



Hierarchical Bayesian Network Based Incremental Model for Flood Prediction

Yirui Wu^{1,2}, Weigang Xu¹, Qinghan Yu¹, Jun Feng^{1(✉)}, and Tong Lu²

¹ College of Computer and Information, Hohai University, Nanjing, China
{wuyirui,fengjun}@hhu.edu.cn, weigangxu@163.com, 1140833939@qq.com

² National Key Lab for Novel Software Technology, Nanjing University, Nanjing, China
lutong@nju.edu.cn

Abstract. To minimize the negative impacts brought by floods, researchers pay special attention to the problem of flood prediction. In this paper, we propose a hierarchical Bayesian network based incremental model to predict floods for small rivers. The proposed model not only appropriately embeds hydrology expert knowledge with Bayesian network for high rationality and robustness, but also designs an incremental learning scheme to improve the self-improving and adaptive ability of the proposed model. Following the idea of a famous hydrology model, i.e., XAJ model, we firstly present the construction of hierarchical Bayesian network as local and global network construction. After that, we propose an incremental learning scheme, which selects proper incremental data to improve the completeness of prior knowledge and updates parameters of Bayesian network to prevent training from scratch. We demonstrate the accuracy and effectiveness of the proposed model by conducting experiments on a collected dataset with one comparative method.

Keywords: Incremental learning · Hierarchical Bayesian network
Flood prediction

1 Introduction

Flood, as one of the most common and largely distributed natural disasters, happens occasionally and brings large damages to life and property. In the past decades, researchers have proposed a quantity of models for accurate, robust and reasonable flood prediction. We generally category models into two types, namely hydrology model [8, 11, 17] and data-driven model [4, 6, 18]. Hydrology models utilize highly non-linear mathematic systems to represent the complex hydrology processes from clues to results. However, such models are extremely sensitive to parameters [16] and require quantity of research efforts of experts to fit them for one specific river. On the contrary, data-driven models use machine learning methods to directly predict the river runoff values based on historical observed and time-varying flood factors. However, floods are complicated natural

phenomena affected by multiple factors. It's hard to guarantee the rationality and robustness by utilizing such data-driven models and not considering physical processes.

In this work, we pay special interests to the problem of flood prediction for small rivers, whose catchments are smaller than 3000 km. Predicting floods with either hydrology models or data-driven models for small rivers could be a hard task, since small rivers are not only complex to model and analyze, but also suffer from shortages of exhaustive historical observation data. It's an intuitive thought that we should properly utilize the strength of hydrology model to improve the accuracy, robustness and rationality of data-driven model. The hydrology expert knowledge behind the hydrology model could relieve the requirement for large amount of data, which solves the problem of not enough data at a certain extent. Moreover, we aim to construct data-driven models with "growth" ability. That is the predicting capability of models could be gradually improved with more captured data. In fact, the floods data collected in small rivers are generally lack of completeness and unevenly distributed. By involving the ability of growth, the constructed model can run ahead and converge to a finalized and robust system during the running period. Moreover, the predicting capability of models are greatly affected by the occurrence of climatic variations, human activities and other environmental changes. Models with growth ability thus should continuously process new information captured from the latest floods and make self-adaptive adjustments to ensure the accuracy of predictions.

Guided by the ideas of expertise and growth, we propose a hierarchical Bayesian network based incremental model. In order to extract the expert hydrology knowledge behind physical models, the entities and relations of the proposed model refer to the physical factors and processes extracted from a famous hydrology model, *i.e.*, the XAJ model [10, 17]. Moreover, we construct an incremental learning scheme to develop the growth ability of the proposed model without changing network structures or training from the scratch.

The main contribution of the paper is to propose a hierarchical Bayesian network based incremental model for flood prediction of small rivers, which not only embeds hydrology process to improve the accuracy, robustness and rationality, but also designs an incremental learning scheme to improve the self-improving and adaptive ability. Owing to the expertise and growth ability of the proposed model, the requirements for size of training dataset could be largely reduced, which coincides with the environment and conditions of predicting floods for small rivers. The proposed method is powerful to discover the inherent patterns between input flood factors and flow rate, especially for regions whose flood formation mechanism is too complex to construct a convinced physical model.

