2

# 3    River Stage Prediction Based on a Distributed Support Vector Regression

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

5          Dept. of Civil and Structural Engineering, Hong Kong Polytechnic University,
6                Hung Hom, Kowloon, Hong Kong, People's Republic of China
7                          (*Email: cekwchau@polyu.edu.hk)

8    **Abstract**:
9          An accurate and timely prediction of river flow flooding can provide time for the
10   authorities to take pertinent flood-protection measures such as evacuation. Various data-
11   derived models including LR (linear regression), NNM (the nearest-neighbor method) ANN
12   (artificial neural network) and SVR (support vector regression), have been successfully
13   applied to water level prediction. Of them, SVR is particularly highly valued, because it has
14   the advantage over many data-derived models in overcoming overfitting of training data.
15   However, SVR is computationally time-consuming when used to solve large-size problems.
16   In the context of river flow prediction, equipped with LR model as a benchmark and genetic
17   algorithm-based ANN (ANN-GA) and NNM as counterparts, a novel distributed SVR (D-
18   SVR) model is proposed in the present study. It implements a local approximation to training
19   data because partitioned original training data is independently fitted by each local SVR
20   model.  ANN-GA and LR models are also used to help determine input variables. A two-step
21   GA algorithm is employed to find the optimal triplets ($C, \varepsilon, \sigma$) for D-SVR model. The
22   validation results reveal that the proposed D-SVR model can carry out the river flow
23   prediction better in comparison with others, and dramatically reduce the training time
24   compared with the conventional SVR model. The pivotal factor contributing to the
25   performance of D-SVR may be that it implements a local approximation method and the
26   principle of structural risk minimization.

27

28   **Keywords**: Water level prediction; D-SVR; Input selection; Parameter optimization

## 29   Introduction

30         As one of a number of nonstructural flood protection measures, an accurate and
31   timely prediction of water levels in the station of interest is of great importance in helping the
32   authorities determine whether to take measures and if do which measures would best mitigate
33   potential flood damage. In the last two decades, with the development of software technology,
34   many approaches affiliated to 'black box' techniques including NNM (nearest neighbor
35   method), ANN (artificial neural network), and SVR (support vector regression) have been
36   widely applied to flood prediction.
37         NNM has been reported in the literature to analyze rainfall-runoff and
38   runoff/discharge processes and has been compared with ARX (autoregressive model with
39   exogenous inputs), or ARMAX (autoregressive moving average model with exogenous
40   inputs). NNM yielded satisfactory results (Yakowitz*, 1987; Karlsson and Yakowitz, 1987;*
41   *Galeati, 1990*). The technique was extended to NNLPW (nearest neighbor linear perturbation

model) for rainfall-runoff prediction (*Shamseldin and O'Connor, 1996*). Feature selection is one of the most important aspects of pattern recognition, as used in the nearest neighbor method. In the context of univariate time series such as discharge, the feature vector can consist of several previous values (*Karlsson and Yakowitz, 1987; Galeati, 1990*).

Since the renaissance of ANNs in the late of 1980s, they have become the preferred prediction approach for many researchers and have been applied to a variety of issues. While some researchers in the literature employed ANNs alone for river flow forecasts (*Prochazka, 1997; Thirumalaiah and Deo, 1998; Sheta and El-Sherif, 1999; Liong et al.,2000; Salas et al., 2000; Qin et al., 2002; Cannon and Whitfield, 2002; Li and Gu, 2003; Huang et al., 2004; Cheng et al., 2005; García-Pedrajas et al., 2006*), many other researchers compared ANNs with traditional statistical techniques for river flow flood predictions. Comparisons between ANNs and AR (autoregressive) approaches appeared in the work of Raman and Sunilkumar (*1995*), Elshorbagy and Simonovic (*2000*), Thirumalaiah and Deo (*2000*) and Kişi (*2003*). Likewise, some studies were focused on comparisons between ANNs and ARMA (*Jain et al., 1999; Abrahart and See, 2000; Castellano-Me´ndeza et al., 2004*). The majority of studies have proven that ANNs are able to outperform traditional statistical techniques. Further, the superiority of ANNs over nonlinear regression in predicting river flows has been attributed to the possible existence of nonlinear dynamics, which are not well captured by the regression technique. A hybrid ANN model developed by Wang et al. (*2006*) was used to predict daily stream flow.

