# Predicting monthly streamflow using data-driven models coupled with data-preprocessing techniques

C. L. Wu, K. W. Chau* and Y. S. Li

Dept. of Civil and Structural Engineering, Hong Kong Polytechnic University,

Hung Hom, Kowloon, Hong Kong, People's Republic of China


*Email: cekwchau@polyu.edu.hk

## ABSTRACT

In this paper, the accuracy performance of monthly streamflow forecast was discussed when employing data-driven modeling techniques on the streamflow series. A Crisp Distributed Support Vectors Regression (CDSVR) model was proposed for monthly streamflow prediction in comparison with four other models, Auto-Regressive Moving Average (ARMA), K-Nearest-Neighbors (KNN), Artificial Neural Networks (ANNs), and Crisp Distributed Artificial Neural Networks (CDANN). With respect to distributed models of CDSVR and CDANN, the Fuzzy-C-Means (FCM) clustering technique first split the flow data into three subsets (low, medium and high levels) according to the magnitudes of the data, and then three single SVRs (or ANNs) were fitted to three subsets. This paper gave a detailed analysis on reconstruction of dynamics which was used to identify the configuration of all models except for ARMA. To improve the model performance, the data preprocessing techniques of Singular Spectrum Analysis (SSA) and/or Moving Average (MA) were coupled with all five models. Some discussions were presented (1) on the number of neighbors in KNN; (2) on the configuration of ANN; and (3) on the investigation of effects of MA and SSA. Two streamflow series from different locations of People's Republic of China (Xiangjiaba and Danjiangkou) were applied for the analysis of forecasting. Forecasts were conducted at four different horizons (one-, three-, six-, and twelve-month-ahead forecasts). The results showed that models fed by pre-processed data performed better than models fed by original data, and CDSVR outperformed other models except for at six-month-ahead horizon for Danjiangkou. For the perspective of streamflow series, the SSA exhibited better effects on Danjingkou data because its raw discharge series was more complex than the discharge of Xiangjiaba. The MA considerably improved the performance of ANN, CDANN, and CDSVR by adjusting the

1   correlation relationship between input components and output of models. It can be also found that the performance

2   of CDSVR deteriorated with the increase of the forecast horizon.

3

4   **KEYWORDS**

5   Monthly streamflow forecast, Distributed Support Vector Regression, Reconstruction of dynamics, Singular

6   Spectrum Analysis, False Nearest Neighbors, Moving average, Fuzzy-C-Means

7

8   ## 1. Introduction

9       As an important issue in hydrology time series, streamflow (or discharge) series

10   prediction was of great concern over the past few decades. Numerous data-driven modeling

11   techniques were proposed for the forecast and simulation of the streamflow series (*Carlson et*

12   *al., 1970; Jain et al. 1999; cannas, et al., 2006; Wang et al, 2006b; Lin et al., 2006*). The linear

13   autoregressive models were based on the assumption that hydrological time series is originated

14   from a stochastic process with an infinite number of degrees of freedom. Linear models such as

15   AutoRegressive (AR), AutoRegressive Moving Average (ARMA), AutoRegressive Integrated

16   Moving Average (ARIMA), and Seasonal ARIMA (SARIMA) had made a great success in

17   river flow prediction (*Carlson et al., 1970; Salas et al., 1985; Haltiner and Salas, 1988; Yu and*

18   *Tseng, 1996; Kothyari and Singh, 1999; Huang et al., 2004; María et al., 2004*). Meanwhile,

19   chaos-based streamflow predictions have gained increasing interests of the hydrology

20   community in the past two decades (*Jayawardena and Lai, 1994; Jayawardena and Gurung,*

21   *2000; Elshorbagy et al., 2002; Sivakumar et al.,2002*) in which a streamflow series was

22   assumed to be derived from a nonlinear deterministic dynamic system. The local prediction

23   technique of K-Nearest-Neighbors (KNN) algorithm was generally used (*Farmer and*

24   *Sidorowich, 1987*). However, some doubts have been raised in terms of the existence of chaos

1     in hydrologic data (*Ghilardi and Rosso, 1990; Koutsoyiannis and Pachakis, 1996; Pasternack,*

2     *1999; Schertzer et al., 2002; Wang et al., 2006a*).

3         Both the above two assumptions tends to be absolute descriptions for a hydrologic

4     streamflow series. Salas et al. (*1985*) suggested that a streamflow process should be treated as

5     an integration of stochastic (or random) and deterministic components. Describing it as either a

6     totally linear stochastic process or fully nonlinear deterministic chaos was not a practical

7     approach (*Elshorbagy et al. 2002*). The model based on either of two assumptions may not be

8     the most suitable. Therefore, various soft computing methods such as Artificial Neural

9     Networks (ANNs) and Support Vector Regression (SVR) have been widely employed to

10    capture the nonlinearity in a streamflow series. For example, comparisons between ANNs and

11    AR approaches appeared in the work of Raman and Sunilkumar (*1995*), Thirumalaiah and Deo

12    (*2000*), and Kişi (*2003; 2005*). Comparisons between ANNs and ARMA were performed in

13    Jain et al. (*1999*), Abrahart and See (*2002*), and Castellano-Me ́ndeza et al. (*2004*). Most of

14    these studies proved that ANNs outperform traditional statistical techniques. Lin et al. (*2006*)

15    proposed to use a SVR model to predict long-term monthly flow discharge series. It was

16    demonstrated that SVR exhibited a better performance than ARMA. Comparison studies

17    between SVR and ANN were also conducted in literature. Dibike et al. (*2001*) and Behzad et al.

18    (*2009*) drew similar conclusions that SVR outperformed ANN in modeling rainfall-runoff

19    process. Laio et al. (*2003*) undertook a comparison of KNN and ANN for flood predictions and

20    found that KNN performed slightly better at short forecast period while the situation was

21    reversed for longer duration. Yu et al. (*2004*) proposed a scheme that combined chaos theory

22    and SVR to predict the daily discharge. Their results showed that KNN performed worse in

23    comparison with ARIMA and SVR. However, Sivakumar et al. (*2002*) found that the

1    performance of the Phase-Space Reconstruction (PSR) approach was consistently better than

2    that of ANN in short-term river flow prediction. Obviously, it is difficult to justify which

3    modeling technique is more suitable for a streamflow forecast from the results in the literature.

4         It is generally accepted that the underlying mechanisms of streamflow generation are

5    likely to be quite different during low, medium and high flow periods. For instance, the base

6    flow mainly contributes to low-flow events whereas intense storm rainfall gives rise to high-

7    flow events. A single global ANN model could not predict the high- and low-runoff events

8    satisfactorily (*Minns and Hall, 1996*). Modular models were therefore proposed where sub-

9    processes are first of all identified and then separate models (also called local or expert model)

10   are established for each of them (*Solomatine, and Ostfeld, 2008*). Depending on the split of

11   training data being soft or crisp, modular models are different. The soft split means the dataset

12   can be overlapped and the overall forecasting output is the weighted-average of each local

13   model (*Zhang and Govindaraju, 2000; Shrestha and Solomatine, 2006; Wang et al.,2006b; Wu*

14   *et al., 2008*). Zhang and Govindaraju (*2000*) examined the performance of modular networks in

15   predicting monthly discharges based on the Bayesian concept. Wu et al. (*2008*) employed a

16   distributed SVR (DSVR) for daily river stage prediction. The DSVR was a modified version of

17   the DSVR designed by Cheng et al. (*2006*) in which the final output is the weighted sum of the

18   outputs of each single SVR (hereafter referred to as Soft DSVR (SDSVR)). On the contrary,

19   there is not any overlap of data in the crisp split and the final forecasting output is based on one

20   of local models (*Corzo and Solomatine, 2007; Jain and Srinivasulu, 2006; See and Openshaw,*

21   *2000; Sivapragasam and Liong, 2005; Solomatine and Xue, 2004*). Solomatine and Xue (*2004*)

22   used M5 model trees and neural networks in a flood-forecasting problem. Sivapragasam and

23   Liong (*2005*) divided the flow range into three regions, and employed different SVR models to

1   predict daily flows in high, medium and low regions. In fact, ANN or SVR cannot extrapolate

2   beyond the range of the data used for training. Poor forecasts/predictions are expected when a

3   new input data is outside of the range of those used for training. This problem in the case of

4   SVR was indeed noticed by Wu et al. (*2008*). The DSVR derived from crisp data split is

5   therefore proposed (hereafter referred to as Crisp DSVR (CDSVR)) to overcome the poor

6   extrapolation in the SDSVR. Similarly, the improved DANN is called Crisp DANN (hereafter

7   referred to as CDANN).

8       Suitable data preprocessing can improve the performance of data-driven models.

9   Besides the conventional rescaling or standardization of training data used in ANN models

10  (*Dawson and Wilby, 2001*), preprocessing methods from the perspective of signal analysis are

11  also crucial because streamflow series may be also viewed as a quasi-periodic signal, which is

12  contaminated by various noises at different flow levels. An important signal decomposition

13  technique, singular spectrum analysis (SSA) was recently introduced to hydrology field by

14  some researchers (*Lisi et al., 1995; Sivapragasam et al., 2001; Marques et al., 2006*). An

15  analysis technique for the series may be particularly significant if it is able to expose important

16  characteristics of the time series in order to attain predictability (*Marques et al., 2006*). The

17  SSA method decomposes a time series into a number of components with simpler structures,

18  such as a slowly varying trend, oscillations and noise. SSA uses the basis functions

19  characterized by data-adaptive nature, which makes the approach suitable for the analysis of

20  some nonlinear dynamics (*Elsner and Tsonis, 1997*). In addition, the issue of lagged predictions

21  in the ANN model has been mentioned by some researchers (*Dawson and Wilby, 1999; Jian*

22  *and Srinivasulu, 2004; de Vos and Rientjes, 2005; Muttil and Chau, 2006*). One of the reasons

23  for lagged predictions is the use of previous observed streamflow data as ANN inputs (*de Vos*

1 *and Rientjes, 2005*). An effective solution is to obtain new inputs of models by moving average

2 over the original streamflow series.

3      The purpose of this study is to investigate the performance of five data-driven models of

4 ARMA, KNN, ANN, CDANN and CDSVR with data-preprocessing techniques of SSA and

5 MA on two real monthly streamflow series. This paper is organized in the following manner.

6 Section 2 presents four sets of streamflow data used in this study and describes the principles of

7 PSR and the identification of its parameters using the correlation integral approach and the

8 FNN approach. This section helps to identify all model inputs except for ARMA. Section 3

9 describes the SSA and the choice of its parameters. The modeling methods are presented in

10 Section 4. The implementation of the forecast models, including data preparation and selection

11 of model parameters, is discussed in Section 5. Forecast results and some discussions are

12 described in Section 6 and conclusions of this study are presented in Section 7.

## 13 **2. Reconstruction of Dynamics**

### 14 **2.1 Streamflow Data**

15      Monthly streamflow series of three watersheds and one river, i.e. Xiangjiaba, Manwan,

16 Danjiangkou, and Yangtze River, were analyzed for reconstruction of dynamics. Considering

17 similarity of flow characteristics amongst the former three (see relevant analysis later), only

18 Xiangjiaba and Danjiangkou were selected for forecasting in this study.

19      The largest watershed, Xiangjiaba, is at the upstream of Yangtze River with average

20 yearly discharge of 4538 $m^3$/s. Monthly streamflow series were taken from the hydrological

21 station near the Xiangjiaba Dam site located in Sichuan Province. The basin area contributed to

1    the streamflow series is around $45.88 \times 10^4$ km$^2$. The period of the data was from January 1940

2    to December 1997.

3        The medium watershed, Manwan, is located in the Lancang River which originates

4    from the Qinghai-Tibet Plateau. Monthly streamflow series were taken from the hydrological

5    station near the Manwan Dam site located in Sichuan Province. The catchment area controlled

6    by the station is $11.45 \times 10^4$ km$^2$, and the average yearly discharge is 1230 m$^3$/s based on a

7    statistic of 30-year data (January 1974 to December 2003).

8        The smallest watershed, Danjiangkou, lies at the upstream of Han River with average

9    yearly discharge of 1203 m$^3$/s. Monthly streamflow data came from the hydrology station at the

10    Danjiangkou Dam site which is located in Hubei Province. The catchment area at the dam site

11    is around $9.5 \times 10^4$ km$^2$. The data range was from January 1930 to December 1981.

12        The last streamflow series is Yangtze River, the largest river in China. The selected

13    monthly streamflow data were from the hydrology station of Cuntan located in the middle

14    stream of the river. The stream flow series spanned from January 1893 to December 2007.

15        Four monthly streamflow series are shown in Fig. 1. Monthly streamflow data in

16    Xiangjiaba, Manwan, and Cuntan are characterized by a smooth process whereas monthly

17    streamflow data in Danjiangkou exhibits more complex oscillations. The linear fits (dotted lines

18    in Fig. 1) verify the consistency of the streamflow series. All series exhibit good consistency

19    because the linear fits are closed to horizontal. Since there was no large-scale hydraulic works

20    such as dams built during the data collection period, the streamflow process is fairly pristine in

21    each case.