2 Related Work

Hydrology Model. The famous XAJ model not only considers the rains and runoffs, but also takes other hydrology processes into account, such as evapora-

tion from water bodies and surface, rain infiltrated and stored by the soil, and so on. We explain the processes of XAJ model with the following four modules:

1. Evaporation module: XAJ model firstly divide the river watershed into several local regions. Evaporation values of local regions are computed based on the soil tension water capability (referring to soil water storage capability) in three layers, *i.e.*, upper, lower and deep soil layers.
2. Runoff generation module: The XAJ model defines local runoff is not produced until the soil water of the local region reaches its maximum of soil tension water capacity, and thereafter the excess rainfall becomes the runoff without further loss. Therefore, the local runoff of XAJ model is calculated according to the rainfall, evaporation and soil tension water capability.
3. Runoff separation module: The local runoff is subdivided into three components, including surface runoff, interflow runoff and groundwater runoff.
4. Runoff routing module: The outflow from each local region is finally routed by the Muskingum successive-reaches model [17] to calculate the outlet flow of the whole river catchment.

Sensitive parameters of the XAJ model need be adjusted by experts' experiences, which makes it difficult to apply on small rivers for predictions.

Data-Driven Model. From the views of computer scientists, floods are directly induced and affected by a set of multiple factors, including rainfall, soil category, the structure of riverway and so on. Early, Reggiani *et al.* [9] construct a modified Bayesian predicting system by involving numerical weather information to address the spatial-temporal variabilities of precipitation during prediction. Later, Cheng *et al.* [1] perform accurate daily runoff forecasting by proposing an artificial neural network based on quantum-behaved particle swarm optimization, which trains the ANN parameters in an alternative way and achieves much better forecast accuracy than the basic ANN model. Recently, Wu *et al.* [14] construct a Bayesian network for flood predictions, which appropriately embeds hydrology expert knowledge for high rationality and robustness. The proposed method is built on it and involves an incremental design over all steps of Bayesian network for fitting to the problem of flood predictions for small rivers.

Impressed by significant ability of deep learning architectures [5, 7, 15], researchers try to utilize deep learning architectures for flood prediction. For example, Zhuang *et al.* [18] design a novel Spatio-Temporal Convolutional Neural Network (ST-CNN) to fully utilize the spatial and temporal information and automatically learn underlying patterns from data for extreme flood cluster prediction. Liu *et al.* [6] propose a deep learning approach by integrating stacked auto-encoders (SAE) and back propagation neural networks (BPNN) for the predictions of stream flow. Most recently, Wu *et al.* [13] propose context-aware attention LSTM network to accurately predict sequential flow rate values based on a set of collected flood factors. However, the above deep learning methods require large datasets to train. Without prior knowledge and inferences extracted from hydrology models, the deep learning based models can't predict floods in a rational sense.

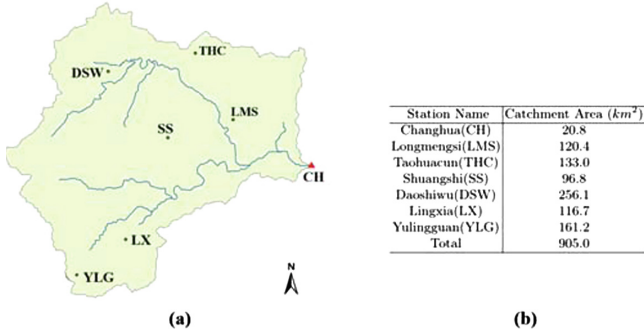


Fig. 1. Illustration of the Changhua watershed, where (a) is the map for various kinds of stations and (b) represents catchment areas corresponding to the listed rainfall stations. Note that we need predict the flow rate values of river gauging station CH and station SS functions as an evaporation station.