SVR, with highly similar structures to ANN, can learn from experimental data. SVR performs structural risk minimization (SRM) that aims at minimizing a bound on the generalization error (*Kecman, 2001*). In this way, it creates a model with a minimized VC-dimension (named after the authors, Vapnik and Chervonenkis), which means good generalization. Since SVR generalization performance does not depend on the dimensionality of input space, it can be used with small data sets. However, ANN is data intensive, and has to cover as many patterns as possible in order to perform well, and the generality of ANN is difficult to control as a result of implementing the empirical risk minimization (ERM) principle. Recently, some applications of SVR have been seen in the prediction of rainfall-runoff process, rainfall, and river flow. For example, Sivapragasam et al. (*2001*) performed one-lead-day rainfall forecasting and runoff forecasting using SVR, in which the input data are pre-processed by singular spectrum analysis, resulting in a high-dimensional input space. Yu et al. (*2004*) proposed a scheme that combined chaos theory and SVM to forecast daily runoff. Bray and Han (*2004*) applied SVM to forecast runoff, focusing on the identification of an appropriate model structure and relevant parameters. Sivapragasam and Liong (*2004*) used the sequential elimination approach to identify the optimal training data set and then performed SVR to forecast the water level. Sivapragasam and Liong (*2005*) divided the flow range into three regions, and employed different SVR models to predict daily flows in high, medium and low regions. Lin et al. (*2006*) presented a SVR model to predict long-term monthly flow discharge series, and a comparison with results of appropriate ARMA and ANN models demonstrated the better performance of SVR. Yu et al. (*2006*) carried out a real-time flood stage forecasting based on SVR in which a hydrological concept of the time of response was employed to identify lags of inputs and a two-step grid search method was used for finding optimal parameters.

However, a major drawback of SVR is that training time tends to increase exponentially with the number of training samples. For example, according to the algorithm

88   presented in this paper below, the time required is about two days for a magnitude of 1000
89   training data whereas it is only 40 minutes for a magnitude of 100 training data. Moreover,
90   using a single model to learn large-size data may well lead to mismatch as there are different
91   noise levels in different input regions (*Cheng et al., 2006b*), which is a normal scenario for
92   those rivers characterized by seasonal flooding.
93   This paper mainly aims at developing a distributed SVR (D-SVR) model with a two-
94   step GA parameter optimization method to carry out a prediction of river flow. In order to
95   evaluate the performance of D-SVR, prediction is also arrived at via linear regression (LR),
96   NNM, and ANN-GA (genetic algorithm-based ANN). As an extension of the previous study
97   (*Chau et al., 2005*), some of the background on LR and ANN-GA will be set aside in the
98   present paper. Thus, the paper is constructed as follows: firstly, the principle of SVR and D-
99   SVR is introduced and following this NNM is briefly described. Secondly, in the section on
100  construction of models, an emphasis is placed to input selection, and parameter $k$ in NNM
101  and parameters ($C, \varepsilon, \sigma$) in D-SVR are optimized. In the results and discussion section,
102  results reveal that D-SVR model outperforms the other three models, but with a larger
103  training time except for the conventional SVR. In the conclusion, it is suggested that
104  nonlinear models may achieve more notable advantages over LR in the case of rainfall-runoff
105  mapping.

106  **SVR and Distributed SVR**
107  Unlike classical adaptation algorithms that work in an $L_1$ or $L_2$ norm and minimize the
108  absolute value of an error or of an error square with ERM, SVR performs SRM (*Kecman,*
109  *2001*). In this way, it creates a model with good generalization. The SRM induction principle
110  and the methodology of SVR are briefly described below (*Gunn, 1998; Dibike et al., 2001;*
111  *Kecman, 2001; Sivapragasam et al., 2001, Liong and Sivapragasam, 2002; Cherkassky and*
112  *Ma, 2004; Yu et al., 2006*).
113  ***Statistical learning theory***
114  We consider here standard regression formulation in general settings for predictive
115  learning. The goal is to estimate an unknown real-valued function in the relationship:
116  $$y = r(X) + \delta \qquad\qquad (1)$$
117  where $\delta$ is independent and identically distributed (i.i.d) zero mean random error (noise), $X$
118  is a multivariate input and $y$ is a scalar output. The estimation is made based on a finite
119  number of samples (training data): ($X_i, y_i$), ($i = 1, \cdots, N$). The training data are i.i.d. samples
120  generated according to some (unknown) joint probability density function
121  $$p(X, y) = p(X)p(y|X) \qquad\qquad (2)$$
122  The unknown function in (1) is the mean of the output conditional probability (aka regression
123  function)
124  $$r(X) = \int y p(y|X) dy \qquad\qquad (3)$$
125  A class of functions $f(X, \omega)$ can be formulated to approximate the relationship between input
126  vector and the output variable, where $\omega$ is the parameter vector of the function. The problem
127  of learning is to select the best function $f(X, \omega_0)$ (learning machine) from $f(X, \omega)$ that can
128  predict the output $y$ as accurately as possible. Generally, the quality of an approximation is

3

129 measured by the loss or discrepancy measure $L(y, f(X, \omega))$. Therefore, the best approximation
130 function is that for which the following expected risk function $R(\omega)$ is as small as possible:

$$R(\omega) = \int L(y, f(X, \omega)) dp(X, y) \qquad (4)$$

132 It is known that the regression function (3) is the one minimizing prediction risk (4) with the
133 familiar squared loss function loss:

$$L(y, f(X, \omega)) = (y - f(X, \omega))^2 \qquad (5)$$

135 Note that the set of functions $f(X, \omega)$, $\omega \in \Lambda$ supported by a learning method may or may not
136 contain the regression function (3). Thus, the problem of regression estimation is the problem
137 of finding the best approximation function that minimizes the prediction risk function

$$R(\omega) = \int (y - f(X, \omega))^2 dp(X, y) \qquad (6)$$

139 using only the training data. This risk function measures the accuracy of the learning
140 method's predictions of unknown target function $r(X)$.

