecasting. The confidence intervals derived from the proposed model generally have good coverage, but still tend to be conservative to predict some extreme drought events.

## 1. Introduction

Drought is part of the natural variability of climate and its recurrence is inevitable and random. It affects nearly everywhere across all climate regions, though its features differ from region to region. Although specific definitions of drought depend on differences in disciplinary perspectives (e.g., meteorology, hydrology, agriculture, and socioeconomy), a typical definition of it originates from a deficiency of precipitation over a prolonged period, resulting in water supply shortage for some activity, group, or environmental sector. Drought reduces crop, livestock, and forest production and can result in widespread famine and death. Drought ranks the first amongst all natural hazards in terms of the number of people directly affected [1–3], and it is threatening nearly 50% of the most populated areas [4].

Drought forecasting plays an important role in taking contingency actions in advance of drought to mitigate its risk and impacts. A variety of forecasting methods have been developed to predict drought occurrence. Statistical run theory is the first attempt to predict drought likelihood [5–10]. Other stochastic models such as Markov chain [11, 12], renewal process [13, 14], and Poisson process [15, 16] also have been long suggested to characterise and predict drought events. Time series model is another widely approach for drought forecasting. For example, Chung and Salas [17] applied a low-order discrete autoregressive moving average model to estimate drought occurrence probabilities and the risk of dependent hydrological processes. Mishra and Desai [18] used seasonal and nonseasonal autoregressive integrated moving average (ARIMA) models to forecast the standardized precipitation index (SPI) in the Kansabati river basin in India. They showed that the predicted SPI agreed with observations with one- to two-month lead time but the prediction performance decreased with increasing lead time. Durdu [19] showed that the method of Mishra and Desai [18] can be used for drought forecast with reasonably accuracy upto two months in the Buyuk Menderes river basin in Turkey. Moreover, many other forecasting methods have been developed to address nonlinearity and nonstationarity of time series of drought events. Kim et al. [20] evaluated the application of the nonparametric kernel smoothing to estimate return periods of drought in arid regions. Kim et al. [21] proposed a conjunction model to forecast drought based on wavelet transformation and neural networks. Mishra and Desai [22] showed that a feed-forward recursive neural network provided better prediction of SPI than the ARIMA model with one-month lead time.

Drought forecasting remains challenging and is subject to great uncertainty partly due to anticipating some parts of hydrologic cycle (e.g., rainfall, soil moisture, and groundwater level). It is, therefore, desired to make single-value predictions as well as reliable uncertainty measures, especially for long lead times. Uncertainty estimation of drought forecast can largely improve decision making on water resources management. All researches mentioned in the previous paragraph only focus on predicting the *mean* of drought stage and did not include any analysis of uncertainty. Consideration of forecast uncertainty has been shown to provide valuable information to decision makers who need to understand the variability of upcoming drought status. For example, Carbone and Dow [23] resampled historical monthly temperature and precipitation and derived an ensemble forecast of the Palmer drought severity index (PDSI) with weighted resamples. Hwang and Carbone [24] applied the resampling strategy suggested by Carbone and Dow [23] to the residuals of a predictive model of drought indices and generate drought ensemble forecasts. They found the predictive model based on nonparametric autoregressive models had good forecast capability of SPI with up to 3-month lead time in terms of mean forecast. They also demonstrated that variability of the forecast ensemble members shared probability density functions that were similar to those of the observations.

The objective of this research is to introduce a statistical method, called *Random Forest* (RF), to predict drought events together with an appropriate estimation on uncertainty measures. The RF method does not restrict a particular relationship within the drought index series or make any distributional assumptions on the model error. An ensemble of equally probable realisations of drought indices is generated by the proposed model and differences amongst the realisations can be used as a measure of uncertainty. Figure 1 illustrates how the RF model is superior over the ARIMA model by an example. This paper is organised as follows. Section 2 describes the study region and data used in this study. The predictive model based on the RF method is introduced in Section 3, followed by the study results reported in Section 4. Summary and further research recommendations are presented in Section 5.
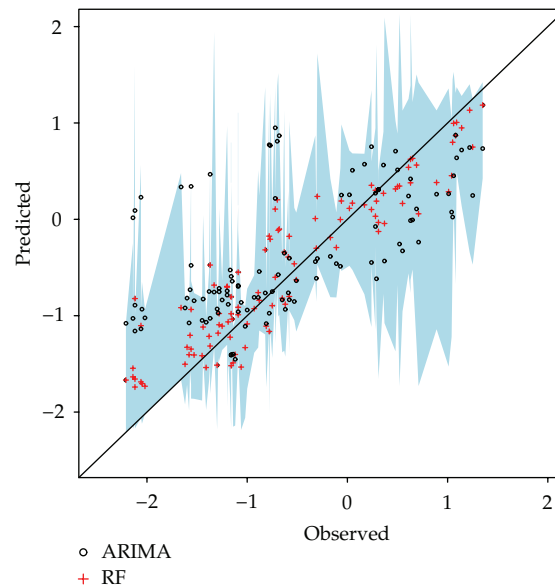
**Figure 1:** Comparison of ARIMA and RF predicted SPI(12) of Beijing.