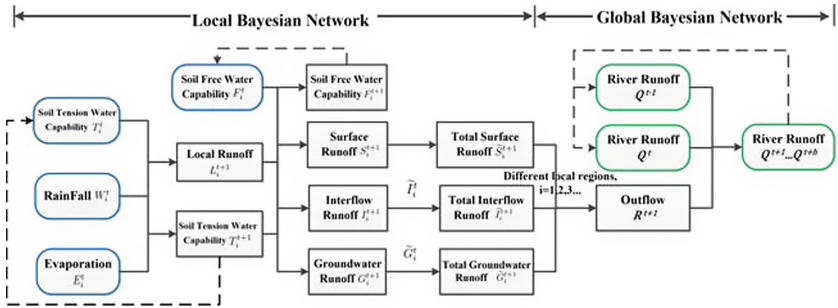


Fig. 2. Illustration of the proposed hierarchical Bayesian network based incremental model, where dotted lines refer to time-varying updating, blue and green rectangles represent incremental inputs and flood predictions, respectively. (Color figure online)

3 The Proposed Method

Take a typical small river, *i.e.*, Changhua, for an example, we show its general information in Fig. 1, where we can notice 7 rainfall stations, 1 evaporation station and 1 river gauging station. In our work, we aim to predict the flow rate values at the river gauging station CH for the next 6 h with the proposed incremental model. The input set of flood factors consists of rainfalls observed at the rainfall stations, evaporation and soil moisture observed at the evaporation station SS and former river runoff observed at CH.

Considering that XAJ model is organized with local and global steps, we follow its conception to design the proposed hierarchical Bayesian network based incremental model as shown in Fig. 2. By inferring probabilistic relations between inputting flood factors and intermediate variables extracted from the XAJ model, we embed the hydrological expert knowledge with the proposed model by first establishing relations and then improving the representations of knowledge with

probabilistic distributions other than function systems. We construct the incremental learning scheme by firstly selecting proper data to improve the generative completeness of the proposed Bayesian network and then updating Conditional Probability Table (CPT) of network, which prevents training from scratch. Note that we calculate the initial value of soil tension water capability T_i^t based on soil moisture measured at the evaporation station. Meanwhile, soil free water capability F_i^t is settled as 0 at the beginning, which will gradually converge to the real value. Note we transform the run-off regression problem to a multi-label classification problem by splitting the observed runoff values of Changhua dataset into 2000 intervals, *i.e.*, assigning 2000 labels to the predictions of run-off values.

3.1 Construction of Hierarchical Bayesian Network

In this subsection, we firstly introduce the theory foundation and novelty by utilizing Bayesian Network for flood predicting. After that, we describe the construction of Hierarchical Bayesian network.

Given data D , we determine the posterior distribution of θ based on Bayesian theory as follows:

$$P(\theta|D) = \frac{L(D|\theta)P(\theta)}{P(D)} \quad (1)$$

where $L(D|\theta)$ is the likelihood function and $P(\theta)$ is the prior distribution of random variable θ . Since the denominator of Eq. 1 is a constant related only to the data set, the choice of prior distribution $P(\theta)$ is important for calculation of the posterior distribution $P(\theta|D)$. Selecting proper $P(\theta)$ generally requires to consider from the measured data and available prior knowledge. The former is named as data-based prior distribution and could be obtained from the existing data and research results, while the latter, named as non-data-based prior distribution, refers to a prior distribution resulted from subjective judgments or theory.

By extracting prior expert hydrology knowledge from the XAJ model and historic observation data, we think Bayesian Network offers an appropriate structure to joint learn the posterior distribution with the prior knowledge. Specifically, the proposed method firstly considers the given observation data D is formed by a set of hydrology attributes $\{X_i|i = 1..n\}$ and the predicting run-off value could be represented as an attribute X_0 as well. Therefore, we could represent the joint distribution of $\{X_i|i = 0..n\}$ as

$$P(X_0, X_1, X_2, \dots, X_n) = \prod_{i=0}^n P(X_i|\zeta(Parents(X_i))) \quad (2)$$

where function $Parents()$ and $\zeta()$ represents the sets of directly precursor attributes and the corresponding joint distribution, respectively. In order to solve Eq. 4 for X_0 , we utilize marginalization [3] operations to convert it as a list of

conditional probabilities. We further adopt Bayesian network and the cooperating CPTs to describe conditional probabilities. During training, we use loopy belief propagation to estimate the parameters of conditional probability table. Due to the loopy structure of the network, it is difficult to check for the convergence. We thus adopt that training is terminated when 10 iterations of gradient decent go not yield averagely improved likelihood over the previous 10.