141   A difficulty arises in the process of calculating (6) because the probability distribution
142 $p(X, y)$ is unknown. Therefore, it is necessary an induction principle for risk minimization.
143 One such principle is the ERM inductive principle. A straightforward method is to replace
144 the expected risk $R(\omega)$ by the empirical risk $R_{emp}(\omega)$

$$R_{emp}(\omega) = \frac{1}{N} \sum_{i=1}^{N} (y_i - f(X_i, \omega))^2 \qquad (7)$$

146   However, the ERM principle does not guarantee that the function $f_{emp}(X, \omega)$ that
147 minimizes the empirical risk $R_{emp}(\omega)$ converges to the true (or best) function $f(X, \omega_0)$ that
148 minimizes the expected risk $R(\omega)$ when the number of training data is limited, such that the
149 sample is small. In other words, a smaller error on the training set does not necessarily imply
150 higher generalization ability (i.e., a smaller error on an independent test set). To make the
151 most out of a limited amount of data, a novel statistical technique called SRM has been
152 developed (*Vapnik 1995, 1998*). The theory of uniform convergence in probability provides
153 bounds on the deviation of the empirical risk from the expected risk. This theory shows that
154 it is crucial to restrict the class of functions that the learning machine can implement to one
155 with a capacity that is suitable for the amount of available training data.

156   The SRM principle theoretically minimizes the expected risk based on the
157 simultaneous minimization of both the empirical risk and the confidence interval $\Omega$.
158 Therefore, SRM can maintain a trade off between the accuracy of the training data and the
159 capacity of the learning machine so as to improve generalization of the model.

160   For $\omega \in \Lambda$ and $N > h$, a typical uniform VC bound on the expected risk (also called
161 generalization bound $R$), which holds with probability $1 - \eta$, has the following form (*Vapnik,
162 1995, 1998*):

$$R(\omega) \leq R_{emp}(\omega) + \Omega(N, h, \eta) \qquad (8)$$

$$\Omega(N, h, \eta) = \sqrt{\frac{h\left(\log \frac{2N}{h} + 1\right) - \log\left(\frac{\eta}{4}\right)}{N}} \qquad (9)$$

165 The parameter $h$ is called the VC-dimension, and it describes the capacity of a set of
166 functions to represent the data set. The VC dimension is a measure of the model complexity

167 and is often proportional to the number of free parameters in the function $f(X,\omega)$.
168 Particularly when $N/h$ is small, a small empirical risk does not guarantee a small value of the
169 actual risk. In this case, in order to minimize the actual risk $R(\omega)$, one has to minimize the
170 right-hand side of the inequality in (8) simultaneously over both terms. In order to do this,
171 one has to make the VC dimension a controlling parameter. Therefore, the SRM inductive
172 principle is intended to minimize the risk functional with respect to both terms: the empirical
173 risk $R_{emp}(\omega)$ and the confidence interval $\Omega$. The VC confidence term in (8) depends on the
174 chosen class of functions, whereas the empirical risk depends on the one particular function
175 chosen by the training procedure. The objective here is to find that subset of the chosen set of
176 functions, such that the risk bound for that subset is minimized. This is done by introducing a
177 "structure" by dividing the entire class of functions into nested subsets (Fig. 1). SRM then
178 consists of finding that subset of functions which minimizes the bound on the actual risk.
179 This is done by simply training a series of machines, one for each subset, where for a given
180 subset the goal of training is simply to minimize the empirical risk. One then takes that
181 trained machine in the series whose sum of empirical risk and VC confidence is minimal
182 (*Burges, 1998*).

183                                  *Fig. 1 should be put here*
184 ***Nonlinear support vector regression***
185        In the real hydrological world, most issues of interest tend to be nonlinear. A linear
186 SVR is extremely limited. In order to deal with the nonlinearity, the input data, $X$, in input
187 space is mapped to a high dimensional feature space via a nonlinear mapping function, $\phi(X)$.
188 Hence, the underlying function becomes
189 $$f(X,\omega) = \omega \cdot \phi(X) + b \qquad (10)$$
190 Therefore, the objective of the SVR is to find optimal $\omega$, $b$ and some parameters in kernel
191 function $\phi(X)$ so as to construct an approximation function of the underlying function.
192        When introducing Vapnik's $\varepsilon$-insensitivity error or loss function (see Fig. 2), the loss
193 function $L_\varepsilon(y, f(X,\omega))$ on the underlying function can be defined as

194 $$L_\varepsilon\big(y, f(X,\omega)\big) = \big|y - f(X,\omega)\big|_\varepsilon = \begin{cases} 0 & if \ \big|y - \big(\omega \cdot \phi(X) + b\big)\big| \le \varepsilon \\ \big|y - \big(\omega \cdot \phi(X) + b\big)\big| - \varepsilon & otherwise \end{cases} \qquad (11)$$

195 where $y$ represents observed value. Fig. 2 presents the concept of nonlinear SVR,
196 corresponding to Eq. (11). Similar to linear SVR (*Kecman, 2001; Yu et al., 2006*), the
197 nonlinear SVR problem can be expressed as the following optimization problem:

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

198
$$subject \quad to \quad \begin{cases} y_i - f(\phi(X_i),\omega) - b \le \varepsilon + \xi_i \\ f(\phi(X_i),\omega) + b - y_i \le \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \ge 0 \end{cases} \qquad (12)$$

199 where, the term of $\frac{1}{2}\|\omega\|^2$ reflects generalization, and the term of $C\sum_{i=1}^{l}(\xi_i + \xi_i^*)$ stands for
200 empirical risk. The objective in Eq. (12) is to minimize them simultaneously, which

201 implements SRM to avoid underfitting and overfitting the training data. $\xi_i$ and $\xi_i^*$ are slack
202 variables, shown in Fig. 2 for measurements "above" and "below" an $\varepsilon$ tube. Both slack
203 variables are positive values. $C$ is a positive constant that determines the degree of penalized
204 loss when a training error occurs.

205      By introducing a dual set of Lagrange Multipliers, $\alpha_i$ and $\alpha_i^*$, the minimization
206 problem can be solved in a dual space. The objective function in dual form can be
207 represented as (*Gunn, 1998*):

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

208
$$subject\ to \quad \begin{cases} \sum_{i=1}^{N}(\alpha_i - \alpha_i^*) = 0 \\ 0 \le \alpha_i^* \le C, \qquad i = 1,\cdots,N \\ 0 \le \alpha_i \le C, \qquad i = 1,\cdots,N \end{cases} \tag{13}$$

209      There is no fixed guideline how to select an appropriate nonlinear function $\phi(X_i)$.
210 Furthermore, the computation of $\left(\phi(X_i)\cdot\phi(X_j)\right)$ in the feature space may be too complex to
211 perform. An advantage of SVR is that the nonlinear function $\phi(X)$ need not be used. The
212 computation in input space can be performed using a ''kernel'' function
213 $K(X_i, X_j) = \left(\phi(X_i)\cdot\phi(X_j)\right)$ to yield inner products in feature space, avoiding having to
214 perform a mapping $\phi(X)$. In utilizing kernel functions, the key issue is to select admissible
215 kernel functions. The admissible kernel function should be any symmetric function in input
216 space which can represent a scalar product in feature space. The Mercer kernel functions
217 belonging to a set of reproducing kernels (*Vapnik, 1999; Kecman, 2001*) can be proven
218 admissible. Therefore, any functions that satisfy Mercer's theorem can be used as a kernel. A
219 couple of commonly used kernels in SVR include: (1) linear $K(X_i, X_j) = X_i \cdot X_j$ ; (2)
220 polynomial with degree $d$ $K(X_i, X_j) = \left[(X_i \cdot X_j) + 1\right]^d$ ;(3) multilayer perceptron
221 $K(X_i, X_j) = \tanh[(X_i \cdot X_j) + b]$ ;(4) Gaussian RBF $K(X_i, X_j) = \exp(-\frac{\|X_i - X_j\|^2}{2\sigma^2})$ . After
222 obtaining parameters $\alpha_i$, $\alpha_i^*$, and $b_0$, the final approximation function of the underlying
223 function is

$$f(X_i) = \sum_{i=1}^{N}(\alpha_k - \alpha_k^*)K(X_k \cdot X_i) + b_0, k = 1,\cdots,n \tag{14}$$

225 where $X_k$ stands for the support vector, $\alpha_k$ and $\alpha_k^*$ are parameters associated with support
226 vector $X_k$ , $N$ and $n$ represent the number of training samples and support vectors,
227 respectively.
228

229                          *Fig. 2 should be put here*
230 **SVR expressed in matrix notation**

231        The standard quadratic optimization problem for an $\varepsilon$-insensitive function can be
232 expressed in matrix notation as (*Gunn, 1998; Kecman, 2001*)

$$minimize \quad L_d(x) = \frac{1}{2}x^T H x + C^T x \qquad (15)$$

234 where, $H$ is Hessian matrix, $x$ stands for Lagrangian Multipliers. They are expressed as

$$H = \begin{bmatrix} G & -G \\ -G & G \end{bmatrix}, \; C = \begin{bmatrix} \varepsilon - Y \\ \varepsilon + Y \end{bmatrix}, \text{ and } x = \begin{bmatrix} \alpha \\ \alpha* \end{bmatrix}$$

236 with constraints

$$x \cdot (1,\cdots,1,-1,\cdots,-1) = 0,$$

$$\alpha_i, \alpha_i* \geq 0, i = 1,\cdots,l.$$

239 $G$ is an $(l,l)$ matrix with entries $G_{ij} = [X_i^T X_j]$ for a linear regression, and $\alpha = [a_1,\cdots,\alpha_l]$,
240 $\alpha* = [a_1*,\cdots,\alpha_l*]$, $\varepsilon - Y = [\varepsilon - y_1,\cdots,\varepsilon - y_l]$, $\varepsilon + Y = [\varepsilon + y_1,\cdots,\varepsilon + y_l]$. (Note that $G_{ij}$, as given
241 above, is a badly conditioned matrix and we rather use $G_{ij} = [X_i^T X_j + 1]$ instead).
242 In the case of the nonlinear regression, the learning problem is again formulated as the
243 maximization of a dual Lagrangian (15). A similar matrix notation as Eq. (15) is expressed.
244 However, $H$ here is with the changed Grammian matrix $G$ that is now given as

$$G = \begin{bmatrix} G_{11} & \cdots & G_{1l} \\ \vdots & G_{ii} & \vdots \\ G_{l1} & \cdots & G_{ll} \end{bmatrix}$$