## 2. Study Region and Data

Flowing through many big cities including Beijing and Tianjin, the Haihe river basin is the largest and the most important water system in northern China. However, the water usage of the Haihe river basin has a tremendous conflict between water supply and demand. The water shortages in this region have especially intensified due to the rapid development of the economy and the explosion of the city population during recent years and cause many severe environmental and ecological problems, such as drought, drying up of river systems, and degradation of lakes and wetlands [25, 26]. The monthly precipitation observations in the period of 1966–2004 recorded from four stations in the Haihe river basin including Beijing, Shijiazhuang, Tangshan, and Tianjin were used in this study (Figure 2).

The SPI [27] is the drought index used in this study and it measures the severity of drought over different time scales. It is a dimensionless index and is only defined by historical rainfall observations. In contrast to the PDSI, the SPI is not adversely affected by topography, but can provide indication or identification of drought a few months sooner [28]. The SPI indicates the severity of drought in a large scale and provides little details for a particular drought event. The critical drought duration is generally characterized by a stochastic process, such as a second-order Markov chain [29]. The 3-month and 12-month, SPI, denoted by SPI(3) and SPI(12), respectively, were considered in the present work. In particular, forecasting using SPI(3) with a one-month lead time and SPI(12) with a six-month lead time were made for short-term and long-term drought forecasting, respectively. Forecasting models were developed based on the data from 1966 to 1995, and predictions one month ahead for SPI(3) or six month ahead for SPI(12) were made from 1996 to 2004. SPI based on a month other than 3 and 12 months can be forecasted.
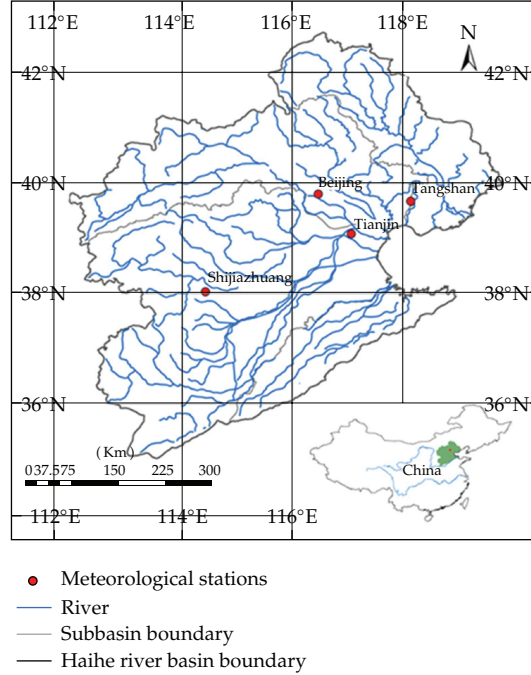
**Figure 2:** Locations of four selected stations in the Haihe river basin.

## 3. Model

Let $\text{SPI}_m$ represent the SPI of a specific month $m$, and the predictive model of $\text{SPI}_m$ can be written as

$$\text{SPI}_m = f\left(\text{SPI}_{m-l}, \text{SPI}_{m-l-1}, \ldots, \text{SPI}_{m-l-(d-1)}\right), \tag{3.1}$$

where $l$ is the forecast lead time and $d$ is the order of nonparametric autoregressive model. In this model, SPI is predicted from the SPI of previous months. If the function $f$ is a linear function, the model defined by (3.1) is an autoregressive model [30]. No specific functional form of $f$ is assumed in the proposed model and this will reduce the risk of model misspecification. We will apply the RF method to estimate $f$ for all months. In the rest of this section, a brief description of RF will be presented.

### 3.1. Regression Trees

Regression trees, often described in graphical and biological terms, use a tree structure to predict new data from training data. In contrast to linear regression, which is a global model and applies a single predictive formula holding over the entire data space, regression trees partition space into smaller regions that have the most homogeneous collection of outcomes. The partitioning is repeated on each derived sub-data space in a recursive manner, called recursive partitioning, and an optimality criterion is used to guide each partition. This process is often likened to the way a tree divides into smaller branches (called nodes) and eventually

to a single leaf (called terminal node), hence the name of the method. The terminal nodes would comprise a collection of closely matched previous observations found from recursive partitioning, and the fitted response is the average of the observations for that node. Figure 3 provides a simple example of regress tree to predict SPI and circles denote nodes and boxes denote terminal nodes.

The basic regression-tree-growing strategy involves in at least three fundamental questions: splitting rules, terminal nodes, and terminal node assignment.