After explaining the theory of Bayesian network, we describe the construction of hierarchical Bayesian network. During the Local Bayesian Network stage, we aim to predict the runoff contribution values in the local regions. We firstly divide the total river watershed into small local regions based on hydrology principles [12] and the locations of rainfall stations. The split results of local regions are represented in Fig. 1(b). We then collect multiple kinds of inputs in each local region, *i.e.* soil moisture T_i^t , rainfall W_i^t and evaporation E_i^t by interpolation based on observed flood factors, where i refers to the index of local region. Next, we follow the first three modules of the XAJ model as discussed in the last section, in order to embed the expert knowledge about hydrology processes into the construction of the local Bayesian network. Finally, the trained local Bayesian network could compute several hydrology intermediate variables, such as surface runoff \tilde{S}_i^{t+1} , interflow runoff \tilde{I}_i^{t+1} and groundwater runoff \tilde{G}_i^{t+1} . In the Global Bayesian Network stage, we utilize the last module of XAJ model to construct the global Bayesian network, which predicts the river runoff for the next h hours $\{Q^t, \dots, Q^{t+h}\}$ based on the output of the local Bayesian network and river runoff Q^{t-1}, Q^t in former times. To sum up, we properly embed the hydrology process and variables of the XAJ model into the hierarchical Bayesian Network.

3.2 Bayesian Network Incremental Learning

In this subsection, we firstly discuss how to select proper incremental data to improve the completeness of the proposed model and then describe steps to update CPTs of the proposed hierarchical Bayesian network.

Incremental data selection is one of the most important factors to improve efficiency of incremental learning. In fact, selecting false labeled samples will bring noise and decrease accuracy of further predictions. Generally, researchers

Algorithm 1. Incremental sample selection algorithm

- Input:** Model trained in the last iteration M , set of incremental samples S
Output: Prior incremental set $P = \emptyset$, undetermined incremental set $U = \emptyset$ and noise set $N = \emptyset$
- 1: **For** each $a_i \in S, c = gt(a_i)$
 - 2: **If** $c \in C_n, N.add(a_i)$
 - 3: **Else** $\beta = M(a_i)$
 - 4: **If** $|\beta - c| < \omega, P.add(a_i)$
 - 5: **ELSEIf** $|\beta - c| < \varepsilon, U.add(a_i)$
 - 6: **ELSE** $N.add(a_i)$
-

select incremental data by calculating model loss, defined as difference values of prediction accuracy between before and after selecting new samples for incremental learning. However, such procedure is rather low in efficiency due to time-consuming calculation. We thus propose a threshold-ruled incremental data selection algorithm for better efficiency, which is presented in Algorithm 1.

In Algorithm 1, function $gt()$ checks the ground-truth classification label from training dataset, function $add()$ adds an incremental sample into different sets, function $M()$ refers to the classification result achieved by hierarchical Bayesian Network in the last iteration, C_n represents the classification labels set in the last iteration, ω and ε are two adaptive parameters to decide the operation on the inputting incremental sample. Specifically, we define $\omega = \tilde{Q} \times 5\%$ and $\varepsilon = \tilde{Q} \times 20\%$ to avoid the induce of noise data, where \tilde{Q} refers to the mean runoff value corresponding to the small river. Note that 20% is originated from the international rule for permissible range of flood prediction system error. After defining the set of P and U based on the inputting data S , we add the samples of P for incremental training at first. After then, we utilize a matrix generated from the normal distribution to expand the data in P by $\tilde{p} = L * p$. For the generated and expanded data \tilde{p} , we further process it as input by Algorithm 1 and utilize the corresponding results of P and U for incremental training at last.