246 where the entries $G_{ij} = \phi^T(X_i)\phi(X_j) = K(X_i)(X_j), i,j = 1,\cdots,l.$ Based on the above matrix
247 form, a SVR programming is easy to make.
248 ***D-SVR Configuration***

249                                 *Fig. 3 should be put here*
250         A primitive idea of D-SVR is to partition the original training set into a couple of
251 subsets and then generate a local SVR for each subset independently. Further, an appropriate
252 data fusion approach (sometimes called aggregation) is employed to combine local
253 predictions into a hybrid output. Fig. 3 displays the configuration of D-SVR. First of all,
254 fuzzy c-means clustering algorithm is employed to split the original training set into $L$
255 training subsets. In the present study, water level variables are characterized by clear
256 seasonal variability, and so the raw training set is clustered into eight sub regions. Thus, each
257 subset further serves for training $L$ SVRs. For a new input $X$, $L$ outputs ( $Y_i$,i=1,$\cdots$,L ) will
258 be generated by the D-SVR model and are associated with $L$ degrees of membership
259 ( $\mu_i$,i=1,$\cdots$,L ). Degree of membership can be determined via the inverse of Square Euclidean
260 Distance between the new input $X$ and $C_i$ which is the center of $i$-th subset. Calculation is
261 formulated as follows:

262
$$\begin{cases} \mu_i = 1 & if\ d_i = 0, \quad i = 1, \cdots, L \\ \mu_i = \left(\dfrac{1}{d_i}\right)\bigg/\left(\sum_{i=1}^{L} \dfrac{1}{d_i}\right) & otherwise \end{cases} \qquad (16)$$

263    where, $d_i = \left\|X - C_i\right\|^2$ and $\sum_{i=1}^{L} \mu_i = 1, \mu_i \in [0,1]$.

264    After $L$ outputs ($Y_i, i=1,\cdots,L$) and their degrees of membership ($\mu_i, i=1,\cdots,L$) are
265    achieved, the combined output $Y$ is

266
$$Y = \sum_{i=1}^{L} \mu_i Y_i \qquad (17)$$

267    However, we found experimentally that there are some drawbacks in this D-SVR:
268    when training data is partitioned into several independent subsets without any overlapping, a
269    large prediction error occurs. Generally, the error is larger than that obtained by using a SVR
270    model alone. Analysis also found that the SVR is weak at extrapolation. When an input is far
271    from its clustering center, the SVR will generate a weird prediction, usually quite large
272    although associated with a small degree of membership. In view of this, we attempted to
273    make the following improvement. We set the nearest neighboring two training subsets to
274    overlap one input region by one in the entire input space, thus the number of training data in
275    all sub models will be twice that of the original training data. Furthermore, only two
276    maximum degrees of membership are activated to contribute to the combined output $Y$.
277    Therefore, the third box in Fig. 3 addresses this task, where $\mu_j$ ($j = 1, 2$) is the first two
278    maximum degree of membership in $\mu_i$ ($i = 1, \cdots, n$), and $Y_j$ ($j = 1, 2$) are corresponding
279    outputs as listed in the fourth box of Fig. 3. Finally, a combined output for D-SVR model is

280
$$Y = \sum_{j=1}^{2} \left(\mu_j \bigg/ \sum_{j=1}^{2} \mu_j\right) Y_j \qquad (18)$$

281    **Nearest-Neighbor Method (NNM)**
282    The following is a brief review of the NN method (*Galeati, 1990; Shamseldin and*
283    *O'Connor, 1996*). Let $\{X(i), i = 1, N\}$ be a set of rainfall measurements or parameters related
284    to the forecasting process being studied (e.g., temperature, soil saturation, etc.) expressed as
285    $X(i) = (P_i, P_{i-1}, P_{i-2}, \cdots, P_{i-m+1})^T$ where $P$ stands for feature information (various hydro-
286    meteorological factors affecting runoff prediction (*Galeati, 1990; Yakowitz, 1987*), $m$ is the
287    number of feature information contributing to feature vector or the vector dimension and
288    $\{Q(i), i = 1, N\}$ a set of discharges. Here, $X$ and $Q$ may be single or multiple variables. For
289    each feature vector $X(i)$, there is an associated discharge $Q(i)$ observed at the same time
290    instant. Thus, the available historical data may be summarized into a set of pairs of feature
291    vectors $X(i)$ and scalar discharges $Q(i)$, as $\{X(i), Q(i) : i = 1, N\}$, where n is the total number
292    of the data in the whole historical record. Thus, the NN prediction of $Q(N+1)$ is obtained as:

293
$$\hat{Q}(N+1) = \frac{1}{k} \sum_{i \in S(X,N)} Q(i+1) \qquad (19)$$

294    where $S(X,N)$ denotes the indices of $k$, the nearest neighbors to the feature vector $X(N)$. The
295    meaning of "nearest neighbors" has to be interpreted according to the Euclidean distance: if