*Splitting Rules*

On what criteria are splits to be made? Creating splits is similar to variable selection in regression. For regression problems, the variable and the location of a split are chosen by the sum of squared error between the observations and the mean of the observations within each node. For example, the root node $SPI_{m-1}$ was split to two branches and a threshold value 2.1 was used as a splitting rule shown as Figure 3.

*Terminal Nodes*

When to stop a tree from growing? The simplest terminal node is the node where all training data are from a single class or a single value of response. This choice may make a tree grow so large that it fits all training data perfectly. This tends to overfit training data and has potentially poor predictions on independent test data. If a tree is too small, the relationship between predictors and predicants may not be extracted completely. Choosing an appropriate tree size is thus of great importance. A common approach is to grow an overly large tree so that a minimum node size is reached and prune the tree back to the optimal size by an independent test data or cross-validation. The example regression tree in Figure 3 has four terminal codes.

*Terminal Node Assignment*

Which value is to be assigned to the terminal nodes? For regression problems, values at the terminal node are generally assigned by the mean of the observations in that node. The assignment rules can be modified to reflect costs of explanatory variables and misspecification, if necessary. The values for terminal codes in Figure 3 are −1, 1.5, 0.2, and 0.5, respectively.

More technical aspects regarding regression trees can be found in Hastie et al. [31], Chapter 9, for example.

### 3.2. Random Forest

A simple regression tree has the disadvantage of being sensitive to training data used to build the tree, especially when the size of training data is small. Little change in training data may result in very different trees and predictions [32]. Ensemble trees are one solution to improve robustness and reliability of regression trees by "randomly" growing a collection of trees from bootstrap samples (i.e., training data randomly drawn with replacement) and aggregating predictions. RF [32–34] is one of the most popular ensemble tress methods and has been extensively applied in medicine, neuroscience, bioinformatics (e.g., [35–37]), and
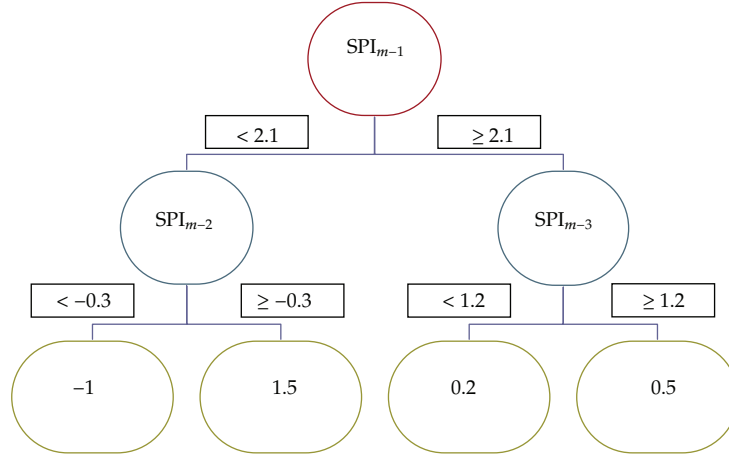
**Figure 3:** A simple example of regression tree.

operations research (e.g., [38]). The approach has been used to only a few applications in environmental science (e.g., [39, 40]). In the following paragraphs, we briefly describe the algorithm of RF in the context of estimating $f_m$ in (3.1).

For a given month $m$, denote a training data as $X = (x_1, x_2, \ldots, x_n)$, where $x_i = (\mathrm{SPI}^i_m, \mathrm{SPI}^i_{m-l}, \mathrm{SPI}^i_{m-l-1}, \ldots, \mathrm{SPI}^i_{m-l-(d-1)})$ and $n$ is the total number of the training data. The function $f_m$ in (3.1) is estimated from the following steps.

(a) Bootstrap sampling: draw $B$ random samples of size $n$ : $X^*_1, X^*_2, \ldots, X^*_B$, where $X^*_j = (x^*_{1,j}, x^*_{2,j}, \ldots, x^*_{n,j})$ with replacement from the entire training set. That is $P(x^*_{i,j} = x_1) = \cdots = P(x^*_{i,j} = x_n) = 1/n$ for any $i = 1, \ldots, n$ and $j = 1, \ldots, B$.

(b) Random-forest tree growing: grow an ensemble of $B$ random-forest trees based on bootstrap samples $X^*_1, X^*_2, \ldots, X^*_B$ by repeating the following substeps for each node until the minimum number of nodes (called the minimum node size) is reached.

    (b.1) Randomly select a subset of predictors of size $p$ : $(\mathrm{SPI}_{m-l-k_1}, \mathrm{SPI}_{m-l-k_2}, \ldots, \mathrm{SPI}_{m-l-k_p})$ out of the complete predictor set of size $d$, where $k_i \in (0, \ldots, d-1)$ for $i = 1, \ldots, p$.