After selection on the proper incremental data, we discuss the updating rule inside the network. When incremental data and the former training date are ruled by the same joint distribution, the training Bayesian network could be adjusted only with the parameters to fit with new data. Following this idea, we define D_0 , D_+ and $D = D_0 + D_+$ as the initial dataset, incremental dataset and total dataset, respectively. We also define the number of dataset as $N_0 = |D_0|$, $N_+ = |D_+|$ and $N = N_0 + N_+$. Supposing that there are n variables X_1, X_2, \dots, X_n and the corresponding possible values $x_i^1, x_i^2, \dots, x_i^{r_i}$, we could use

$$\theta_{ijk} = p(x_i^k | \pi_i^j, \theta_i, G) \tag{3}$$

to represent the parameters of Bayesian network with structure G , where $\pi_i^1, \pi_i^2, \dots, \pi_i^{q_i}$ ($q_i = \prod_{x_m \in \pi_i} r_j, m \neq i$) are the father node set for node X_i . After adding samples for incremental learning, we thus could calculate the modified parameters as

$$\theta_{ijk}(D, G) = \frac{\theta'_{ijk}(D_0, G) + N_{ijk}(D_+, G)}{\theta'_{ij}(D_0, G) + N_{ij}(D_+, G)} \tag{4}$$

where $\theta'_{ij}(D_0, G) = \sum_{k=1}^{r_i} \theta'_{ijk}(D_0, G)$, $N_{ij}(D_+, G) = \sum_{k=1}^{r_i} N_{ijk}(D_+, G)$ and the network parameters can be defined as

$$\begin{cases} \sum_{k=1}^n \theta_{ijk} = 1 \\ \theta_{ij} = \bigcup_{k=1}^{r_i} \theta_{ijk} \\ \theta_i = \bigcup_{j=1}^{q_i} \theta_{ij} \\ \theta = \bigcup_{i=1}^n \theta_i \end{cases} \tag{5}$$

4 Experimental Results

4.1 Dataset and Measurements

We collect hourly data of floods happened from 1998 to 2010 in Changhua river as our dataset. The floods happened from 1998 to 2003 and from 2009 to 2010 are used as the basic training and testing dataset respectively, meanwhile floods happened from 2004 to 2008 are adopted as the incremental datasets, which are divided into five parts and marked with D_1 to D_5 , respectively. We analysis the runoff values of Changhua dataset and find the values are unevenly distributed in a fixed interval, which proves the supposition for data of small rivers, *i.e.*, incomplete and highly uneven. Therefore, it's necessary to involve the incremental learning to improve the performance of flood prediction in small rivers.

To better evaluate performance of the proposed method, we adopt several quality measurements for evaluation of classification results, which could be represented as

$$FN = \frac{N_{non}}{N_{all}} \quad (6)$$

$$k - FC = \frac{N_{k,correct}}{N_{all}} \quad (7)$$

where N_{all} is the total number of testing samples and N_{non} refers to number of none deciding testing samples, which can't be assigned with labels by the proposed model due to the lack of complete prior knowledge, *i.e.*, related probability inferences. $N_{k,correct}$ refers to the number of testing samples, whose run-off prediction values are close with ground-truth values. The difference value between the prediction and ground-truth should be smaller than value represented by k splitting intervals, where k is define as 1 in our experiment. Note that FN is designed to show the ability to acquire new knowledge during the process of incremental learning, meanwhile k-FC is used to evaluate the ability for accurately flood prediction. Higher FN and k-FC value implies better performance.

4.2 Performance Analysis

We show the improvements on FN and 1-FC measurement with the proposed method in Fig. 3. We can observe great decrements of FN values during the period of incremental learning, especially for the first, third and fifth increment. This is due to the completeness of the prior knowledge is gradually increased with more training samples and the proposed method is efficient in extracting such knowledge by incremental learning. The reason for different decrement values lies in the fact that the dataset is split based on year other than the amount of new knowledge. For 1-FC, we can view an obvious decrement in prediction accuracy with larger perdition hours, which implies the task of flood prediction becomes harder when predicting for a relatively long time. With the incremental learning, we find the prediction accuracy is improved, especially for the first and fifth increment. The most obvious improvements are labeled by blue rectangles in Fig. 3, which refer to the fifth incremental learning for prediction in 4 and 5 h.