296     $d(n)$ represents a vector of coordinates $d_1, d_2, \ldots, d_m$, the differences between the current
297     feature vector and past data, the Euclidean distant is defined as:

298
$$\|d\| = \left( \sum_{i=1}^{m} d_i^2 \right)^{1/2} \tag{20}$$

299     Therefore, if $i$ is in $S$ and $j$ is not in $S$, then $\|X(N) - X(i)\| \leq \|X(N) - X(j)\|$. Intuitively

300     speaking, the forecast $\hat{Q}(N+1)$ by the $k$ nearest neighbor method is the sample average of

301     succeeding runoff of the $k$ nearest neighbors in the database.

302         As an example from the work of Karlsson and Yakowitz (*1987*) displayed in Fig.4,

303     for simplicity, it is supposed that the feature vector depends only on three values of past

304     discharges ($m = 3$) i.e. $X(N) = [Q(N), Q(N-1), Q(N-2)]$, and it is assumed that $k = 4$. The

305     NN algorithm searches through all the consecutive triplets of the historical record for the four

306     triplets closest (in a Euclidean sense) to the present feature vector. The predicted discharge is

307     the mean of successive outflows (shown in Fig.4 as circles) from the four closest historical

308     events.

309         Standardization of $X$ and $Q$ is usually necessary because it eliminates the units from

310     components or elements and reduces any differences in the range of values amongst

311     components such as rainfall and discharge with their different units and scales. In order to

312     reflect the relative importance because the more recent measurements in the feature vector

313     generally have a greater weight towards predicted values, the Euclidean distance can be

314     computed as a weighted Euclidean norm, i.e., $\|d\|_w = \left( \sum_{i=1}^{m} w_i \cdot d_i^2 \right)^{1/2}$ where

315     $w = (w_1, w_2, \cdots, w_m)$ is a fixed sequence of positive numbers (weights). In the present study, an

316     equivalent weight is assigned to each dimension in the feature vector because all variables

317     are water levels.

318         Thus, the prediction model is $\hat{Q}(N+1) = \frac{1}{k} \sum_{i \in S(X,N)} Q(i+1)$. In order to reflect the

319     relative contribution to prediction value, each of all $k$ neighbors is set to a weight factor $\omega_i$

320     which is based on the Euclidean distance. The prediction model becomes

321
$$\hat{Q}(N+1) = \frac{1}{k} \sum_{i \in S(X,N)} Q(i+1) \cdot \omega_i \tag{21}$$

322     where, $\omega_i = \|d_i\|^{-2} \Big/ \sum_{i=1}^{k} \|d_i\|^{-2}$. Then, an optimal $k$ has to be determined by calibration.

323     Generally, the data set is divided into two parts: one is used to construct the NN-predictors

324     (constructing patterns); the other is used to calibrate parameters. Objective function

325     optimizing $k$ is set up as $J(k) = \sum (Q(i+1) - \hat{Q}(i+1))^2, i = 1, \cdots N$, where $Q(i+1)$ is observed

326     value.

327                                 *Fig. 4 should be put here*

328 **Construction of Models**

329 *Study Area*

9

330        The channel reach studied is in the middle stream of the Yangtze River, which is the
331 largest river in China. It passes through Wuhan City, which is the capital of the Hubei
332 Province (see Fig. 5). The flow of the Yangtze River is quite unsteady and exhibits a seasonal
333 behavior. The flow is low during the winter months, and peak flow occurs during August and
334 September. A hydrological year is often classified into a flooding period and a nonflooding
335 period, which are from June to October and from November to May, respectively. The water
336 level at the Luo-Shan station can be as low as 17.35 m during the nonflooding period and as
337 high as 31.04 m during the flooding period. The average water levels are 20.8 and 27.1 m
338 during the nonflooding and flooding periods, respectively. The purpose of this study is to
339 predict water levels of the downstream station, Han-Kou, by known water levels of the
340 upstream station, Luo-Shan. The lateral inflow is neglected, because it is very small in
341 comparison with the discharge of the main stream.

342                                        *Fig. 5 should be put here*

### *Data Preparation*

344        A remarkable property of ANNs or SVRs is their ability to handle nonlinear, noise,
345 and non-stationary data. However, with suitable data preparation beforehand, it is possible to
346 improve the performance further (*Maier and Dandy, 2000; Bray and Han, 2004*). Data
347 preparation involves a number of processes such as data collection, data division and data-
348 preprocessing. Here, data division and data standardization belonging to data preprocessing
349 will be covered.

350        Many research papers have discussed data division in the process of application of
351 ANN (*ASCE, 2000; Chau et al., 2005*). Typically, ANNs are unable to extrapolate beyond
352 the range of the data used for training. Consequently, poor forecasts/predictions can be
353 expected when the validation data contains values outside of the range of those used for
354 training. It is also imperative that the training and validation sets are representative of the
355 same population. Often statistical properties (mean, variance, range) from them are compared
356 in order to measure the representatives. The similar data handling can be applied to SVR in
357 order to obtain the same baseline of comparison. Taking the same data splitting way as that
358 in Chau et al. (*2005*), the data are randomly divided into three sets: training, testing, and
359 validation. While 75% of the data are used for training, 25% are used for validation. The
360 training data are further divided into 2/3 for the training set and 1/3 for the testing set.