1    **2.2 Phase Space Reconstruction**

2    The most frequently used reconstruction method for a univariate or multivariate time

3    series is the delay-time method (*Takens, 1981; Farmer and Sidorowich,1987; Sauer et*

4    *al. ,1991*). A dynamic univariate time series $\{x_1, x_2, \cdots, x_N\}$, may be reconstructed into a series

5    of delay vectors of the type $\mathbf{Y}_t = \{x_t, x_{t+\tau}, x_{t+2\tau}, \cdots, x_{t+(m-1)\tau}\}$, $t = 1, 2, \cdots, N-(m-1)\tau$, where

6    $\mathbf{Y}_t \in \mathbf{R}^m$, $\tau$ is the delay time as a multiple of the sampling period and $m$ is the embedding

7    dimension. Under ideal conditions of time series of infinite length, all the reconstructions

8    would be analogous and topologically equivalent to the real system. Owing to the shortness of

9    real time series and the inevitable presence of dynamical noise, optimal reconstruction is

10   involved in the choice of $m$ and $\tau$ (*Laio et al., 2003*).

11   The temporal evolution of the dynamic system is given as a mapping $\mathbf{Y}(t) \mapsto \mathbf{Y}(t+T)$

12   (or $\mathbf{Y}_t \mapsto \mathbf{Y}_{t+T}$). The functional relationship between the current state $\mathbf{Y}(t)$ at time $t$ and the

13   predicted state $\mathbf{Y}^F(t+T)$ at time $t+T$ can be written as follows:

14   $$\mathbf{Y}^F(t+T) = f(\mathbf{Y}(t)) + e_t \tag{1}$$

15   where $e_t$ is a typical noise term. In the form of time series, it can be expressed

16   as    $[x_{t+T}^F, x_{t+T+\tau}^F, \cdots, x_{t+T+(m-2)\tau}^F, x_{t+T+(m-1)\tau}^F] = f([x_t, x_{t+\tau}, \cdots x_{t+(m-2)\tau}, x_{t+(m-1)\tau}]) + e_t$    .    Therefore,

17   predicting future trajectory by current trajectory becomes viable once the function $f(\bullet)$ is

18   determined. In practice, the expression is often defined as:

19   $$x_{t+T+(m-1)\tau}^F = f(\mathbf{Y}(t)) + e_t \tag{2}$$

20   where only the last component in $\mathbf{Y}^F(t+T)$ is indicated since normally the prediction of this last

21   component is of concern (*Laio et al., 2003*).

1 **2.3 Determination of Parameters $(\tau, m)$**

2     The correlation exponent method and False Nearest Neighbors (FNN) method were

3 employed to identify the parameter pair $(\tau, m)$ in this study.

4 **(1) Correlation exponent**

5     The method identifies the two parameters for the perspective of verifying the existence

6 of chaos. The diagnosis of the existence of chaos can begin if the phase space has been

7 reconstructed. The PSR requires two parameters $(\tau, m)$, where $m$ can be identified from the

8 plot of $d_2$ (the correlation dimension) versus $m$ when $d_2$ reaches a saturation value $D_2$ (*Tsonis*,

9 *1992*). In other words, $D_2$ must be obtained prior to the achievement of $m$. However, $m$ is

10 subject to $\tau$ (often called decorrelation time) and should be first determined. Also, some

11 researchers thought that $\tau$ and $m$ should not be determined separately. For example, Liong et al.

12 (*2002*) used Genetic Algorithm (GA) to optimize the triplet $(m, \tau, k)$. In the present study, we

13 tend to adopt a widely accepted method to obtain $\tau$ as follows.

14     The $\tau$ can be defined when the AutoCorrelation Function (ACF) attains the value of

15 zero or below a small value, or the Average Mutual Information (AMI) reaches the first

16 minimum (*Tsonis, 1992; Kantz and Schreiber, 2004*). The calculations of ACF and AMI are

17 discussed in detail in the works of Fraser and Swinney (*1986*), Tsonis (*1992*), and Abarbanel et

18 al. (*1993*). Fig. 2 displays the ACF and AMI for the four monthly streamflow series. The ACF

19 for each series was first attained zeros at the same lag time of 3. Since the AMI gave the same

20 estimates of $\tau$ as the ACF, $\tau$ was consistently chosen at the lag time of 3 for each case.

21     With the chosen $\tau$, the correlation dimension can be computed by the correlation

22 integral according the formula of Grassberger-Procaccia algorithm (*Grassberger and*

23 *Procaccia,1983*). This original formula was modified by Theiler (*1986*) for the estimation of

1    the correlation integral in a time series which poses serious problems of temporal correlations.

2    Thus, the modified correlation integral $C(r)$ for a $m$-dimension phase space is defined as:

3
$$C(r) = \frac{2}{N_{pairs}} \sum_{i=1}^{N} \sum_{j=i+w+1}^{N-i} H(r - \|\mathbf{Y}_i\text{-}\mathbf{Y}_j\|) \qquad (3)$$

4    where $N_{pairs} = (N-w+1)(N-w)$, $w$ is the Theiler window excluding those points which are

5    temporally correlated, $\mathbf{Y}_i$ and $N$ are the same as in Section 2.2, $r$ is the radius of a ball centered

6    on $\mathbf{Y}_i$, H is the Heaviside step function with $H(u)=1$ if $u>0$ and $H(u)=0$ if $u \le 0$. The

7    correlation integral only counts the pairs ($\mathbf{Y}_i, \mathbf{Y}_j$) whose distance, in a Euclidean sense, is

8    smaller than $r$. In the limit of an infinite amount of data ($N \to \infty$) and sufficiently small $r$, the

9    relation of $C(r) \propto r^{D_2}$ between $C(r)$ and $r$ is expected when $m$ exceeds the correlation

10   dimension of the chaos system. The correlation dimension $D_2$ and the correlation exponent $v$

11   can be defined as:

12
$$D_2 = \lim_{\substack{r \to 0 \\ N \to \infty}} v, \text{ where } v = \frac{\partial \ln C(r)}{\partial \ln r} \qquad (4)$$

13   Since $D_2$ is unknown before conducting the computation, the convergence of the correlation

14   dimension $D_2$ in $m$ must be examined.

15        The procedure of the computation is first to plot $\ln C(r)$ versus $\ln r$ with a given $m$.

16   Then, the potential scaling region is determined wherever the slope (i.e. the correlation

17   exponent $v$) of the curve for the given $m$ is approximately constant. The constant slope can be

18   estimated by a straight line fitting of the scaling region. In general, the best way to define the

19   scaling region is to produce another figure which demonstrates the slope of the $\ln C(r)$ as a

20   function of $\ln r$. If a scaling region exists, a plateau should be shown in the figure. This plateau

21   provides an estimate for $d_2$, a correlation dimension of the possible attractor for the present $m$.

1   If $d_2$ converges to a finite value $D_2$ (i.e. saturation value) after repeating the above procedure for

2   successively higher $m$, a true attractor of dimension $D_2$ is formed and the system may be

3   considered as chaos. Meanwhile, $m$ can be identified as the value that corresponds to the first

4   occurrence of the saturation value $D_2$ in the plot of $d_2$ versus $m$.

5       The graphs of the correlation exponent $v$ versus $\ln r$ for the four streamflow series are

6   shown in Figs. 3 respectively. Fig. 3 demonstrates that the scaling region cannot be identified

7   for any $m$ ($m < 20$) in Danjiangkou catchment whereas the scaling region can be determined in

8   other three cases. It can be noticed that the scaling region becomes ambiguous at $m = 8$ for

9   Manwan and Xiangjiaba whereas a narrow scaling region can still be defined even at $m = 20$ for

10  Cuntan.

11      However, an accurate estimation of $d_2$ requires a minimum number of points. Some

12  researches claim that the size should be $10^A$ (*Procaccia, 1988*) or $10^{2+0.4m}$ (*Tsonis, 1992*), where

13  $A$ is the greatest integer smaller than $d_2$ and $m$ ($m < 20$) is the embedding dimension used for

14  estimating $d_2$ with an error less than 5%. Other research found that smaller data size is needed.

15  For instance, the minimum data points for reliable $d_2$ is $10^{d_2/2}$ (*Ruelle, 1990; Essex and*

16  *Nerenberg, 1991*), or $\sqrt{27.5}^{d_2}$ (*Hong and Hong, 1994*) and empirical results of dimension

17  calculations are not substantially altered by going from 3000 or 6000 points to subsets of 500

18  points (*Abraham et al., 1986*).

19      Based on the results of Fig. 3, the relationship between correlation dimension $d_2$ and

20  embedding dimension $m$ is depicted in Fig. 4. The saturation values $D_2$ for three streamflow

21  series (Xiangjiaba, Cuntan and Manwan) are at interval of (1.5, 2). Generally, a sufficient

22  condition for the smallest $m$ is that $m$ is an integer larger than $2D_2$. The associated $m$ is

23  therefore set the value of 4 for the three series. With the potential values on $D_2$ or $m$, some

1    criteria such as $10^A$ can be satisfied whereas other criteria such as $10^{2+0.4m}$ cannot be satisfied.

2    The latter criteria means that few hydrologic records can be assessed for $m > 5$ attractors since

3    as many as 10,000 points require 27 years of daily records or around 900 years of monthly

4    records. Thus, the three monthly series of Xiangjiaba, Cuntan and Manwan may be treated as

5    chaotic with suggested variable $m = 4$. Danjiangkou series cannot be identified as non-chaotic

6    or random process because its data size is as small as 624. Furthermore, the phase portraits of

7    four streamflow series are portrayed in Fig. 5 where $(\tau, m)$ is (3,3). Obviously, the state spaces

8    in the 3-dimensional maps are clearly unfolded for Xiangjiaba, Cuntan, and Manwan whereas

9    no clear trajectory is revealed for Danjiangkou.

10    **(2) FNN**

11    The correlation integral method appears to be data intensive and certainly subjective.

12    For simplicity, the FNN method is commonly employed for the PSR of a hydrologic

13    streamflow series (*Wang et al., 2006b*). The FNN algorithm was originally developed for

14    determining the number of time-delay coordinates needed to recreate autonomous dynamics

15    directly from properties of the data itself (*Kennel et al., 1992; Abarbanel et al., 1993*). The

16    following discussion outlines the basic concepts of the FNN algorithm. Suppose the point

17    $\mathbf{Y}_i = \left\{ x_i, x_{i+\tau}, x_{i+2\tau}, \cdots, x_{i+(m-1)\tau} \right\}$ has a neighbor $\mathbf{Y}_j = \left\{ x_j, x_{j+\tau}, x_{j+2\tau}, \cdots, x_{j+(m-1)\tau} \right\}$, the criterion that

18    $\mathbf{Y}_j$ is viewed as a false neighbor of $\mathbf{Y}_i$ is:

19
$$\frac{\left| x_{i+m\tau} - x_{j+m\tau} \right|}{\left\| \mathbf{Y}_i \text{-} \mathbf{Y}_j \right\|} > R_{tol} \tag{5}$$

20    where $\| \ \|$ stands for the distance in a Euclidean sense, $R_{tol}$ is some threshold with the common

21    range of 10 to 30 (set as 15 in this study). For all points $i$ in the vector state space, Eq. (5) is

22    performed and then the percentage of points which have FNNs is calculated. The algorithm is

1    repeated for increasing $m$ until the percentage of FNNs drops to zero, or some acceptable small

2    number such as 1%, where $m$ is the target $m$.

3        Fig. 6 demonstrates the Percentage of FNNs (FNNP) for the four streamflow series with

4    $R_{tol} = 15$ and $\tau = 3$. The identified $m$ is 4 for Manwan and 5 for other three cases. The FNN

5    technique seems to be unable to distinguish random process from deterministic system because

6    a similar $m$ was found when the FNN was applied to a random series with the same data size as

7    Cuntan or Danjiakou. Thus, the phase space reconstructed by the FNN does not mean the

8    unfolding of a potential attractor in a dynamic system.