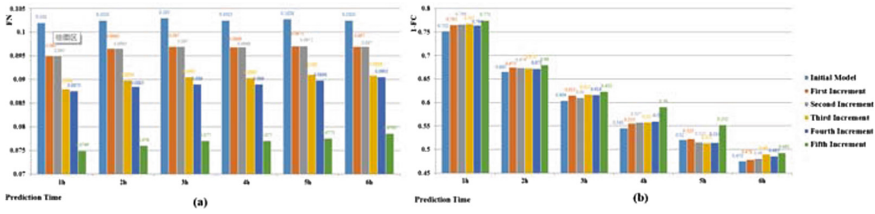


Fig. 3. Illustration of the improvements on FN and 1-FC with the proposed incremental learning scheme, where blue rectangles represent the obvious improvement on 1-FC with the fifth increment. (Color figure online)

This fact proves that the proposed method is better at predicting in a relatively long time.

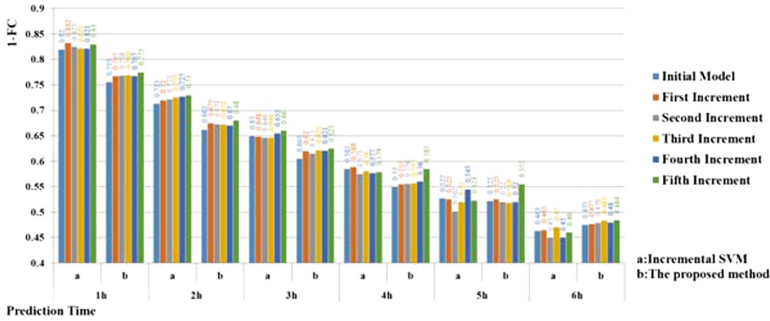


Fig. 4. Comparison of 1-FC values on Changhua dataset computed by the proposed method and incremental SVM.

In Fig. 4, we compare the 1-FC values computed by the proposed method and incremental SVM [2]. Since SVM could predict without complete prior knowledge, it’s meaningless to compare FN. We implement the incremental SVM according to the instructions given in their paper. From Fig. 4, we can find the prediction accuracy achieved by the proposed method is lower than that achieved by the incremental SVM when predicting for 1 h, 2 h, 3 h and 4 h. However, the proposed method gets better performance when predicting for 5 h and 6 h, which proves the proposed method is better than incremental SVM at predicting in a relatively long time. With Incremental learning, we can find improvements achieved by either incremental SVM or the proposed method. However, the increase values gained by the proposed method are more impressive than that gained by the incremental SVM, especially when predicting for 4 h and 5 h. This proves the proposed method is more efficient than incremental SVM for tasks of incremental learning, especially for long time flood predicting.

5 Conclusion

In this paper, we propose a hierarchical Bayesian network based incremental model to predict floods for small rivers. The proposed model not only appropriately embeds hydrology expert knowledge with Bayesian network for high rationality and robustness, but also designs an incremental learning scheme to improve the self-improving and adaptive ability of the proposed model. By involving power of incremental learning, the proposed model could be gradually improved with more collected data, which makes it fit with various application scenarios. Experiment results on Changhua dataset show the proposed method outperforms several comparative methods and achieves promising prediction results on small rivers. Our future work includes the exploration on other hydrology purposes with the proposed method, for example mid-term flood predicting.

Acknowledgement. This work was supported by National Key R&D Program of China under Grant 2018YFC0407901, the Natural Science Foundation of China under Grant 61702160, Grant 61672273 and Grant 61832008, the Fundamental Re-search Funds for the Central Universities under Grant 2016B14114, the Science Foundation of Jiangsu under Grant BK20170892, the Science Foundation for Distinguished Young Scholars of Jiangsu under Grant BK20160021, Scientific Foundation of State Grid Corporation of China (Research on Ice-wind Disaster Feature Recognition and Prediction by Few-shot Machine Learning in Transmission Lines), and the open Project of the National Key Lab for Novel Software Technology in NJU under Grant K-FKT2017B05.