361        In the present study we extract 1,448 input-output data pairs of the following format
362 from the data record:

363                                   $[X(t-4), X(t-2), X(t), Y(t+1)]$

364 which shows that the water level of Y at Han-Kou for the next day can be mapped by water
365 levels of X at Luo-Shan at the present day, two-day ahead and four-day ahead. A detailed
366 description for the mapping format can be found in the section on inputs selection. It was
367 ensured that the data used for training, testing, and validation represents the same population
368 so there is no need to extrapolate beyond the range of their training data. Table 1 shows the
369 statistical parameters, including the mean, standard deviation, minimum, maximum, and
370 range, for the training, testing, and validation sets.

371                                       *Table 1 should be put here*

372        Generally, original data for different variables span different ranges. In order to
373 ensure that all variables receive equal attention during the training process, they should be
374 normalized. In this regard, it is not true for this case as shown in Table 1. However, due to

375 restricted domain of independent variables of transfer functions in ANN and kernel functions
376 in SVR, the raw data normalization is required. Additionally, normalization will improve the
377 condition number of the Hessian in the optimization problem (*Gunn, 1998*). All data are
378 scaled to the interval 0.1–0.9. The advantage of using [0.1, 0.9] rather than [0, 1] is that
379 extreme (high and low) water levels, occurring outside the range of the calibration data, may
380 be accommodated (*Hsu et al., 1995*). The scaling and reserve scaling processes are
381 formulated below:

382
$$X_{norm} = 0.1 + 0.8 \times \left( \frac{X_i - X_{min}}{X_{max} - X_{min}} \right) \quad (22)$$

383
$$Y_{norm} = 0.1 + 0.8 \times \left( \frac{Y_i - Y_{min}}{Y_{max} - Y_{min}} \right) \quad (23)$$

384
$$\hat{Y}_i = Y_{min} + \left( \frac{1.0}{0.8} \right) \times (\hat{Y}_{i,norm} - 0.1) \times (Y_{max} - Y_{min}) \quad (24)$$

385 where $X_{norm}$ and $Y_{norm}$ denote scaled appearance of the raw data $X_i$ and $Y_i$, $\hat{Y}_{i,norm}$ stands for

386 the scaled prediction corresponding to $Y_i$, and $\hat{Y}_i$ is the prediction of $Y_i$ in original scale.

### *Inputs Selection*

388 In model development the selection of appropriate input variables is important since it
389 provides the basic information about the system being modeled. However, determining
390 appropriate inputs is not an easy task. Generally, input determination can be divided into two
391 broad stages (*Bowden et al., 2005*). In the first stage, the objective is to reduce the
392 dimensionality of the original set of inputs, resulting in a set of independent inputs, which are
393 not necessarily related to the model output. As a matter of fact, the addition of unnecessary
394 variables would create a more complex model than is required. Moreover, the complex
395 model is susceptible to overfitting of training data. Therefore, it is imperative that variables
396 are independent of each other as system inputs. This subset of inputs can then be used in the
397 second stage to determine which of these inputs are related in some way to the output.
398 Bowden et al. (*2005)* presented a comprehensive review of approaches on input
399 determination in the water resources and those approaches are broadly classified into five
400 groups. In the present paper, a mixed approach is employed to find optimal inputs.
401 Usually, the number of input variables is not known a priori. A firm understanding of
402 the hydrologic system under consideration plays an important role in the successful
403 implementation of the model. For the present case, the travel time of flood between Luo-
404 Shan and Han-Kou is determined to be about 24 hrs using the Muskingum method. In other
405 words, the flood at Han-Kou has a phase lag of approximately one day with that at Luo-Shan.
406 So X(t) as an input is reasonable. In order to reduce the dimensionality of inputs, an
407 autocorrelation analysis on water levels on Luo-Shan was performed and is shown in Fig. 6.
408 An extreme good autocorrelation exists in water level series and any one input at least in the
409 first ten lags cannot be deleted according to this chart. A linear relation on water levels exists
410 between Luo-Shan and Han-Kou. A stepwise linear model analysis on inputs (Luo-Shan
411 water levels) and output (Han-Kou water level) can help determine optimal inputs from a
412 viewpoint of the linear relationship. Fig. 7 is the result of a stepwise linear model. The
413 optimal linear mapping format between two hydrology stations is with three inputs $X_{10}$, $X_8$,
414 and $X_6$(corresponding to X (t), X (t-2), and X (t-4)) and one output Y (t+1).