9        In summary, the parameter pair $(\tau, m)$ in PSR is (3, 4) for Xiangjiaba, Cuntan, and

10    Manwan and is (3, 5) for Danjiangkou when KNN is used for the streamflow prediction. The

11    determination of $m$ also means the achievement of the model inputs in three soft computing

12    models of ANN, CDANN, and CDSVR. For the purpose of forecasting, the value of the lagged

13    time $\tau$ is set to one in the three soft computing models because we tolerate having some

14    information redundancy in preference to losing any useful information. For instance, the

15    previous four month data is used as inputs of ANN if the one-month-ahead forecast is

16    performed for Xiangjiaba, and the previous five month data is the inputs of ANN for

17    Danjiangkou.

18    ## 3. Singular Spectrum Analysis

19    **3.1 Theory of SSA**

20        The SSA employed in the study can be referred to in Vautard et al. (*1992*) and Elsner

21    and Tsonis *(1997)*. For the sake of simplicity, a univariate time series is considered for

22    explanation of the SSA. In general, four steps are involved in the implementation of SSA. The

1     first step is to construct the 'trajectory matrix'. The 'trajectory matrix' results from the method

2     of delays. In the method of delays, the coordinates of the phase space will approximate the

3     dynamic of the system by using lagged copies of the time series. Therefore, the 'trajectory

4     matrix' can reflect the evolution of the time series with a careful choice of $(\tau, m)$ window. In

5     the context of SSA, the $m$ is usually called the window length (or singular number). For time

6     series $\{x_1, x_2, \cdots, x_N\}$, the 'trajectory matrix' is denoted by

7

$$
\mathbf{X} = \frac{1}{\sqrt{N}}
\begin{pmatrix}
x_1 & x_{1+\tau} & x_{1+2\tau} & \cdots & x_{1+(m-1)\tau} \\
x_2 & x_{2+\tau} & x_{2+2\tau} & \cdots & x_{2+(m-1)\tau} \\
x_3 & x_{3+\tau} & x_{3+2\tau} & \cdots & x_{3+(m-1)\tau} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
x_{N-(m-1)\tau} & x_{N-(m-2)\tau} & x_{N-(m-3)\tau} & \cdots & x_N
\end{pmatrix}
\tag{6}
$$

8     where $m$ is the embedding dimension, $\tau$ is the lagged (or delay) time. The matrix dimension is

9     $n \times m$ where $n = N - (m-1)\tau$. The next step is the singular value decomposition (SVD) of the

10     trajectory matrix $\mathbf{X}$. Let $\mathbf{S} = \mathbf{X}^{\mathrm{T}}\mathbf{X}$ (called the lagged-covariance matrix). With SVD, $\mathbf{X}$ can be

11     written as $\mathbf{X} = \mathbf{DLE}^{\mathrm{T}}$ where $\mathbf{D}$ and $\mathbf{E}$ are left and right singular vectors of $\mathbf{X}$, and $\mathbf{L}$ is a

12     diagonal matrix of singular values. $\mathbf{E}$ consists of orthonormal columns, and is also called the

13     'empirical orthonormal functions' (EOFs). Substituting $\mathbf{X}$ into the definition of $\mathbf{S}$ yields the

14     formula of $\mathbf{S} = \mathbf{EL}\mathbf{D}^{\mathrm{T}}\mathbf{DLE}^{\mathrm{T}} = \mathbf{EL}^2\mathbf{E}^{\mathrm{T}}$. Further $\mathbf{S} = \mathbf{E} \wedge \mathbf{E}^{\mathrm{T}}$ since $\mathbf{L}^2 = \wedge$ where $\wedge$ is a diagonal

15     matrix consisting of ordered values $0 \le \lambda_1 \le \lambda_2 \le \cdots \lambda_m$. Therefore, the right singular vectors of

16     $\mathbf{X}$ are the eigenvectors of $\mathbf{S}$ (*Elsner and Tsonis, 1997*). In other words, the singular vectors $\mathbf{E}$

17     and singular values of $\mathbf{X}$ can be respectively attained by calculating the eigenvectors and the

18     square roots of the eigenvalues of $\mathbf{S}$.

1      The first two steps make up the decomposition stage of SSA, and the next two steps

2      belong to the recovering stage. The third step is to calculate the principal components ($a_i^k$'s) by

3      projecting the original time record onto the eigenvectors as follows:

4
$$a_i^k = \sum_{j=1}^{m} x_{i+(j-1)\tau} e_j^k , \quad \text{for} \quad i=1,2\cdots, N-(m-1)\tau \tag{7}$$

5      where $e_j^k$ represents the jth component of the kth eigenvector. As known, each principal

6      component is a filtered process of the original series with length $N-(m-1)\tau$, not length $N$ as

7      desired, which poses a problem in real-time prediction.

8      The last step is to generate reconstruction components (RCs) whose lengths are the

9      same as the original series. The generation of each RC depends on a convolution of one

10     principal component with the corresponding singular vector, given by

11
$$x^k_{i+(j-1)\tau} = a_i^k e_j^k \quad \text{for } k=1,2\cdots, m \tag{8}$$

12     where $i=1,2\cdots, N-(m-1)\tau$, and $j=1,2\cdots, m$, or given by Vautard et al. (*1992*). Therefore,

13     The $m$ RCs can be achieved and their sum is a complete recover of the original series if all $m$

14     principal components and their associated eigenvectors are employed in the process of signal

15     reconstruction. To reflect significant oscillatory modes, the original record can be filtered by

16     choosing $p\,(<m)$ RCs from all $m$ RCs.

17     **3.2 Determination of Parameters $\left(\tau, m\right)$**

18     The choice of the parameter pair ($\tau, m$) in SSA are not as strict as in PSR. The

19     $\tau$ and $m$ from SSA and PSR are generally different although some authors considered them as

20     same in application (e.g., *Sivapragasam et al., 2001*). The choice of $m$ represents a compromise

21     between information content and statistical confidence. As a basic requirement, the

1   chosen $m$ should be able to resolve obvious different oscillations hidden in the original signal.

2   In other words, some leading eigenvalues should be identified. Figs. 7 and 8 show the

3   sensitivities of the eigenvalue decomposition to the window length $m$ and the lag

4   time $\tau$ respectively. Results from Fig. 7 show that 5 leading eigenvalues stand out for

5   Xiangjiaba, Cuntan, and Manwan whilst 7 leading eigenvalues can be identified for

6   Danjiangkou when $m \geq 30$. Each leading eigenvalue corresponds to a distinctive oscillation.

7   Results from Fig. 8 suggest that the curve of eigenvalues of each series is insensitive to the

8   lagged time $\tau$ when $m = 30$. The $\tau$ was consistently chosen as 1. The $m$ of 30 may be too large

9   for the analysis below. For simplicity, the final parameter pair $(\tau, m)$ in SSA were (1, 5) for

10  Xiangjiaba, Cuntan, and Manwan and (1, 7) for Danjiangkou.

11      Taking Danjiangkou as an example, Fig. 9 presents 7 RCs and the original series

12  excluding the validation data. The RC1 represents the lowest frequency oscillation and the RC7

13  reflects the highest frequency oscillation. The high-frequency components stand for noise

14  signals to some extent. Potentially, they may be filtered out so as to clean the original series.

15  Fig. 10 depicts AMI and CCF between RCs and the original streamflow series. The last plot

16  presents the average of the AMI and CCF from the previous plots. The average value indicates

17  an overall correlation either being positive or negative. As shown in Fig. 10, the value of CCF

18  from each RC is significantly different with the increase of the lag time. It means that the

19  contribution of each RC to the output of models is changeable at different prediction levels.

20  Therefore, the number $p\,(<m)$ of RCs at different prediction levels will be chosen by trial and

21  error according to the minimum performance of model (Root Mean Square Error was taken in

22  this study).

## 1  4. Description of Models

2      Fig. 11 depicts the framework of implementation of all five data-driven models in this

3  study. They are ARMA, KNN, ANN, CDANN, and CDSVR from the top to the bottom. The

4  configuration of each model except for ARMA was introduced as follows. In view of the

5  similarity of the flow characteristics among Xiangjiaba, Cuntan, and Manwan, two streamflow

6  series, Xiangjiaba and Danjiangkou, were therefore chosen for forecasting in the rest of this

7  paper.

8  **4.1 KNN**

9      The KNN was applied to estimate the function of $f(\bullet)$ in Eq. (2) in the study. The KNN

10  algorithm depends on observations that are in some finite neighborhood of the point of estimate.

11  The basic idea behind the KNN is to break up domain into local neighborhoods and fit

12  parameters in each neighborhood separately. To predict $x^{F}_{t+T+(m-1)\tau}$ based on $\mathbf{Y}(t)$ and the past

13  ones, $k$ nearest neighbors of $\mathbf{Y}(t)$ should be found. Jayawardena and Lai (*1994*) summarized

14  five local prediction approaches according to the $k$ being equal to 1 or larger than 1: zeroth

15  order approximation when $k=1$; for $k>1$, the other four are respectively average $k$ neighbours,

16  weighted average of $k$ neighbours, exponential weighting $k$ neighbours, and linear

17  approximation via fitting $k$ pairs $(\mathbf{Y}(t'_i),\mathbf{Y}(t'_i+\tau))$. In the current study, the weighted averaging

18  $k$ neighbors approach is adopted and the weight function is proportional to the inverse of

19  square Euclidean distance between $\mathbf{Y}(t)$ and $\mathbf{Y}(t')$, which can be referred to Wu et al. (*2008*).

1 **4.2 ANN**

2     The single ANN used in this study was a static feed-forward multilayer perceptron

3 (MLPs). The static ANN is able to capture the dynamics of a system in the network model by

4 using delayed time inputs. The architecture design of the ANN consists of the number of

5 hidden layers and the number of neurons in input layer, hidden layers and output layer. ANNs

6 with one hidden layer are commonly used in hydrologic modeling (*Dawson and Wilby, 2001;*

7 *de Vos and Rientjes, 2005*) since these networks are considered to provide enough complexity

8 to accurately simulate the dynamic and nonlinear-properties of the hydrologic process.

9 Therefore, a three-layer static ANN was finally chosen for the present study, which comprises

10 the input layer with $m$ neurons, the hidden layer with $h$ nodes, and the output layer with one

11 node. As mentioned previously, $m$ is 4 for Xiangjiaba and 5 for Danjiankou. The hyperbolic

12 tangent was chosen as the activation function for the neurons of the hidden layer, which usually

13 requires rescaling both the input and outputting data to [-1, 1]. Based on the form of Eq. (2), the

14 prediction model is mathematically expressed as

15 $$x_{t+T+(m-1)\tau}^{F} = f(\mathbf{Y}(t), w, \theta, m, h) = \theta_0 + \sum_{j=1}^{h} w_j^{out} \tanh(\sum_{i=1}^{m} w_{ji} x_{t+(i-1)\tau} + \theta_j) \qquad (9)$$

16 where $x_{t+(i-1)\tau}$ are the $m$ components of the state $\mathbf{Y}(t)$; $w_{ji}$ are the weights defining the link

17 between the *ith* node of the input layer and the *jth* of the hidden layer; $\theta_j$ are biases associated

18 to the *jth* node of the hidden layer; $w_j^{out}$ are the weights associated to the connection between

19 the *jth* node of the hidden layer and the node of the output layer; and $\theta_0$ is the bias at the output

20 node.

21     As mentioned in Section 2.3, $m$ is 4 for Xiangjiaba and 5 for Danjiankou for a unique

22 predicted output. The ensuing task is to optimize the size $h$ of the hidden layer with the

1    determined input and output neurons. In terms of CDANN and CDSVR, they employed the

2    same input and output for each single model of them as the ANN.

3    **4.3 Distributed Models**

4    **(1) FCM** (*Bezdek, 1981; Wang et al., 2006b*)

5        The Fuzzy C-Means (FCM) method partitions a set of N vector $\mathbf{Y}_j, j=1,...,N$, into c

6    fuzzy clusters, and each data point belongs to a cluster to a degree specified by a membership

7    grade $u_{ij}$, between 0 and 1. We define a matrix $\mathbf{U}$ comprising the elements $u_{ij}$, and assume that

8    the summation of degrees of belonging for a data point is equal to 1, i.e., $\sum_{i=1}^{c} u_{ij} = 1$,

9    $\forall j = 1, \cdots, N$. The goal of the FCM algorithm is to find c cluster centers so that the cost

10   function of dissimilarity measure is minimized. The cost function can be defined by

11
$$J(\mathbf{U}, v_1, \cdots v_c) = \sum_{i=1}^{c} J_i = \sum_{i=1}^{c} \sum_{j=1}^{N} u_{ij}^l d_{ij}^2 \tag{10}$$

12   where $v_i$ is the cluster center of the fuzzy subset $i$ ; $d_{ij} = \|v_i - \mathbf{Y}_j\|$ is the Euclidean distance

13   between the *i*th cluster center and *j*th data point; and $l \geq 1$ is a weighting exponent, taken as 2

14   here so as to match the square Euclidean distance. The necessary conditions for Eq. (10) to

15   reach its minimum are:

16
$$v_i = \sum_{j=1}^{N} u_{ij}^l \mathbf{Y}_j \bigg/ \sum_{j=1}^{N} u_{ij}^l \tag{11}$$

17
$$u_{ij} = \left[ \sum_{k=1}^{c} \left( \frac{d_{ij}}{d_{kj}} \right)^{2/(l-1)} \right]^{-1} \tag{12}$$

1    The FCM algorithm is an iterative procedure that satisfies Eqs. (11) and (12) to

2    minimize Eq. (10). Implementation of the algorithm can be referred to Bezdek (*1981*) and

3    Wang et al. (*2006b*) for details.

4    The number c of clusters was taken to be 3 in this study. The three subsets were

5    expected to represent three different magnitudes of flows, i.e., low flow, medium flow, and

6    high flow.