    (b.2) Pick the best split predictor $\mathrm{SPI}_{m-l-k}$ together with the best split value in the sense of minimising the mean square error from the subset of predictors selected at (b.1).

    (b.3) Split the node based on $\mathrm{SPI}_{m-l-k}$ into two branches according to the best split predictor/value selected at (b.2).

(c) Ensemble averaging: the RF tree from the $j$th bootstrap sample $X^*_j$ provides one prediction $\mathrm{SPI}^{j*}_m$, where $j = 1, \ldots, B$. We use $(\mathrm{SPI}^{1*}_m, \ldots, \mathrm{SPI}^{B*}_m)$, as an ensemble forecast of SPI for month $m$ and $\mathrm{SPI}^*_m = (1/B) \sum_{j=1}^{B} \mathrm{SPI}^{j*}_m$ is reported as the mean forecast.

The confidence interval of $\mathrm{SPI}_m^{j*}$ at a nominal level of $\alpha$ (e.g. 0.95) is defined by $[\mathrm{SPI}_L(\alpha), \mathrm{SPI}_U(\alpha)]$, where

$$P\left\{\mathrm{SPI}_m^{j*} \leq \mathrm{SPI}_L(\alpha)\right\} = P\left\{\mathrm{SPI}_m^{j*} \geq \mathrm{SPI}_U(\alpha)\right\} = \frac{1-\alpha}{2}. \tag{3.2}$$

From the algorithm described above, only two parameters $p$ (the number of predictors randomly selected at each node) and $B$ (the number of ensemble trees) are required to specify to implement the RF method. In this study, we choose $p = d/3$, where $d$ is the number of total predictors (i.e., the total number of previous SPI used in prediction) and $B$ is chosen to be 500. An RF $R$ package (available at http://cran.r-project.org/web/packages/randomForest) is used in this study.

## 4. Results

In order to demonstrate model performance, the ARIMA model is considered as a baseline and is compared with the model based on RF. As described in Section 2, the number of the previous SPI used in the predictive model is determined by the best fitted ARIMA model selected from the Akaike information criterion (AIC). The performance difference between the RF-based model and the best fitted ARIMA reflects the validity of the assumption of the ARIMA model, such as linearity and stationarity. Three well-known error statistics were calculated to measure the difference between the observed and predicted SPI series, including bias, mean absolute error (MAE), and root mean-squared error (RMSE) and they are defined by

$$\mathrm{Bias} = \frac{1}{N}\sum_{i=1}^{N}(\mathrm{EST}_i - \mathrm{OBS}_i), \tag{4.1}$$

$$\mathrm{MAE} = \frac{1}{N}\sum_{i=1}^{N}|\mathrm{EST}_i - \mathrm{OBS}_i|, \tag{4.2}$$

$$\mathrm{RMSE} = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(\mathrm{EST}_i - \mathrm{OBS}_i)^2}, \tag{4.3}$$

where $\mathrm{EST}_i$ and $\mathrm{OBS}_i$ denote the $i$th estimated and observed values, respectively. Furthermore, according to the weather classification of McKee et al. [27], dry days are defined by the days with SPI values falling below −1. We thus considered another two error statistics based on dry days: RMSE at dry days and the proportion of dry days detected (i.e., SPI predictions less than −1 at dry days). These two additional error statistics are intended to evaluate whether the RF is a more effective and efficient drought prediction tool than the ARIMA.

All error statistics of one month ahead SPI(3) prediction and six month ahead SPI(12) prediction for each stations are presented in Table 1. We have the following findings. (1) In general, the RF performed consistently better than the ARIMA. In particular, all error statistics from the RF were smaller than those from the ARIMA, except that the biases of SPI(3) predictions obtained from both methods for Beijing were almost equal. (2) RMSE at dry days was greater than overall RMSE for all predictions, indicating the difficulty of predicting drought

**Table 1:** Error statistics of SPI forecast based on ARMA model and the RF model.

| Station | Error statistics | SPI(3) | | SPI(12) | |
|---|---|---|---|---|---|
| | | ARMA | RF | ARMA | RF |
| Beijing | BIAS | 0.033 | 0.034 | 0.245 | 0.092 |
| | MAE | 0.637 | 0.415 | 0.606 | 0.3 |
| | RMSE | 0.791 | 0.526 | 0.796 | 0.385 |
| | RMSE (dry days) | 0.963 | 0.662 | 0.968 | 0.44 |
| | % dry days detected | 27.3% | 54.5% | 31.1% | 71.1% |
| Shijiazhuang | BIAS | −0.019 | 0.002 | 0.034 | −0.022 |
| | MAE | 0.557 | 0.289 | 0.774 | 0.396 |
| | RMSE | 0.736 | 0.382 | 1.015 | 0.507 |
| | RMSE (dry days) | 1.065 | 0.531 | 1.551 | 0.674 |
| | % dry days detected | 35.0% | 60.0% | 7.1% | 42.9% |
| Tangshan | BIAS | 0.074 | 0.05 | 0.214 | 0.087 |
| | MAE | 0.656 | 0.434 | 0.796 | 0.378 |
| | RMSE | 0.814 | 0.54 | 1.034 | 0.475 |
| | RMSE (dry days) | 1.036 | 0.651 | 1.247 | 0.553 |
| | % dry days detected | 33.3% | 50.0% | 17.8% | 60.0% |
| Tianjin | BIAS | 0.069 | 0.015 | 0.211 | 0.078 |
| | MAE | 0.597 | 0.288 | 0.523 | 0.264 |
| | RMSE | 0.735 | 0.366 | 0.624 | 0.331 |
| | RMSE (dry days) | 1.041 | 0.526 | 0.897 | 0.467 |
| | % dry days detected | 16.7% | 50.0% | 0.0% | 19.0% |