References

1. Cheng, C., Niu, W., Feng, Z., Shen, J., Chau, K.: Daily reservoir runoff forecasting method using artificial neural network based on quantum-behaved particle swarm optimization. *Water* **7**(8), 4232–4246 (2015)
2. Diehl, C.P., Cauwenberghs, G.: SVM incremental learning, adaptation and optimization. In: *Proceedings of International Joint Conference on Neural Networks*, vol. 4, pp. 2685–2690 (2003)
3. Friedman, N., Geiger, D., Goldszmidt, M.: Bayesian network classifiers. *Mach. Learn.* **29**(2–3), 131–163 (1997)
4. Han, S., Coulibaly, P.: Bayesian flood forecasting methods: a review. *J. Hydrol.* **551**, 340–351 (2017)
5. Jing, P., Su, Y., Nie, L., Bai, X., Liu, J., Wang, M.: Low-rank multi-view embedding learning for micro-video popularity prediction. *IEEE Trans. Knowl. Data Eng.* **30**(8), 1519–1532 (2018)
6. Liu, F., Xu, F., Yang, S.: A flood forecasting model based on deep learning algorithm via integrating stacked autoencoders with BP neural network. In: *Proceedings of IEEE International Conference on Multimedia Big Data*, pp. 58–61 (2017)
7. Nie, L., Zhang, L., Yan, Y., Chang, X., Liu, M., Shaoling, L.: Multiview physician-specific attributes fusion for health seeking. *IEEE Trans. Cybern.* **47**(11), 3680–3691 (2017)
8. Paquet, E., Garavaglia, F., Garçon, R., Gailhard, J.: The schadex method: a semi-continuous rainfall-runoff simulation for extreme flood estimation. *J. Hydrol.* **495**, 23–37 (2013)

9. Reggiani, P., Weerts, A.: Probabilistic quantitative precipitation forecast for flood prediction: an application. *J. Hydrometeorol.* **9**(1), 76–95 (2008)
10. Ren-Jun, Z.: The Xinanjiang model applied in China. *J. Hydrol.* **135**(1–4), 371–381 (1992)
11. Rogger, M., Viglione, A., Derx, J., Blöschl, G.: Quantifying effects of catchments storage thresholds on step changes in the flood frequency curve. *Water Resour. Res.* **49**(10), 6946–6958 (2013)
12. Villarini, G., Mandapaka, P.V., Krajewski, W.F., Moore, R.J.: Rainfall and sampling uncertainties: a rain gauge perspective. *J. Geophys. Res. Atmos.* **113**(D11) (2008)
13. Wu, Y., Liu, Z., Xu, W., Feng, J., Shivakumara, P., Lu, T.: Context-aware attention LSTM network for flood prediction. In: *Proceedings of International Conference on Pattern Recognitions* (2018)
14. Wu, Y., Xu, W., Feng, J., Shivakumara, P., Lu, T.: Local and global Bayesian network based model for flood prediction. In: *Proceedings of International Conference on Pattern Recognition* (2018)
15. Wu, Y., Yue, Y., Tan, X., Wang, W., Lu, T.: End-to-end chromosome Karyotyping with data augmentation using GAN. In: *Proceedings on International Conference on Image Processing*, pp. 2456–2460 (2018)
16. Yao, C., Zhang, K., Yu, Z., Li, Z., Li, Q.: Improving the flood prediction capability of the Xinanjiang model in ungauged nested catchments by coupling it with the geomorphologic instantaneous unit hydrograph. *J. Hydrol.* **517**, 1035–1048 (2014)
17. Zhao, R., Zhuang, Y., Fang, L., Liu, X., Zhang, Q.: The Xinanjiang model. In: *Proceedings Oxford Symposium Hydrological Forecasting*, vol. 129, pp. 351–356 (1980)
18. Zhuang, W.Y., Ding, W.: Long-lead prediction of extreme precipitation cluster via a spatiotemporal convolutional neural network. In: *Proceedings of the 6th International Workshop on Climate Informatics: CI* (2016)