7    **(2) SVR**

8    A nonlinear SVR was used in this study. A brief introduction of SVR was presented

9    here. The details of SVR can be referred to Kecman (*2001*), Yu et al. (*2006*) and Wu et al.

10    (*2008*). According to equation (2), the underlying function $f(\bullet)$ in the context of a nonlinear

11    SVR is given by

$$x^{F}_{t+T+(m-1)\tau} = f(\mathbf{Y}(t), \omega) = \omega \cdot \phi(\mathbf{Y}(t)) + b \qquad (13)$$

13    where the input data $\mathbf{Y}(t)$ in the input space is mapped to a high dimensional feature space via

14    a nonlinear mapping function $\phi(\mathbf{Y}(t))$. The objective of the SVR is to find optimal $\omega, b$ and

15    some parameters in kernel function $\phi(\mathbf{Y}(t))$ so as to construct an approximation function of

16    the $f(\bullet)$.

17    When introducing Vapnik's $\varepsilon$-insensitivity error or loss function, the loss function

18    $L_{\varepsilon}(y, f(\mathbf{Y}(t), \omega))$ on the underlying function can be defined as

$$L_{\varepsilon}(y, f(\mathbf{Y}(t), \omega)) = \left| y - f(\mathbf{Y}(t), \omega) \right|_{\varepsilon} = \begin{cases} 0 & if \left| y - (\omega \cdot \phi(\mathbf{Y}(t)) + b) \right| \le \varepsilon \\ \left| y - (\omega \cdot \phi(\mathbf{Y}(t)) + b) \right| - \varepsilon & otherwise \end{cases} \qquad (14)$$

20    where $y$ represents observed value. Similar to linear SVR (*Kecman, 2001; Yu et al., 2006*), the

21    nonlinear SVR problem can be expressed as the following optimization problem:

1

$$minimize \quad R_{W,\xi_i,\xi_i^*} = \frac{1}{2}\|\omega\|^2 + C\sum_{i=1}^{N}(\xi_i + \xi_i^*)$$

$$subject \quad to \quad \begin{cases} y_i - f(\phi(\mathbf{Y}_i),\omega) - b \leq \varepsilon + \xi_i \\ f(\phi(\mathbf{Y}_i),\omega) + b - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0 \end{cases} \quad (15)$$

2 where $\mathbf{Y}_i$ represents $\mathbf{Y}(i)$ for simplicity, the term of $\frac{1}{2}\|\omega\|^2$ reflects generalization, and the term

3 of $C\sum_{i=1}^{N}(\xi_i + \xi_i^*)$ stands for empirical risk. The objective in Eq. (15) is to minimize them

4 simultaneously, which implements SVR to avoid underfitting and overfitting the training data.

5 $\xi_i$ and $\xi_i^*$ are slack variables for measurements "above" and "below" an $\varepsilon$ tube. Both slack

6 variables are positive values. $C$ is a positive constant that determines the degree of penalized

7 loss when a training error occurs.

8 By introducing a dual set of Lagrange multipliers, $\alpha_i$ and $\alpha_i^*$, the minimization

9 problem can be solved in a dual space. The objective function in dual form can be represented

10 as (*Gunn, 1998*)

11

$$maximize \quad L_d\left(\alpha,\alpha^*\right) = -\varepsilon\sum_{i=1}^{N}\left(\alpha_i^* + \alpha_i\right) + \sum_{i=1}^{N}\left(\alpha_i^* - \alpha_i\right)y_i - \frac{1}{2}\sum_{i,j=1}^{N}(\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j)\left(\phi(\mathbf{Y}_i)\cdot\phi(\mathbf{Y}_j)\right)$$

$$subject \quad to \quad \begin{cases} \sum_{i=1}^{N}(\alpha_i - \alpha_i^*) = 0 \\ 0 \leq \alpha_i^* \leq C, \qquad i = 1,\cdots,N \\ 0 \leq \alpha_i \leq C, \qquad i = 1,\cdots,N \end{cases} \quad (16)$$

12 By using a "kernel" function $K(\mathbf{Y}_i,\mathbf{Y}_j) = \left(\phi(\mathbf{Y}_i)\cdot\phi(\mathbf{Y}_j)\right)$ to yield inner products in

13 feature space, the computation in input space can be performed. In the present study, Gaussian

14 radial basis function (RBF) was adopted in the form of $K(\mathbf{Y}_i,\mathbf{Y}_j) = \exp(-\|\mathbf{Y}_i - \mathbf{Y}_j\|^2/2\sigma^2)$.

21

1   Once parameters $\alpha_i$, $\alpha_i^*$, and $b_0$ are obtained, the final approximation function of the $f(\bullet)$

2   becomes

3
$$f(\mathbf{Y}_i) = \sum_{i=1}^{N} (\alpha_k - \alpha_k^*) K(\mathbf{Y}_k \cdot \mathbf{Y}_i) + b_0 \, , k = 1, \cdots, s \qquad (17)$$

4   where $\mathbf{Y}_k$ stands for the support vector, $\alpha_k$ and $\alpha_k^*$ are parameters associated with support

5   vector $\mathbf{Y}_k$, $N$ and $s$ represent the number of training samples and support vectors respectively.

6   Three parameters $(C, \varepsilon, \sigma)$, however, have to be optimized first before Eq. (17) is used to

7   perform the forecasting.

8   **(3) CDANN and CDSVR**

9        As depicted in Fig. 11, CDANN (or CDSVR) consists of three single ANNs (or SVR).

10  A basic idea behind the type of model is that the training set is first split into three subsets by

11  the FCM, and then each subset is fitted by a single ANN (or SVR). The three subsets are

12  associated with subset1, subset2 and subset3 in Fig. 11. Generally, the final output of the

13  distributed model is a weighted sum of the outputs of all single models when making a forecast

14  from a new input vector, and the weights are determined according to the distance between the

15  new input vector and each cluster center (*Cheng et al., 2006; Wu et al.,2008*). As mentioned in

16  Introduction, the DANN (or DSVR) dependent upon weighted output was called SDANN (or

17  SDSVR). They tend to generate poor forecasts because ANN (or SVR) is unable to extrapolate

18  beyond the range of the data used for training. In the CDANN (or CDSVR), only one single

19  ANN (or SVR) model, whose clustering center is the nearest to the new input, is triggered for

20  the forecasting. As a consequence, the final output of the distributed models is directly from the

21  output of the triggered ANN (or SVR). Therefore, the validation set is necessarily split into

22  three parts in advance before the forecasting or simulation by CDANN (or CDSVR).

22

1 ## 5. Application of Models

2 **5.1 Data preparation**

3    Streamflow series data were divided into three parts: model training, cross-validation

4 and validation. The last ten years' streamflow data, called validation set, were set aside for

5 validation. This validation set would not be used until all model development and training was

6 finished completely. Of the remaining data, the first two-thirds called training set was for

7 model training, and the other one-third called cross-validation set was for the purpose of

8 confirming and validating the initial analysis. Models based on soft computing method have

9 difficulties to extrapolate beyond the range of the data used for training. As a consequence, it is

10 imperative that the training and validation sets are representative of the same population.

11 Statistical properties (mean, deviation, range) from them are compared in order to measure the

12 representative. According to the statistical properties of these data sets, the division of data is

13 satisfied since the statistical parameters of the training sets for all the streamflow series are

14 closed to the cross-validation sets and the testing sets.

15    With suitable data preparation beforehand, it is possible to improve the performance of

16 models although these soft computing models (ANN, CDANN, and CDSVR) are very powerful

17 to handle nonlinear, noisy, and non-stationary data. Data-preprocessing can ensure that all

18 variables receive equal attention during the training process (*Dawson and Wilby, 2001*). Also,

19 data-preprocessing may be a need to improve the training efficiency of ANN (or SVR)

20 (*Sudheer et al., 2003*). Data-preprocessing techniques include signal filtering, data-

21 transformation, rescaling or standardization. In this study, the signal filtering and data-

22 transformation were respectively the SSA and Moving Average (MA) over input signals. The

1  rescaling operation was applied to models of ANN, CDANN, and CDSVR, which rescaled the

2  input and output data to [-1, 1].

**5.2 Evaluation of model performances**

4  The measures of evaluating model performance that are used in the present paper

5  comprise Root Mean Square Error (RMSE), the Nash-Sutcliffe Coefficient of Efficiency (CE)

6  (*Nash and Sutcliffe, 1970*), and the Persistence Index (PI) (*Kitanidis And Bras, 1980*). They are

7  respectively formulated as: $RMSE = \sqrt{\dfrac{1}{N}\sum_{i=1}^{N}(T_i - \hat{T}_i)^2}$ ), $CE = 1 - \sum_{i=1}^{N}(T_i - \hat{T}_i)^2 \Big/ \sum_{i=1}^{N}(T_i - \bar{T})^2$

8  and $PI = 1 - \sum_{i=1}^{N}(T_i - \hat{T}_i)^2 \Big/ \sum_{i=1}^{N}(T_i - T_{i-L})^2$. In these equations, $N$ =number of observations, $\hat{T}_i$ =forecasted

9  streamflow, $T_i$ =observed streamflow, $\bar{T}$ =average observed streamflow, and $T_{i-L}$ is the

10  streamflow estimate from a so-call persistence model (or called naïve model) that basically

11  takes the last streamflow observation (at time $i$ minus the lead time $L$) as a prediction. CE and

12  PI values of 1 stand for perfect fits.

**5.3 Optimization scheme of model parameters**

14  The ARMA method was used as one of the benchmark models in the present study.

15  Parameters (*P, Q*) in ARMA were estimated by trial and error with each of them varied over

16  the range of 0 to 12, excluding the unfeasible case where both of them are simultaneously equal

17  to 0, in view of the predominant periodicity of 12 months (Fig. 2). The best pair of (*P, Q*) is

18  associated with the minimal value of Akaike Information Criterion (AIC).

19  The number of $k$ in KNN in this study was chosen as $k = m+1$ which was suggested

20  by Sugihara and Mary (*1990*). As a result, a fixed $k = 5$ and $k = 6$ were respectively adopted in

21  Xiangjiaba (recall that $m = 4$) and Danjiangkou (recall that $m = 5$).

1    The optimal size $h$ of the hidden layer was found by systematically increasing the

2    number of hidden neurons from 1 to 10 until the network performance on the cross-validation

3    set was no longer improved significantly. The optimization procedure can be begun if ANN is

4    equipped with necessary training algorithm in advance. In our preliminary investigations (Wu

5    et al., Data-driven models for monthly streamflow time series prediction, submitted to Journal

6    of Hydrological Process, 2009; hereinafter referred to as *Wu et al., submitted manuscript, 2009*),

7    we tested performance of three training algorithms of Levenberg-Marquardt (LM), Particle

8    Swarm Optimization (PSO), and Levenberg Marquart combined with Genetic algorithm (LM-

9    GA). The results showed that the LM-GA performed best. Thus, the LM-GA algorithm was

10   used for training in this study. Finally, a 4-8-1 configuration of the ANN model with the LM-

11   GA algorithm was designed for one-month-ahead forecast of Xiangjiaba. Similarly, a 5-5-1

12   configuration of the ANN model was found for one-month-ahead forecast of Danjiangkou. For

13   other prediction levels, the configurations of ANN and CDANN were also found by the same

14   trial and error method.

15   The CDSVR had the same inputs and output of model as ANN and CDANN. To

16   conduct forecasting, three parameters ($C, \varepsilon, \sigma$) in CDSVR is necessarily determined when

17   Gaussian RBF was adopted as the kernel function. Some techniques to determine these

18   parameters have been recommended in literature. For instance, Cherkassky and Ma (*2004*),

19   from the perspective of analytic 'rule-of-thumb' selection, suggested selection of parameter

20   $C$ based on the formula: $\dfrac{C}{N} = \max(\left|\bar{y} + 3\sigma_y\right|, \left|\bar{y} - 3\sigma_y\right|)$, where $\bar{y}$ and $\sigma_y$ are the mean and the

21   standard deviation of the training data $y$, and $N$ is the data size of training data. $\varepsilon$ is suggested

22   as $\varepsilon = 3\sigma_a \sqrt{\ln N / N}$, where $\sigma_a{}^2$ is the variance of additive noise $\delta$. Global evolutionary

1    techniques are widely reported for optimization of these parameters. Hsu et al. (*2003*)

2    recommend a "grid-search" on $C$ and parameter $\sigma$ in RBF using cross-validation. Yu et al.

3    (*2004*) and Lin et al. (*2006*) employed Shuffled Complex Evolution (SCE-UA) to search these

4    parameters. A two-step grid search method was applied to find the optimal parameter triplet

5    ($C,\varepsilon,\sigma$) by Yu et al. (*2006*) and Wu et al. (*2008*). In the present study, the two-step GA was

6    found the best in terms of effectiveness and efficiency of optimization compared to two-step

7    PSO and two-step SCEUA.  Thus, the two-step GA method was employed in this study.

8          There are other two important parameters, the number $p$ of chosen reconstructed

9    components and the memory length of moving average. They were chosen by trial and error

10    and the final values of them were associated with the minimum RMSE.