events accurately. The difference between RMSE and RMSE at dry days of the RF-based model was smaller than that of the ARIMA for each prediction. This suggests that the RF-based model is more robust in predicting dry events. (3) The RF-based model is even more robust for longer term prediction. The longer-term drought forecast typically involves more uncertainty and thus is more challenging to predict. The ARIMA indeed lost the predictive capability for SPI(12) prediction. For example, none of dry days in Tianjin indicated by SPI(12) was forecasted by the ARIMA. Instead, the accuracy of SPI prediction based on the RF was less affected by a longer lead time. In particular, at three out of four stations (except for Shijiazhuang) the RF led to comparable and even smaller prediction errors indicated by five error statistics.

For a graphical illustration, the SPI(3) and SPI(12) predictions by the RF method together with the associated 95% confidence intervals are shown in Figures 4 and 5, respectively. Most predictions of SPI agreed with observations very well. For example, from Figure 4, a few extreme drought events with SPI(3)< −1.5 for Shijiazhuang were well forecasted by the RF method but not ARIMA. However, a number of extreme drought events identified by SPI < −1.5 still fell outside of the 95% confidence interval. It is evident that the confidence intervals for the one month ahead SPI(3) prediction were more narrow than that for the six month ahead SPI(12), because more uncertainty is expected for the forecast with a longer lead time.
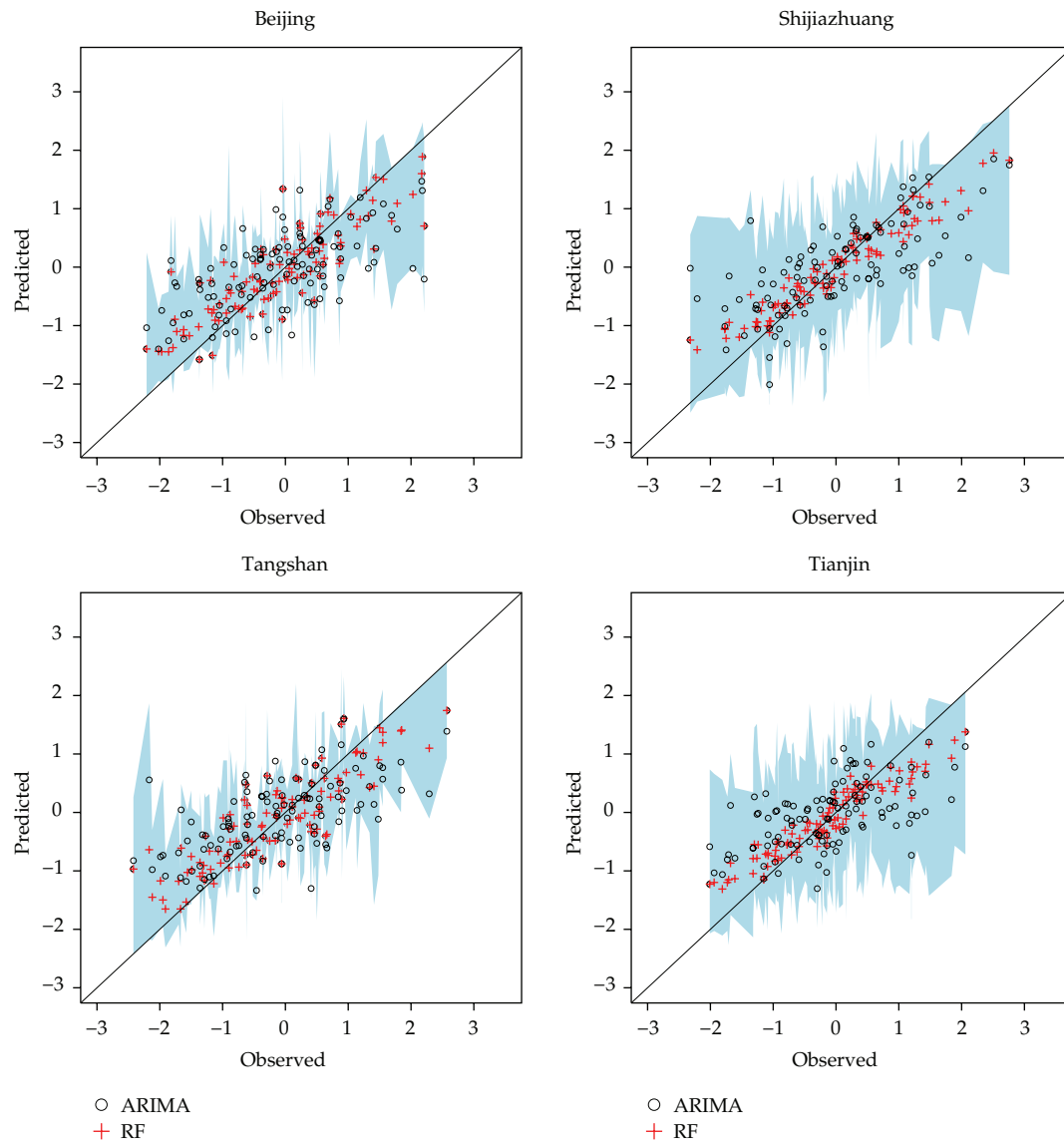
**Figure 4:** Observed and predicted SPI(3) by the RF and ARIMA and the 95% confidence intervals constructed by the RF.