## 6.  Results and Discussion

12    **5.1 Main results**

13          Table 1 shows the results of one-month-ahead prediction of all five models without data

14    preprocessing of MA and SSA using the monthly data from Xiangjiaba and Danjiangkoufor

15    respectively. These models made better accuracy for Xiangjiaba data than Danjiangkou data.

16    ARMA gave the best performance for Xiangjiaba and the second best performance for

17    Danjiangkou as compared to the other models, which is inconsistent with the conclusion in

18    literature that ANN mostly outperforms ARMA (*Jain et al. 1999; Abrahart and See, 2002;*

19    *Castellano-Meʹndeza et al. 2004*). This is most likely due to the ANN (CDANN or CDSVR)

20    model not being optimal considering that the model input was identified by reconstruction of

21    dynamics.

1    Figs. 12 and 13 present the scatter plots and the plots of the observed and forecasted

2    time series of one-month-ahead forecast for Xiangjiaba and Danjiangkou using the ANN and

3    KNN models. The biggest drawback is that each model mismatched quite a number of peak

4    flows and the ANN model always underestimated the peak flows of both Xiangjiaba and

5    Danjiangkou. Fig. 13 shows that the ANN model's timing of the peaks is markedly lagged.

6    This prediction lag effect is the result of using previously observed values as ANN inputs. The

7    ANN gives the most weight to the latest discharge input for calculating the forecast because the

8    forecast depends on the autocorrelation of the monthly time series. Therefore, the MA over

9    inputs was suggested to overcome the lag effect by de Vos and Rientjes (*2005*). Compared with

10   Fig. 12, Fig. 13 also demonstrates that it is difficult for any powerful model to fit a signal

11   contaminated by complex noises.  However, it is possible to improve the model performance by

12   using the SSA technique to clean the original signal.  This was explored in more detail in Fig.

13   14 by using KNN model to conduct four different lead forecasts. For Xiangjiaba data, the SSA

14   can significantly improve RMSE of the KNN model at one- and three-month forecast horizons.

15   For Danjiangkou data, the RMSE of the KNN model greatly decreases at all forecast horizons

16   except for twelve-month horizon.

17   Fig. 15 depicts the results of one-month-ahead forecast of the ANN model with the help

18   of MA and SSA for Xiangjiaba and Danjiangkou. Compared with the scatter plots from ANN

19   model in Figs. 12 and 13, the current scatter plots with low spread, and the low RMSE and high

20   CE and PI indicate better model performances. The peak flows still cannot be fitted perfectly.

21   Also, there is an overall underestimate of the discharge of Danjiangkou by the ANN model

22   because the scatter plot is wholly skew to the axis of observed values, which implies the

23   identified ANN model is not optimal.

1    The results of one-month-ahead forecast of CDSVR with the aid of MA and SSA are

2    presented in Fig. 16. The scatter plot with perfect match of the diagonal indicates excellent

3    model performance.  The plot of observed and forecasted discharges shows that the peak flows

4    are also fitted perfectly.  In particular, the whole forecasted error generated by ANN model in

5    Fig. 15 for Danjiangkou is also overcome by the CDSVR model. A further discussion can be

6    found in the later analysis of errors.

7    Figs. 17 and 18 show the validation data and the comparison of absolute prediction

8    errors of various models for Xiangjiaba and Danjiangkou. Figs. 17a and 18a show that the

9    extreme forecast errors of ARMA model occurred at the last points, which indicates the ARMA

10   model with the current technique of SSA is not viable for a real-time prediction. Figs. 17b and

11   18b compare three soft computing models of ANN, CDANN and CDSRV. Error curves of

12   ANN and CDANN have very similar trends with obvious underestimates at most points. The

13   error curve of CDANN displays that the peak flows were mostly well simulated, and the errors

14   tend to be random. Fig. 19 presents one-step-ahead absolute forecasting error distributions of

15   KNN, ANN, CDANN, and CDSVR using Danjiangkou data. The histograms of errors for KNN

16   and CDSVR are characterized by a quasi-normal distribution with mean near zero, which

17   indicates that errors for them basically pass the test for "whiteness". However, the histogram of

18   errors for ANN is evidently skewed to the left which means that the ANN model was not

19   optimal. The histogram of errors for CDANN is also skewed although there is a bit

20   improvement compared to the ANN model. The issue often occurs in the application of ANN.

21   In practice, it is not easy to obtain the optimal configuration of ANN although ANN is able to

22   map any complex relationship.

1    Table 2 shows the one-month-ahead forecast performance in terms of RMSE, CE, and

2    PI, and relevant parameters including model parameters and data-preprocessing parameters of

3    all five models. The data-preprocessing parameters were decided by trial and error. Compared

4    to the results in Table 1, data-preprocessing techniques of MA and/or SSA can considerably

5    improve the performance of all models. Prediction was also extended to longer prediction

6    horizons, namely $T = 3$ months, 6 months, and 12 months. The static prediction method was

7    employed for multiple step prediction, and hence relevant parameters of each model at

8    different forecast horizons have to be optimized again. The results of three-, six-, and twelve-

9    month-ahead forecasts are given in Tables 3-5. The CDSVR model outperforms the other four

10   models at all prediction horizons for Xiangjiaba, and the CDSVR model performs the best

11   among all models at all prediction horizons expect for at six-month-ahead forecast horizon for

12   Danjiangkou. At six-month-ahead forecast horizon, the KNN method performs the best.  The

13   data-preprocessing parameters from Tables 2-5 show that MA cannot affect the model

14   performance of ARMA and KNN whereas SSA imposes a great impact on the model

15   performance of KNN. From the perspective of catchment data of the current study, MA is

16   suitable for each catchment data whilst SSA is more effective for Danjiangkou data.

17   **5.2 Discussion**

18       The following part concerns some discussions about forecasting models and the effects

19   of data-preprocessing techniques.

20   **(1) About parameter $k$ in KNN**

21       The parameter of $k$ in KNN model poses a great impact on the performance of KNN.

22   The choice of $k$ should ensure the reliability of the solution although a preliminary value of $k$

23   is based on $k = m+1$. The check of robustness of $k$ in terms of RMSE of validation data was

1    presented in Fig. 20, where $k$ was in the interval of [2, 40]. Adopting the value of $k$ as $m+1$

2    was reasonable because its RMSE was only 2.3% larger than the minimum RMSE

3    (corresponding to $k=4$ ) for Xiangjiaba, and 2.5% larger than the minimum RMSE

4    (corresponding to $k=20$) for Danjiangkou.

5    **(2) About models of ANN, CDANN, and CDSVR**

6           The performance of ANN (CDANN or CDSVR) is considerably affected by the

7    identified configuration. Recall that the three soft computing models had the same inputs and

8    output of model, which was derived from the reconstruction of dynamics. The purpose is to

9    fairly compare these models. In terms of the one-month-ahead forecast for Danjiangkou, the

10   configuration of ANN with 12-5-1 achieved better performance where the number of input

11   neurons was optimized by trial and error. Obviously, the model inputs from the reconstruction

12   of dynamics are not the optimal. However, the basic conclusion in this study remains

13   unchanged because the three soft computing models are compared based on the same model

14   inputs.

15          The distributed models improved the performance of ANN by using local fitting instead

16   of global fitting but at the expense of the increase of the training time. The CDANN is actually

17   three single ANNs. Hence, it has the same disadvantages as the ANN that the forecast results

18   tend to be unstable. The CDSVR, however, exhibited more stable forecast results and better

19   generalization but need longer training time. In view of the fact that training time will

20   exponentially increase with the number of training samples (*Wu et al., 2008*), the SVR model is

21   suitable for samples less than hundreds. Conversely, the ANN model prefers more samples to

22   fewer samples so as to improve the generalization of model. Therefore, ANN and SVR are not

23   always alternative, and are sometimes complementary.

1    Tables 2-5 show that the number of clustering centers for CDSVR was 3 for one-

2    month-ahead forecast and 3 or 6 for longer forecast horizon. It was found that the performance

3    of CDSVR can be improved at some forecast horizons if the number of clustering centers

4    increases. In the meantime, the increase of clustering centers can remarkably promote the

5    training efficiency of the CDSVR model.

6    **(3) About the investigation of effects of MA and SSA**

7    To investigate the effect of MA on the ANN model, Fig. 21 compares the Cross-

8    Correlation Functions (CCFs) between inputs and output from ANN with/without MA for one-

9    month-ahead forecast. For Xiangjiaba whose original data has high autocorrelation, the MA

10   decreases the correlation relationship between the latest two input components and the output

11   of model whereas the MA increases the correlation relationship between other three input

12   components and the output of model. For Danjiangkou whose raw data has low autocorrelation,

13   the MA increases the correlation relationship between each of seven input components and the

14   output of model.

15   Fig. 22 investigates the combined effect of MA and SSA on the ANN model by

16   computing the CCFs between inputs and output of model for a six-month-ahead forecast. Fig.

17   22 shows that there is a common characteristic for Xiangjiaba and Danjiangkou that the

18   combined effect of MA and SSA is to mainly improve the correlation relationship between the

19   more previous input components and the output of model. In fact, the MA and SSA decrease

20   the contribution of the latest component to the output of model, which potentially overcomes

21   the timing error of forecast.

31

## 7. Conclusions

In this paper, five data-driven models, ARMA, KNN, ANN, CDANN, and CDSVR, were adopted to conduct multiple-step monthly discharge forecasts by two catchment data, Xiangjiaba and Danjiangkou. In order to improve the model performance, two data-preprocessing techniques of the MA and SSA were equipped to each model. The results of the one-month-ahead forecasts showed that the MA and/or SSA can considerably improve the performance of each model. The results of the multiple-step forecasts demonstrated that CDSVR stood out from other models except for at the six-month-ahead forecast horizon on Danjiangkou data. In terms of the CDSVR model, the model performance deteriorated with the increase of the forecast horizon. For the perspective of models, the SSA had positively effects on ARMA and KNN whereas the MA had pivotal influences on the improvement of model performance of ANN, CDANN and CDSVR. For the perspective of two different catchments, the MA is suitable for each catchment whilst SSA is more effective for the data of Danjiangkou.

Some discussions were presented (1) on the number of neighbors in KNN; (2) on the configuration of ANN; and (3) on the investigation of effects of MA and SSA. First of all, adopting the number of neighbors $k$ as $m+1$ was considered to be reasonable. Secondly, compared to CDSVR, the poor forecasts from ANN were partly due to the identified configuration of ANN not being optimal. The CDANN model should be expected to have model performance as good as that of CDSVR because it also adopted the local fitting (or distribution) approach. However, obtaining three ANN's configurations with good generalization for CDANN is difficult in view of fewer training data in this study. Conversely, the CDSVR model preferred fewer training data which exponentially improves its training speed but do not obviously affect its generalization due to the implementation of structural risk

1    minimization (SRM) in SVR (*Kecman,2001*).    Finally, the techniques of MA and SSA

2    tremendously improved the model performance of ANN (CDANN or CDSVR) by adjusting the

3    correlation relationship between input components and output of models (i.e., increasing the

4    correlation relationship between the earlier input components and the output of model).

5

# References

Abraham, N.B., A.M. Albano, B. Das, G. de Guzman, S. Yong, R.S. Gioggia, G.P. Puccioni, and J.R. Tredicce (1986), Calculating the dimension of attractors from small data. *Phys. Lett. A*, 114, 217–221.

Abarbanel, H. D. I., R. Brown, J. J. Sidorowich, and L.S. Tsimring (1993), The analysis of observed chaotic data in physical systems, *Reviews of Modern Physics*, 65(4), 1331-1392.

Abrahart, R.J., and L. See (2002), Multi-model data fusion for river flow forecasting: an evaluation of six alternative methods based on two contrasting catchments, *Hydrology and Earth System Sciences*, 6(4), 655-670.

Behzad, M., K. Asghari, M. Eazi, and M. Palhang (2009), Generalization performance of support vector machines and neural networks in runoff modeling, *Expert Systems with Applications*, 36 (4), 7624-7629.

Bezdek, J.C. (1981), *Pattern Recognition with Fuzzy Objective Function Algorithms*. Plenum Press, New York.

Cannas, B., A. Fanni, L. See, and G. Sias, (2006), Data preprocessing for river flow forecasting using neural networks: Wavelet transforms and data partitioning, *Physics and Chemistry of the Earth*, 31, 1164-1171.

Carlson, R.F., A.J.A. MacCormick, and D.G. Watts (1970), Application of linear models to four annual streamflow series, *Water Resour. Res.*, 6(4), 1070-1078.

Castellano-Méndeza, M., W. González-Manteigaa, M. Febrero-Bande, J. Manuel Prada-Sáncheza, and R. Lozano-Calderón (2004) Modeling of the monthly and daily behavior of the runoff of the Xallas river using Box–Jenkins and neural networks methods, *Journal of Hydrology*, 296, 38–58.

Cheng, J., J.S. Qian, and Y.N. Guo (2006), A Distributed Support Vector Machines Architecture for Chaotic Time Series Prediction. Neural Information Processing (ICONIP 2006), Proc. 13th International Conference, Part 1, Hong Kong, China, October 3-6, 2006, Springer.

Cherkassky, V., and Y. Ma (2004), Practical selection of SVM parameters and noise estimation for SVM regression, *Neural Networks*, 17, 113–126.

Corzo, G. and Solomatine, D. P. (2007). Baseflow separation techniques for modular artificial neural network modelling in flow forecasting. *Hydrological Science Journal*, 52 (3), 491–507.

Dawson, C. W. and Wilby, R. L.,(1999), A comparison of artificial neural networks used for river flow forecasting, *Hydrol. Earth Sys. Sci.*,3, 529–540.

Dawson, C. W., and R. L. Wilby (2001). Hydrological Modeling Using Artificial Neural Networks, *Progress in Physical Geography*, 25(1), 80-108.

De Vos, N.J. and T.H.M. Rientjes (2005), Constraints of artificial neural networks for rainfall -runoff modeling: trade-offs in hydrological state representation and model evaluation, *Hydrology and Earth System Sciences*, 9, 111-126.

Dibike, Y. B., S. Velickov, D. Solomatine, and M. B. Abbott, (2001). Model Induction with support vector machines: Introduction and Application. *Journal of Computing in Civil Engineering*, 15(3), 208‐216.

Elshorbagy, A., S.P. Simonovic, and U.S. Panu (2002), Estimation of missing stream flow data using principles of chaos theory, *Journal of Hydrology*, 255, 123–133.

Elsner, J. B. and A. A. Tsonis (1997), *Singular Spectrum Analysis: A New Tool in Time Series Analysis*. Plenum Press, New York.

1    Essex,C., and M.A.H. Nerenberg (1991), Proc. R. Soc. London, Series A 435, 287.

2    Farmer, J. D., and J. J. Sidorowich (1987), Predicting chaotic time series, *Phys. Rev. Lett.*, 59(4), 845–848.

3    Fraser, A.M. and H.L. Swinney (1986), Independent coordinates for strange attractors from mutual information,
4    *Physical Review A*, 33(2), 1134-1140.

5    Ghilardi, P., and R. Rosso (1990), Comment on chaos in rainfall, *Water Resour. Res.* 26 (8), 1837–1839.

6    Grassberger, P., and I. Procaccia (1983), Measuring the strangeness of strange attractors, *Physica 9D*, 189-208.

7    Haltiner, J.P., and J.D. Salas (1988), Short-term forecasting of snowmelt discharge using ARMAX models, *Water
8    Resources Bulletin*, 24(5), 1083-1089.

9    Hong, S.Z., and S.M., Hong (1994), An amendment to the fundamental limits on dimension calculations. *Fractals
10    *2 (1), 123–125.

11    Hsu, C. W., C.C. Chang, and C. J. Lin (2003), A practical guide to support vector classification, Available at
12    http://www.csie.ntu.edu.tw/cjlin/papers/guide/guide.pdf. [Accessed on 20/6/07].

13    Huang, W.R., B. Xu, and A. Hilton (2004), Forecasting Flows in Apalachicola River Using Neural Networks,
14    *Hydrological Processes*, 18, 2545-2564.

15    Jain, S. K., A. Das, and D. K. Drivastava (1999). Application of ANN for reservoir inflow prediction and
16    operation. *J. Water. Resour. Plann. Manage.*, 125(5), 263–271.

17    Jain, A., and S. Srinivasulu (2006), Integrated approach to model decomposed flow hydrograph using artificial
18    neural network and conceptual techniques, *Journal of Hydrology*, 317, 291–306.

19    Jain, A., and S. Srinivasulu (2004), Development of effective and efficient rainfall-runoff models using integration
20    of deterministic, real-coded genetic algorithms and artificial neural network techniques, *Water Resour. Res.*, 40,
21    W04302.

22    Jayawardena, A.W., and A.B. Gurung (2000), Noise reduction and prediction of hydro-meteorological time series:
23    dynamical systems approach vs. stochastic approach, *Journal of Hydrology*, 228, 242–264.

24    Jayawardena, A.W., and F. Lai (1994), Analysis and prediction of chaos in rainfall and stream flow time series,
25    *Journal of hydrology*, 153, 23-52.

26    Kantz, H., and T. Schreiber (2004). *Nonlinear time series analysis (2nd edition)*. Springer.

27    Kecman, V. (2001), *Learning and soft computing: support vector machines, neural networks, and fuzzy logic
28    models*, MIT press, Cambridge, Massachusetts.

29    Kennel, M. B., R. Brown, H. D. I. Abarbanel (1992), Determining embedding dimension for phase space
30    reconstruction using geometrical construction, *Phy. Rev. A.*, 45(6), 3403-3411.

31    Kişi, O. (2003). River flow modeling using artificial neural networks, *Journal of Hydrologic Engineering*, 9(1),
32    60-63.

33    Kişi, O. (2005). Daily river flow forecasting using artificial neural networks and auto-regressive models, *Turkish
34    J. Eng. Env. Sci,* 29, 9-20.

35    Kitanidis, P. K. and R. L. Bras (1980), Real-time forecasting with a conceptual hydrologic model, 2, applications
36    and results, *Wat. Resour. Res.*, 16 (6), 1034–1044.

Kothyari, U.C., and V.P. Singh (1999), A multiple-input single-output model for flow forecasting, *Journal of Hydrology*, 220, 12–26.

Koutsoyiannis, D., and D. Pachakis (1996), Deterministic chaos versus stochasticity in analysis and modeling of point rainfall series. *J. Geophys. Res*. 101 (D21), 26441–26451.

Laio, F., A. Porporato, R. Revelli, and L. Ridolfi (2003), A comparison of nonlinear flood forecasting methods, *Water Resour. Res.*, 39(5), 1129, doi:10.1029/2002WR001551.

Lin, J.Y., C.T. Cheng, and K.W. Chau (2006), Using support vector machines for long-term discharge prediction, *Hydrological Sciences–Journal*, 51(4), 599-611.

Liong, S.Y., and C. Sivapragasam (2002), Flood stage forecasting with support vector machines, *Journal of American Water Resour*, 38(1),173 -186.

Lisi, F., O. Nicolis, and M. Sandri (1995), Combining singular-spectrum analysis and neural networks for time series forecasting, *Neural Processing Letters*, 2(4), 6-10.

María, C.M., G.M. Wenceslao, F.B. Manuel, M.P.S. José, and L.C. Román (2004), Modelling of the monthly and daily behaviour of the discharge of the Xallas river using Box–Jenkins and neural networks methods, *Journal of Hydrology*, 296,38–58.

Marques, C.A.F., J. Ferreira, A. Rocha, J. Castanheira, P. Gonçalves, N. Vaz., and J.M. Dias (2006), Singular spectral analysis and forecasting of hydrological time series, *Physics and Chemistry of the Earth*, 31,1172-1179.

Minns, A. W., and M. J. Hall (1996), Artificial neural networks as rainfall-runoff models, *Hydrol. Sci.*, 41(3), 399–417, 1996.

Muttil, N. and K.W. Chau (2006), Neural network and genetic programming for modelling coastal algal blooms. *Int. J. Environment and Pollution*, Vol. 28, Nos. 3/4, pp.223–238.

Nash, J. E. and J. V. Sutcliffe (1970), River flow forecasting through conceptual models; part I – a discussion of principles, *Journal of Hydrology*, 10, 282–290, 1970.

Pasternack, G.B. (1999). Does the river run wild? Assessing chaos in hydrological systems. *Adv. Water Resour*. 23, 253–260.

Procaccia, I. (1988). Complex or just complicated? *Nature*, 333, 498–499.

Raman, H., and N. Sunilkumar (1995), Multivariate modeling of water resources time series using artificial neural networks, *Hydrological Sciences Journal*, 40(2):145-163.

Ruelle, D. (1990), Proc. R. Soc. London, Ser. A 427,241.

Salas, J.D., G.Q. Tabios III, and P. Bartolini (1985), Approaches to multivariate modeling of water resources time series, *Water Resources Bulletin*, 21(4), 683-708.

Sauer, T., J. A. Yorke, and M. Casdagli (1991), Embedology, *J. Stat. Phys.*, 65, 579– 616.

Schertzer, D., I. Tchiguirinskaia, S. Lovejoy, P. Hubert, H. Bendjoudi, and M. Larchevêque (2002), Which chaos in the rainfall runoff process? A discussion on Evidence of chaos in the rainfall-runoff process by Sivakumar, *Hydrological Sciences-Journal,* 47 (1), 139–147.

See, L., and Openshaw, S. (2000). A hybrid multi-model approach to river level forecasting. *Hydrological Sciences Journal*, 45 (3), 523–536.

Shrestha, D. L. & Solomatine, D. P. (2006). Experiments with AdaBoostRT, an improved boosting scheme for regression. *Neural Comput.*, 17.

Sivakumar, B., A. W. Jayawardena, and T. M. K. Fernando (2002), River flow forecasting: use of phase-space reconstruction and artificial neural networks approaches, *Journal of Hydrology*, 265(1), 225-245.

Sivapragasam, C., and S. Y. Liong (2005). Flow categorization model for improving forecasting. *Nordic Hydrology*, 36 (1), 37–48.

Sivapragasam, C., S.Y. Liong, and M.F.K. Pasha (2001), Rainfall and discharge forecasting with SSA-SVM approach, *Journal of Hydroinformatics*, 3(7), 141–152.

Solomatine, D. P., & Ostfeld, A. (2008). Data-driven modelling: Some past experiences and new approaches. *Journal of Hydroinformatics*, 10(1), 3-22.

Solomatine, D. P. and Xue, Y. I. (2004). M5 model trees and neural networks: application to flood forecasting in the upper reach of the Huai River in China. *Journal of Hydrological Engineering*, 9 (6), 491–501.

Sudheer, K. P., P. C. Nayak, and K. S. Ramasastri (2003), Improving Peak Flow Estimates in Artificial Neural Network River Flow Models, *Hydrological Processes*, 17, 677-686.

Sugihara, G and R.M. May (1990), Nonlinear forecasting as a way of distinguishing chaos from measurement error in time series, *Nature* 344, 734– 741.

Takens, F. (1981), Dynamical systems and Turbulence, *Lecture note in Mathematics*, Vol. 898, Springer, New York.

Theiler, J., (1986), Spurious dimension from correlation algorithms applied to limited time-series data. *Phys. Rev. A* 34 (3), 2427–2432.

Thirumalaiah, K., and M. C. Deo (2000). Hydrological forecasting using neural networks. *Journalof Hydrologic Engineering*, Vol. 5, No. 2, 180-189.

Vautard, R., P. Yiou, and M. Ghil (1992), Singular-spectrum analysis: a toolkit for short, noisy and chaotic signals, **Physica D** 58, 95–126.

Wang,W., P.H.A.J.M., Van Gelder, J.K. Vrijling, and J. Ma (2006a). Testing for nonlinearity of streamflow processes at different timescales, *Journal of Hydrology*, 322, 247-268.

Wang, W., P.H.A.J.M. van Gelder, J.K. Vrijling, and J. Ma (2006b), Forecasting Daily Streamflow Using Hybrid ANN Models, *Journal of Hydrology*, 324, 383-399.

Wu, C.L., K.W. Chau, and Y.S. Li (2008), River stage prediction based on a distributed support vector regression, *Journal of Hydrology*, 358, 96-111.

Yu, P.S., S.T. Chen, and I.F. Chang (2006), Support vector regression for real-time flood stage forecasting, *Journal of hydrology*, 328,704-716.

Yu, X.Y., S.Y. Liong, and V. Babovic (2004), EC-SVM approach for real-time hydrologic forecasting, *Journal of Hydroinformatics*, 6(3), 209-233.

Yu, P.S., and T.Y. Tseng, (1996), A model to forecast flow with uncertainty analysis, *Hydrological Sciences-Journal*, 41(3), 327-344.