## 5. Conclusion

In this study, a drought forecast model based on RF is proposed to predict SPI from the SPI in previous months. Unlike traditional time series models such as the ARIMA, the forecast model is built on a nonparametric framework and is more flexible to capture the underlying relationship. The RF-based model has another advantage of generating ensemble of drought forecast rather than a mean prediction. A confidence interval based on the ensemble forecast can be served as a measure of forecast uncertainty. The performance of the proposed forecast model has been demonstrated by its applications to four stations in the Haihe river basin,
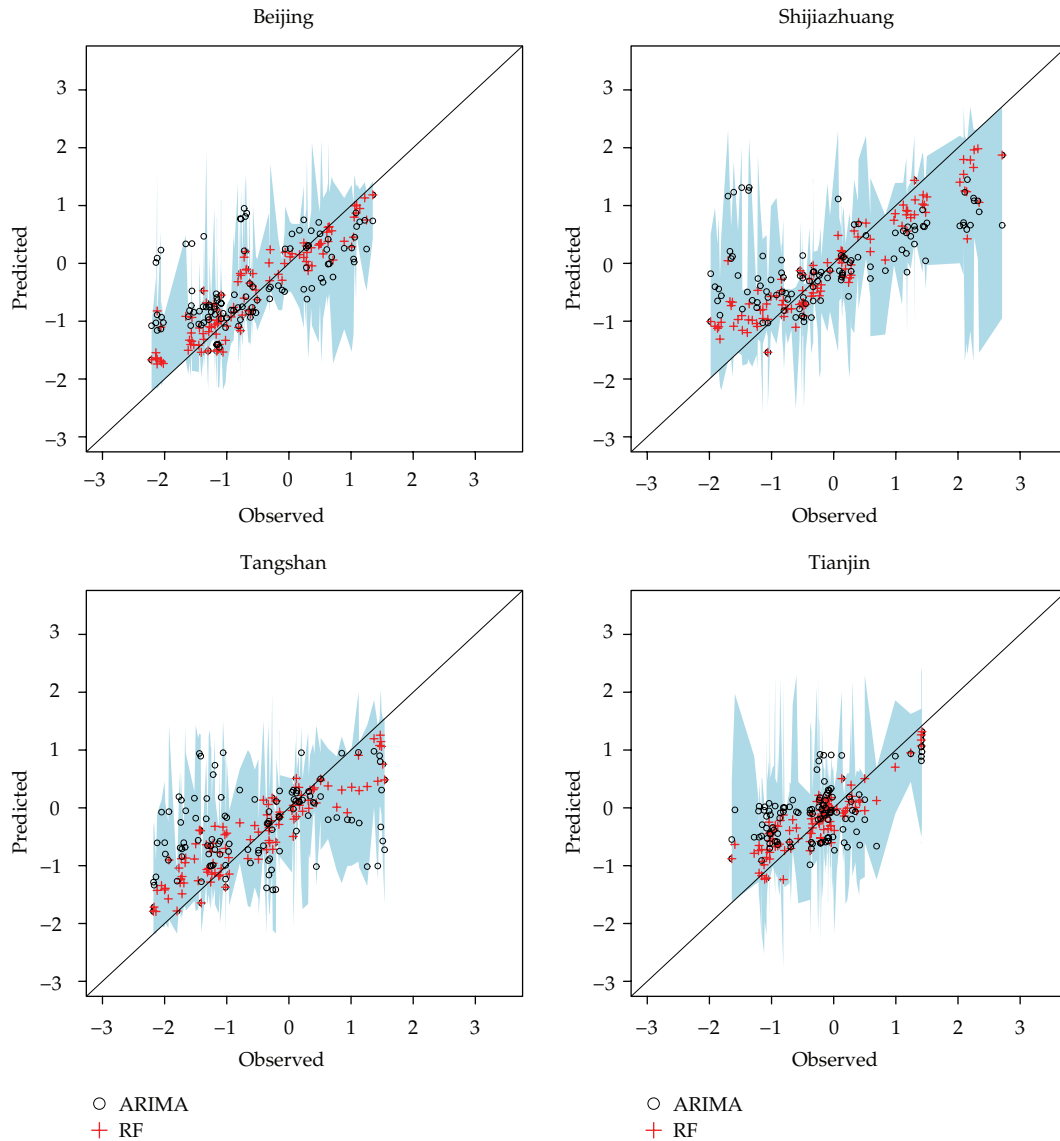
**Figure 5:** Observed and predicted SPI(12) by the RF and ARIMA and the 95% confidence intervals constructed by the RF.

China. Compared with the ARIMA, the RF-based predictive model is more reliable and efficient for both short- and long-term drought forecasting. The 95% confidence interval derived from the ensemble forecast covers nearly all observations reasonably well, though a few extreme drought events are identified outside the specified range. Further potential improvement of the drought forecast skill may be made by introducing useful climate indices and the outputs from climate models to the RF-based predictive model.

## Acknowledgments

## References

[1] G. O. P. Obasi, "WMO's role in the international decade for natural disaster reduction," *Bulletin of the American Meteorological Society*, vol. 75, no. 1, pp. 655–1661, 1994.

[2] K. Hewitt, *Regions at Risk. A Geographical Introduction to Disasters*, Addison-Wesley Longman, Essex, UK, 1997.

[3] D. A. Wilhite, *Drought: A Global Assessment*, Natural Hazards and Disasters Series, Routledge, London, UK, 2000.

[4] USDA, *Major World Crop Areas and Climatic Profiles*, Agricultural Handbook, no. 664, World Agricultural Outlook Board, U.S. Department of Agriculture, 1994.

[5] Y. Yevjevich, "An objective approach to definitions and investigations of continental hydrologic droughts," Hydrological Paper, Colorado State University, Fort Collins, Colo, USA, 1967.

[6] J. Saldariaga and V. Yevjevich, "Application of run-lengths to hydrologic series," Hydrological Paper, Colorado State University, Fort Collins, Colo, USA, 1970.

[7] Z. Sen, "Wet and dry periods of annual streamflow series," *Journal of the Hydraulics Division*, vol. 102, no. 10, pp. 1503–1514, 1976.

[8] Z. Şen, "Run-sums of annual flow series," *Journal of Hydrology*, vol. 35, no. 3-4, pp. 311–324, 1977.

[9] L. A. Moyé, A. S. Kapadia, I. M. Cech, and R. J. Hardy, "The theory of runs with applications to drought prediction," *Journal of Hydrology*, vol. 103, no. 1-2, pp. 127–137, 1988.

[10] L. A. Moyé and A. S. Kapadia, "Predictions of drought length extreme order statistics using run theory," *Journal of Hydrology*, vol. 169, no. 1–4, pp. 95–110, 1995.

[11] V. K. Lohani and G. V. Loganathan, "An early warning system for drought management using the Palmer drought index," *Journal of the American Water Resources Association*, vol. 33, no. 6, pp. 1375–1386, 1997.

[12] V. K. Lohani, G. V. Loganathan, and S. Mostaghimi, "Long-term analysis and short-term forecasting of dry spells by Palmer drought severity index," *Nordic Hydrology*, vol. 29, no. 1, pp. 21–40, 1998.

[13] H. A. Loaiciga and R. B. Leipnik, "Stochastic renewal model of low-flow streamflow sequences," *Stochastic Hydrology and Hydraulics*, vol. 10, no. 1, pp. 65–85, 1996.

[14] H. A. Loaiciga, "On the probability of droughts: the compound renewal model," *Water Resources Research*, vol. 41, Article ID W01009, 8 pages, 2005.

[15] T. J. Chang, "Effects of drought on streamflow characteristics," *Journal of Irrigation and Drainage Engineering*, vol. 116, no. 3, pp. 332–341, 1990.

[16] A. C. Cebrián and J. Abaurrea, "Drought analysis based on a marked cluster Poisson model," *Journal of Hydrometeorology*, vol. 7, no. 4, pp. 713–723, 2006.

[17] C. H. Chung and J. D. Salas, "Drought occurrence probabilities and risks of dependent hydrologic processes," *Journal of Hydrologic Engineering*, vol. 5, no. 3, pp. 259–268, 2000.