Zhang, B., and R. S. Govindaraju (2000), Prediction of watershed runoff using Bayesian concepts and modular neural networks. *Water Resources Research*, 36(3), 753-762.
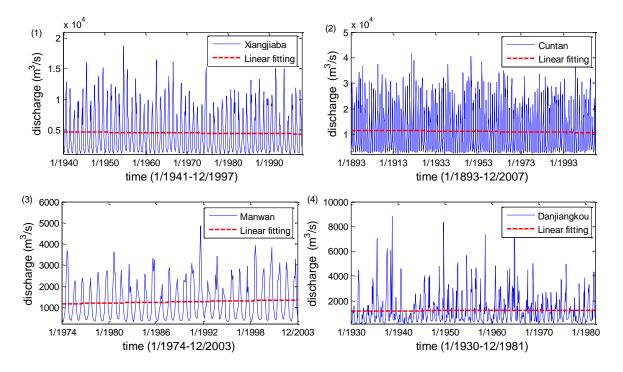
1



2

3    Figure 1. Monthly discharge series of (1) Xiangjiaba, (2) Cuntan, (3) Manwan, and (4) Danjiangkou

4



5

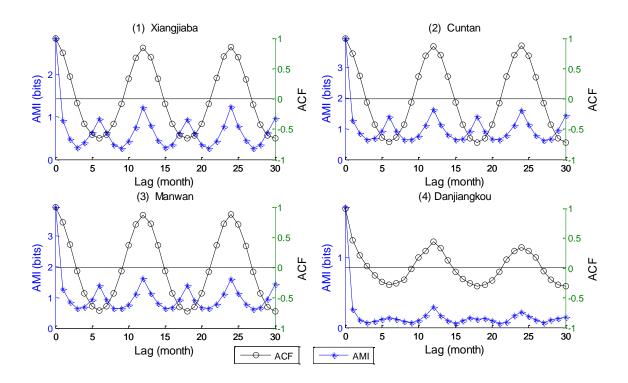1        Figure 2. ACF and AMI of (1) Xiangjiaba, (2) Cuntan, (3) Manwan, and (4) Danjiangkou
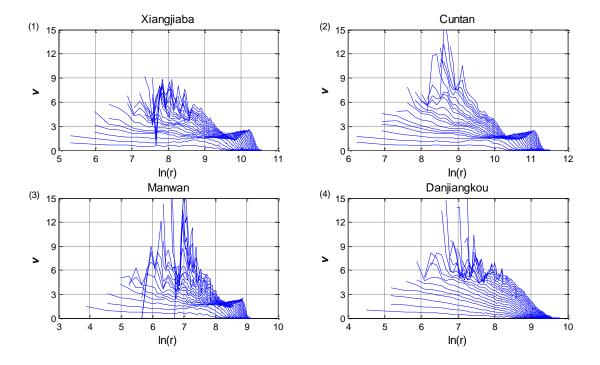


2

3    Figure 3. The estimation of correlation dimension ($d_2$) for (1) Xiangjiaba, (2) Cuntan (3) Manwan, and (4)

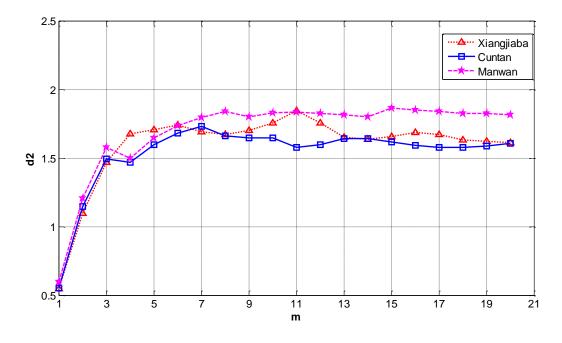4    Danjiangkou. $\tau = 3$ and $m$ is at the interval of [1,20] (increasing from bottom to top in each pane).



5

6    Figure 4. Relationship between $d_2$ and $m$ for monthly discharges of Xiangjiaba, Cuntan, and Manwan

1



2

3    Figure 5. Phase portraits for (1) Xiangjiaba, (2) Cuntan (3) Manwan, and (4) Danjiangkou when $m = 3$ and $\tau = 3$.

4



5

6    Figure 6. FNNP for Xiangjiaba, Cuntan, Manwan, and Danjiangkou when $R_{tol} = 15$ and $\tau = 3$.

7

1

2    Figure 7. Singular Spectrum for (1) Xiangjiaba, (2) Cuntan (3) Manwan, and (4) Danjiangkou with different m



3

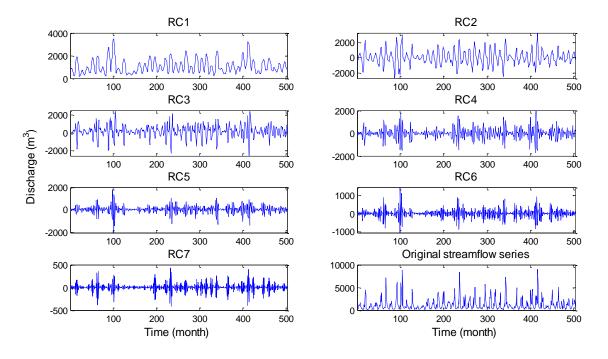4    Figure 8. Singular Spectrum for (1) Xiangjiaba, (2) Cuntan (3) Manwan, and (4) Danjiangkou with different $\tau$

Figure 9. Reconstructed components (RCs) and original streamflow series of Danjiangkou



Figure 10. AMI and CCF between the original flow series and each RC (Danjiangkou)
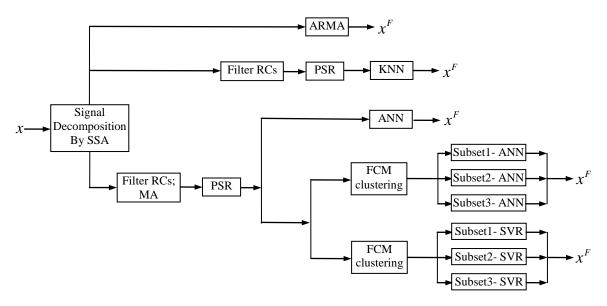
1

2  Figure 11. Framework of modelling implementation for monthly streamflow predictions ( $x^F$ stands for forecasts)
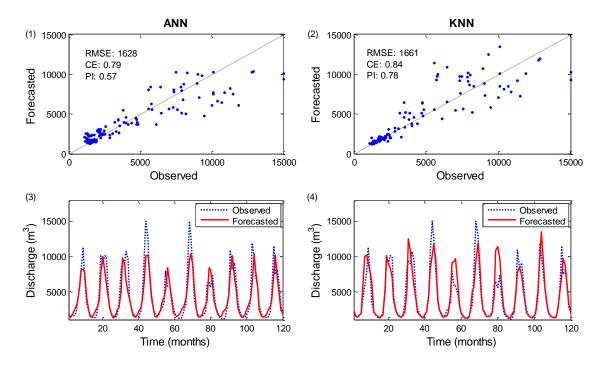


3

4  Figure 12. Observed and forecasted Xiangjiaba's discharges for one-month-ahead forecast (ANN modeling results

5      were presented in (1) and (3), and KNN modeling results were presented (2) and (4))
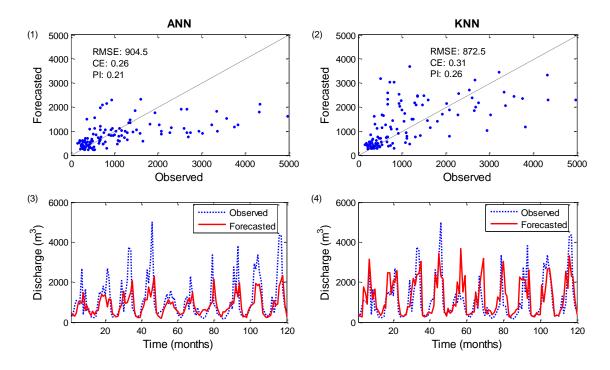
**ANN**

**KNN**

(1) RMSE: 904.5
CE: 0.26
PI: 0.21

(2) RMSE: 872.5
CE: 0.31
PI: 0.26

(3)

(4)

1

2   Figure 13. Observed and forecasted Danjiangkou's discharges for one-month-ahead forecast (ANN modeling

3       results were presented in (1) and (3), and KNN modeling results were presented (2) and (4))
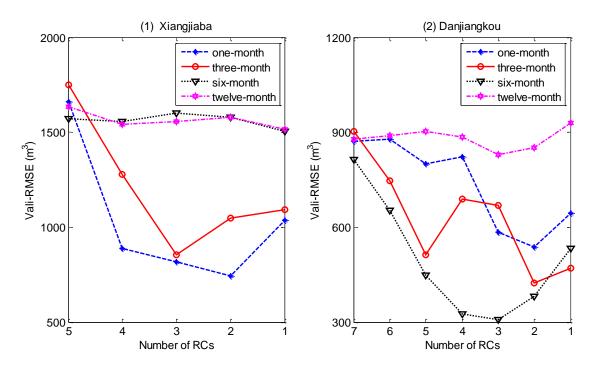
(1) Xiangjiaba

(2) Danjiangkou

4

5   Figure 14. Performance curves of different prediction horizons by KNN model with various RCs as inputs
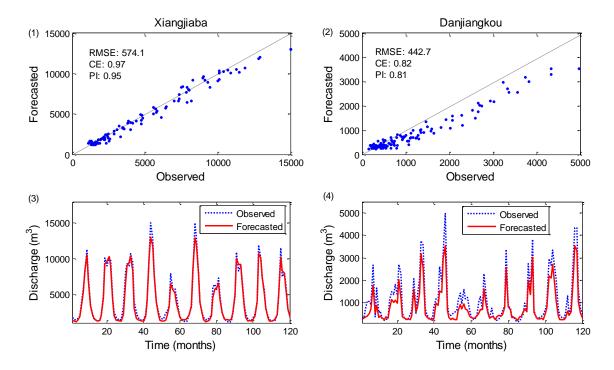
6

1

2    Figure 15. Observed and forecasted discharges for one-month-ahead forecast from ANN model with the data

3          preprocessing of MA and SSA ((1) and (3) for Xiangjiaba; (2) and (4) for Danjiangkou)
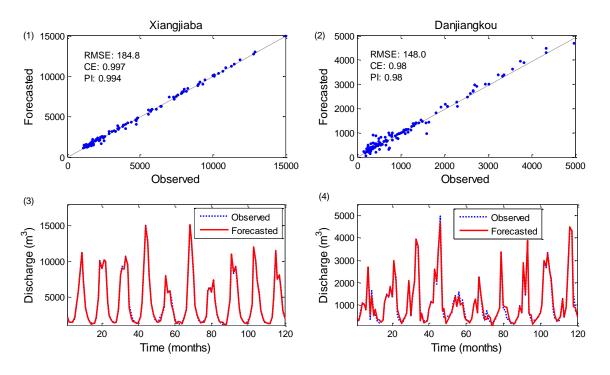


4

5    Figure 16. Observed and forecasted discharges for one-month-ahead forecast from CDSVR model with the data

6          preprocessing of MA and SSA ((1) and (3) for Xiangjiaba; (2) and (4) for Danjiangkou)
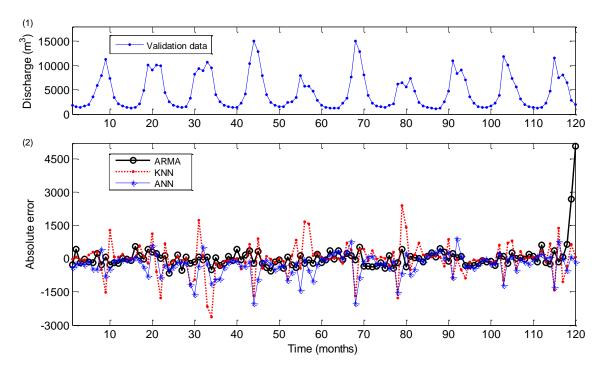
1



2

3       Figure 17a. Forecasting results with MA and SSA for Xiangjiaba (1) validation data (2) absolute error (=

4                              forecast - target) of forecasts from ARMA, KNN and ANN



5

1        Figure 17b. Forecasting results with MA and SSA for Xiangjiaba (1) validation data (2) absolute error (=

2                    forecast - target) of forecasts from ANN, CDANN and CDSVR

3

4

5        Figure 18a. Forecasting results with MA and SSA for Danjiangkou (1) validation data (2) absolute error (=

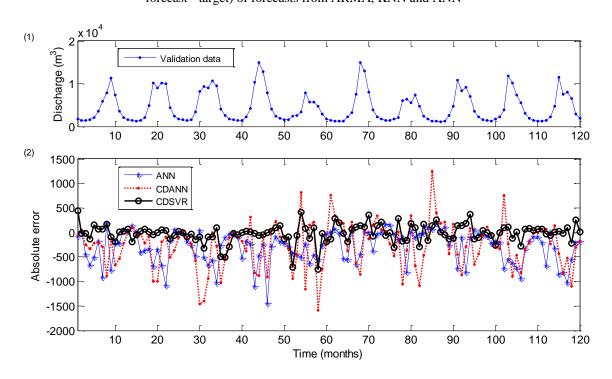6                    forecast - target) of forecasts from ARMA, KNN and ANN

1

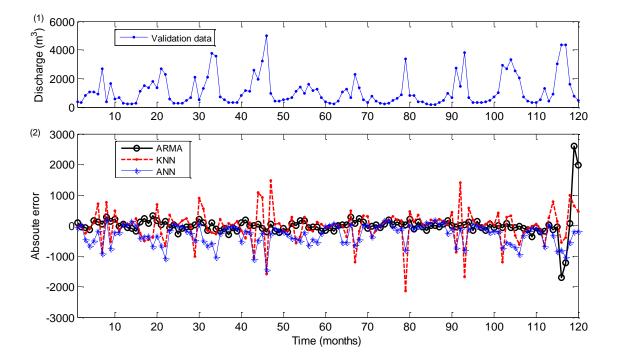2    Figure 18b. Forecasting results with MA and SSA for Danjiankou (1) validation data (2) absolute error (=

3    forecast - target) of forecasts from ANN, CDANN and CDSVR

4



5

1     Figure 19. Histograms of absolute prediction errors of (1) KNN, (2) ANN, (3) CDANN, and (4) CDSVR for

2                                Danjiangkou



3

4          Figure 20. The check of robustness of $k$ in KNN method for (1) Xiangjiaba and (2) Danjiangkou



5

1    Figure 21. CCFs between inputs and output from ANN with/without MA for one-month-ahead forecast: (1)

2    Cuntan and (2) Danjiangkou.



3

4    Figure 22. CCFs between inputs and outputs from ANN with/without MA and SSA for six-month-ahead forecast:

5    (1) Cuntan and (2) Danjiangkou.