[18] A. K. Mishra and V. R. Desai, "Drought forecasting using stochastic models," *Stochastic Environmental Research and Risk Assessment*, vol. 19, no. 5, pp. 326–339, 2005.

[19] Ö. F. Durdu, "Application of linear stochastic models for drought forecasting in the Büyük Menderes river basin, Western Turkey," *Stochastic Environmental Research and Risk Assessment*, vol. 24, no. 8, pp. 1145–1162, 2010.

[20] T. W. Kim, J. B. Valdes, and C. Yoo, "Nonparametric approach for estimating return periods of droughts in arid regions," *Journal of Hydrologic Engineering*, vol. 8, no. 5, pp. 237–246, 2003.

[21] T. W. Kim and J. B. Valdes, "Nonlinear model for drought forecasting based on a conjunction of wavelet transforms and neural networks," *Journal of Hydrologic Engineering*, vol. 8, no. 6, pp. 319–328, 2003.

[22] A. K. Mishra and V. R. Desai, "Drought forecasting using feed-forward recursive neural network," *Ecological Modelling*, vol. 198, no. 1-2, pp. 127–138, 2006.

[23] G. J. Carbone and K. Dow, "Water resource management and drought forecasts in South Carolina," *Journal of the American Water Resources Association*, vol. 41, no. 1, pp. 145–155, 2005.

[24] Y. Hwang and G. J. Carbone, "Ensemble forecasts of drought indices using a conditional residual resampling technique," *Journal of Applied Meteorology and Climatology*, vol. 48, no. 7, pp. 1289–1301, 2009.

[25] W. Wang, S. Peng, T. Yang, Q. Shao, J. Xu, and W. Xing, "Spatial and temporal characteristics of reference evapotranspiration trends in the Haihe River Basin, China," *Journal of Hydrologic Engineering*, vol. 16, no. 3, pp. 239–252, 2011.

[26] W. Wang, Q. Shao, S. Peng et al., "Spatial and temporal characteristics of changes in precipitation during 1957–2007 in the Haihe River basin, China," *Stochastic Environmental Research and Risk Assessment*, vol. 25, no. 7, pp. 881–895, 2011.

[27] T. B. McKee, N. J. Doesken, and J. Kliest, "The relationship of drought frequency and duration to time scales," in *Proceedings of the 8th Conference on Applied Climatology, Anaheim, CA*, pp. 179–184, American Meteorological Society, Boston, Mass, USA, 1993.

[28] M. J. Hayes, M. D. Svoboda, D. A. Wilhite, and O. V. Vanyarkho, "Monitoring the 1996 drought using the standardized precipitation index," *Bulletin of the American Meteorological Society*, vol. 80, no. 3, pp. 429–438, 1999.

[29] Z. Sen, "Statistical analysis of hydrologic critical droughts.," *Journal of the Hydraulics Division*, vol. 106, no. 15134, pp. 99–115, 1980.

[30] G. E. P. Box and G. Jenkins, *Time Series Analysis: Forecasting and Control*, Holden-Day, San Francisco, Calif, USA, 1976.

[31] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, Springer Series in Statistics, Springer, New York, NY, USA, 2nd edition, 2009.

[32] L. Breiman, "Bagging predictors," *Machine Learning*, vol. 24, no. 2, pp. 123–140, 1996.

[33] L. Breiman, "Arcing classifiers," *Annals of Statistics*, vol. 26, no. 3, pp. 801–849, 1998.

[34] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.

[35] P. M. Granitto, F. Gasperi, F. Biasioli, E. Trainotti, and C. Furlanello, "Modern data mining tools in descriptive sensory analysis: a case study with a random forest approach," *Food Quality and Preference*, vol. 18, no. 4, pp. 681–689, 2007.

[36] C. Lehmann, T. Koenig, V. Jelic et al., "Application and comparison of classification algorithms for recognition of Alzheimer's disease in electrical brain activity (EEG)," *Journal of Neuroscience Methods*, vol. 161, no. 2, pp. 342–350, 2007.

[37] B. A. Goldstein, A. E. Hubbard, A. Cutler, and L. F. Barcellos, "An application of random forests to a genome-wide association dataset: methodological considerations and new findings," *BMC Genetics*, vol. 11, article 49, 2010.

[38] B. Larivière and D. van den Poel, "Predicting customer retention and profitability by using random forests and regression forests techniques," *Expert Systems with Applications*, vol. 29, no. 2, pp. 472–484, 2005.

[39] P. M. Kuhnert, A. K. Henderson, R. Bartley, and A. Herr, "Incorporating uncertainty in gully erosion calculations using the random forests modelling approach," *Environmetrics*, vol. 21, no. 5, pp. 493–509, 2010.

[40] L. Firth, M. L. Hazelton, and E. P. Campbell, "Predicting the onset of Australian winter rainfall by nonlinear classification," *Journal of Climate*, vol. 18, no. 6, pp. 772–781, 2005.