6

7

8

9

10

11

12

13

14

15

16

17

1        Table 1 One-month-ahead forecasting results using various models without data
2        preprocessing of MA and/or SSA

| Watershed | Model | RMSE | CE | PI | Model parameters | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | (P,Q) | $(\tau, m, k)$ | (I-h-O)[*a] | $(C, \varepsilon, \sigma)$ |
| Xiangjiaba | | | | | | | | |
| | ARMA | 1293.4 | 0.87 | 0.73 | (12,11) | | | |
| | KNN | 1661.2 | 0.78 | 0.55 | | (3,4,5) | | |
| | ANN | 1627.8 | 0.79 | 0.57 | | | (4-8-1) | |
| | CDANN | 1427.1 | 0.84 | 0.70 | | | (4-8/8/8-1) | |
| | CDSVR | 1407.1 | 0.84 | 0.68 | | | | 3*[b] |
| Danjiangkou | | | | | | | | |
| | ARMA | 762.3 | 0.47 | 0.44 | (11,10) | | | |
| | KNN | 872.5 | 0.31 | 0.26 | | (3,5,6) | | |
| | ANN | 904.5 | 0.26 | 0.21 | | | (5-5-1) | |
| | CDANN | 895.3 | 0.27 | 0.42 | | | (5-5/5/5-1) | |
| | CDSVR | 690.3 | 0.57 | 0.54 | | | | 3*[c] |

3   [*a] (I-h-O) stands for the number of neurons in the input layer, hidden layer and output layer;
4   3*[b] three triplet parameters $(C, \varepsilon, \sigma)$ for Xiangjiaba are (46.96, 0.0056, 0.3278), (737.66, 0.0062, 0.8065),
5   and (0.66128, 0.0248, 0.4205);
6   3*[c] three triplet parameters $(C, \varepsilon, \sigma)$ for Danjiangkou are (773.95, 0.0002, 0.8264), (3.75, 0.0092, 0.1124),
7   and (3.99, 0.0068, 0.5342).

8        Table 2 One-month-ahead forecasting results using various models with data
9        preprocessing of MA and/or SSA

| Watershed | Model | RMSE | CE | PI | Model parameters | | | | Data-preprocessing parameters | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | $(P_i, Q_i)$ | $(\tau, m, k)$ | (I-h-O) | $(C, \varepsilon, \sigma)$ | Memory length of MA (months) | $p$ from m RCs of SSA |
| Xiangjiaba | | | | | | | | | | |
| | ARMA | 586.2 | 0.97 | 0.94 | *[a] | | | | 1 | 5/5 |
| | KNN | 744.4 | 0.96 | 0.91 | | (3,4,5) | | | 1 | 2/5 |
| | ANN | 574.1 | 0.97 | 0.95 | | | (4-8-1) | | 6 | 5/5 |
| | CDANN | 540.8 | 0.98 | 0.95 | | | (4-8/8/8-1) | | 6 | 5/5 |
| | CDSVR | 184.8 | 1.00 | 0.99 | | | | 3*[c] | 6 | 5/5 |
| Danjiangkou | | | | | | | | | | |
| | ARMA | 376.6 | 0.87 | 0.86 | *[b] | | | | 1 | 7/7 |
| | KNN | 538.1 | 0.74 | 0.72 | | (3,5,6) | | | 1 | 2/7 |
| | ANN | 442.7 | 0.82 | 0.81 | | | (5-5-1) | | 3 | 7/7 |
| | CDANN | 298.7 | 0.92 | 0.82 | | | (5-5/5/5-1) | | 3 | 7/7 |
| | CDSVR | 148.0 | 0.98 | 0.98 | | | | 3*[d] | 3 | 7/7 |

10   [*a] $(P_i, Q_i)$ denotes optimal parameter pairs for each RC, i = 1,···,5 for Xiangjiaba, and $P_i$ and $Q_i$ are
11   respectively (10,9,10,11,9) and (12,12,12,12,12);
12   [*b] $(P_i, Q_i)$ denotes optimal parameter pairs for each RC, i = 1,···,7 for Danjiangkou, and $P_i$ and $Q_i$ are
13   respectively (5,5,10,8,7,10,5) and (10,10,11,10,12,10,12);

3*c three triplet parameters ($C, \varepsilon, \sigma$) for Xiangjiaba are (454.51, 0.0031, 0.5911), (53.26, 0.0082, 0.3299), and (731.14, 0.0003, 0.7788);

 3*d three triplet parameters ($C, \varepsilon, \sigma$) for Danjiangkou are (485.64, 0.0023, 0.4321), (0.38, 0.0007, 0.2839), and (805.86, 0.0006, 8.4053).

Table 3 Three-month-ahead forecasting results using various models with data preprocessing of MA and/or SSA

| Watershed | Model | RMSE | CE | PI | Model parameters | | | | Data-preprocessing parameters | |
| | | | | | $(P_i, Q_i)$ | $(\tau, m, k)$ | (I-h-O) | $(C, \varepsilon, \sigma)$ | Memory length of MA (months) | $p$ from m RCs of SSA |
|---|---|---|---|---|---|---|---|---|---|---|
| Xiangjiaba | | | | | | | | | | |
| | ARMA | 606.01 | 0.97 | 0.99 | *a | | | | 1 | 5/5 |
| | KNN | 853.87 | 0.94 | 0.97 | | (3,4,5) | | | 1 | 3/5 |
| | ANN | 795.83 | 0.95 | 0.98 | | | (4-8-1) | | 8 | 5/5 |
| | CDANN | 601.85 | 0.97 | 0.95 | | | (4-8/8-1) | | 8 | 5/5 |
| | CDSVR | 435.05 | 0.98 | 0.97 | | | | 3*c | 8 | 5/5 |
| Danjiangkou | | | | | | | | | | |
| | ARMA | 450.13 | 0.82 | 0.91 | *b | | | | 1 | 7/7 |
| | KNN | 423.42 | 0.84 | 0.92 | | (3,5,6) | | | 1 | 2/7 |
| | ANN | 533.49 | 0.74 | 0.87 | | | (5-5-1) | | 5 | 4/7 |
| | CDANN | 407.85 | 0.85 | 0.69 | | | (5-5/3/2-1) | | 5 | 4/7 |
| | CDSVR | 305.64 | 0.92 | 0.91 | | | | 6*d | 5 | 4/7 |

*a ($P_i, Q_i$) are the same as one-month-ahead forecast;

*b ($P_i, Q_i$) are the same as one-month-ahead forecast;

3*c three triplet parameters ($C, \varepsilon, \sigma$) for Xiangjiaba are (437.77, 0.0006, 0.3696), (10.01, 0.0050, 0.32593), and (36.25, 0.0053, 0.3683);

6*d six triplet parameters ($C, \varepsilon, \sigma$) for Danjiangkou are (916.01, 0.0025, 2.579), (34.01, 0.0003, 8.5518), (221.38, 0.0004, 4.1752), (0.56, 0.0006, 0.2186), (64.41, 0.0005, 4.9333), and (149.37, 0.0024, 0.3234).

1        Table 4 Six-month-ahead forecasting results using various models with data
2        preprocessing of MA and/or SSA

| Watershed | Model | RMSE | CE | PI | Model parameters | | | | | Data-preprocessing parameters | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | $(P_i,Q_i)$ | $(\tau,m,k)$ | (I-h-O) | $(C,\varepsilon,\sigma)$ | | Memory length of MA (months) | $p$ from m RCs of SSA |
| Xiangjiaba | | | | | | | | | | | |
| | ARMA | 1519.30 | 0.82 | 0.94 | *[a] | | | | | 1 | 5/5 |
| | KNN | 1503.00 | 0.82 | 0.94 | | (3,4,5) | | | | 1 | 1/5 |
| | ANN | 1535.70 | 0.81 | 0.94 | | | (4-8-1) | | | 6 | 4/5 |
| | CDANN | 1193.60 | 0.89 | 0.81 | | | (4-8/8/6-1) | | | 6 | 4/5 |
| | CDSVR | 1137.50 | 0.90 | 0.79 | | | | 3*[c] | | 6 | 4/5 |
| Danjiangkou | | | | | | | | | | | |
| | ARMA | 326.04 | 0.90 | 0.96 | *[b] | | | | | 1 | 7/7 |
| | KNN | 308.82 | 0.91 | 0.97 | | (3,5,6) | | | | 1 | 3/7 |
| | ANN | 495.14 | 0.78 | 0.91 | | | (5-5-1) | | | 6 | 4/7 |
| | CDANN | 414.10 | 0.84 | 0.75 | | | (5-3/6/2-1) | | | 6 | 4/7 |
| | CDSVR | 413.68 | 0.84 | 0.83 | | | | 6*[d] | | 6 | 4/7 |

3    *[a] $(P_i,Q_i)$ are the same as one-month-ahead forecast;
4    *[b] $(P_i,Q_i)$ are the same as one-month-ahead forecast;
5    3*[c] three triplet parameters $(C,\varepsilon,\sigma)$ for Xiangjiaba are (437.77, 0.0006, 0.3696), (10.01, 0.0050,
6    0.32593), and (36.25, 0.0053, 0.3683);
7    6*[d] six triplet parameters $(C,\varepsilon,\sigma)$ for Danjiangkou are (916.01, 0.0025, 2.579), (34.01, 0.0003, 8.5518),
8    (221.38, 0.0004, 4.1752), (0.56, 0.0006, 0.2186), (64.41, 0.0005, 4.9333), and (149.37, 0.0024, 0.3234).

9        Table 5 Twelve-month-ahead forecasting results using various models with data
10       preprocessing of MA and/or SSA

| Watershed | Model | RMSE | CE | PI | Model parameters | | | | | Data-preprocessing parameters | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | $(P_i,Q_i)$ | $(\tau,m,k)$ | (I-h-O) | $(C,\varepsilon,\sigma)$ | | Memory length of MA (months) | $p$ from m RCs of SSA |
| Xiangjiaba | | | | | | | | | | | |
| | ARMA | 2157.10 | 0.63 | 0.09 | *[a] | | | | | 1 | 5/5 |
| | KNN | 1516.20 | 0.82 | 0.55 | | (3,4,5) | | | | 1 | 1/5 |
| | ANN | 1657.20 | 0.78 | 0.46 | | | (4-8-1) | | | 1 | 5/5 |
| | CDANN | 1440.60 | 0.83 | 0.86 | | | (4-8/8/8-1) | | | 2 | 5/5 |
| | CDSVR | 1258.70 | 0.87 | 0.74 | | | | 6*[c] | | 4 | 5/5 |
| Danjiangkou | | | | | | | | | | | |
| | ARMA | 848.35 | 0.35 | -0.13 | *[b] | | | | | 1 | 7/7 |
| | KNN | 851.94 | 0.34 | -0.14 | | (3,5,6) | | | | 1 | 3/7 |
| | ANN | 817.83 | 0.39 | -0.05 | | | (5-7-1) | | | 1 | 4/7 |
| | CDANN | 634.35 | 0.63 | 0.64 | | | (5-5/3/2-1) | | | 1 | 4/7 |
| | CDSVR | 627.10 | 0.64 | 0.62 | | | | 6*[d] | | 4 | 4/7 |

1    *[a] ($P_i$,$Q_i$) are the same as one-month-ahead forecast;

2    *[b] ($P_i$,$Q_i$) are the same as one-month-ahead forecast;

3    6*[c] six triplet parameters ($C, \varepsilon, \sigma$) for Xiangjiaba are (0.67,    0.0064,    0.0685), (39.71, 0.0062, 1.3351),

4    (473.02, 0.0028, 0.9397), (42.88, 0.0294, 0.4906), (7.24, 0.0004, 0.5580), and (0.67, 0.0056, 0.0632).

5    6*[d] six triplet parameters ($C, \varepsilon, \sigma$) for Danjiangkou are (342.6, 0.0068, 2.1952), (5.30, 0.0363, 3.7493),

6    (167.4, 0.0837, 7.060), (0.9507, 0.0008, 0.0837), (0.61, 0.0059, 0.2515), and (29.35, 0.0042, 0.9685).