415    Obviously, autocorrelation analysis and stepwise linear regression analysis cannot
416 capture any nonlinearity among inputs and between inputs and output. Further, sensitivity
417 analyses (computing the contribution to variance) (*Nord and Jacobsson, 1998*) and weights
418 analyses (*Muttil and Chau, 2006*) on inputs based on ANN are carried out to extract
419 nonlinear information.  Notably, as Nord and Jacobsson (1998) reported in the conclusion of
420 their paper, due to the random starting conditions, important inputs remain changeable. In
421 addition, according to their evaluation criteria, results from both methods on the ranking of
422 inputs are not always consistent during training. An improvement is to employ the ANN-GA
423 model, which is with the architecture 3-3-1 described in the section of results, to obtain the
424 relative optimal initial weights and biases for an ANN model.
425    When ANN is initialized by weights and biases from GA optimization, a more stable
426 ANN model can be achieved and has a good generalization. Due to the fixed initial weights
427 and biases, evaluation results on input importance are steady, but results from two
428 approaches are still inconsistent. Based on the approach of Nord and Jacobsson (1998), $X_9$,
429 $X_{10}$, $X_8$, $X_3$, and $X_6$ are ranked in the top five important inputs. When adding $X_3$, $X_9$
430 respectively to initial linear model based on $X_6$, $X_8$, and $X_{10}$, several models are generated.
431 The performances of these models are listed in Table 2. From the perspective of AIC and
432 RMSE from LR and ANN-GA, choosing $X_6$, $X_8$, and $X_{10}$, i.e. X (t), X (t-2), and X (t-4), as
433 the optimal inputs is tenable. Finally, the optimal linear regression (LR) model is
434    $$Y(t+1) = 1.18X(t) - 0.398X(t-2) + 0.229X(t-4) - 5.08 \qquad (25)$$

435    *Fig. 6 should be put here*

436    *Fig. 7 should be put here*

437    *Table 2 should be put here*
## Parameters Tuning Strategy of D-SVM
439    Obtaining optimal $\alpha_i$ and $\alpha_i^*$ in Eq. (13) depends heavily on these parameters that
440 dominate the nonlinear SVR including the cost constant $C$, the radius of the insensitive
441 tube $\varepsilon$, and the kernel parameters. In the present study, the Gaussian RBF is employed as
442 kernel function. So these parameters consist of a triplet $(C, \varepsilon, \sigma)$, whose components are
443 mutually dependent, and so changing the value of one parameter changes other parameters.
444 Therefore, a simultaneous or global optimization scheme such as GA can be helpful (*Cheng*
445 *et al., 2006a*). Due to lack of any a priori knowledge for their bounds, a two-step GA search
446 algorithm is recommended here, which is inspired by a two-step grid search method (*Hsu et*
447 *al., 2003*). First, a coarse range search was used to achieve the best region of these three-
448 dimensional grids. In the present study, coarse range partitions for $C$ are $[10^{-2}\ 10^0]$, $[10^0\ 10^2]$,
449 $[10^2\ 5.0\times10^2]$, and $[5.0\times10^2\ 10^3]$. Coarse range partitions for $\varepsilon$ are $[10^{-4}\ 10^{-3}]$, $[10^{-3}\ 10^{-2}]$, $[10^{-2}\ 10^{-1}]$, and $[10^{-1}\ 10^0]$, and coarse range partitions for $\sigma$ are $[10^{-3}\ 10^{-2}]$, $[10^{-2}\ 10^{-1}]$, $[10^{-1}\ 10^0]$,
451 and $[10^0\ 10^2]$. There are $4^3$ grids, and one of them is selected as intervals of parameters for
452 the next step. Then, in the second step a further GA search for the triplets $(C, \varepsilon, \sigma)$ will be
453 carried out in the selected intervals.
454    In order to avoid overfitting of training data, testing data and training data were
455 evaluated at the same time according to GA's fitting degree function (i.e., RMSE), and
456 weighted average of their fitting degrees was used as the fitting degree of each population in
457 the process of GA operation.

458    ***Evaluation of Performance***

459        Many measures for model evaluation have been documented in the literature of
460    hydrology application (*Legates and McCabe, 1999*; *Elshorbagy and* Simonovic*, 2000;*
461    Luchetta and Manetti, 2003; *Goswami et al., 2005*). Several conventional measures such as
462    correlation coefficient ($r$ or $R^2$), efficiency coefficient ($E$), index of agreement ($d$), RMSE,
463    and so on, were critically reviewed by Legates and McCabe (*1999*), and the review suggested
464    that correlation coefficient is inappropriate for model evaluation. Legates and McCabe (*1999)*
465    suggested a complete assessment of model performance should include at least one
466    'goodness-of-fit' or relative error measure (e.g., $E$ or $d$ ) and at least one absolute error
467    measure (e.g., RMSE or MAE) with additional supporting information. Herein, two
468    conventional evaluation criteria in hydrology, RMSE (root mean square error) and $E$
469    (efficiency coefficient), are used to measure performances of models based on training data,
470    testing data and validation data.

471    (1) RMSE

472
$$RMSE = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(Y_i - \hat{Y}_i)^2} \qquad (26)$$

473    (2) E

474        Nash and Stucliffe (*1970*) defined the model coefficient of efficiency which ranges
475    from minus infinity to 1.0, with higher values indicating better agreement, as:

476
$$E = 1 - \frac{\sum_{i=1}^{N}(Y_i - \hat{Y}_i)^2}{\sum_{i=1}^{N}(Y_i - \bar{Y})^2} \qquad (27)$$

477

478    where $\hat{Y}_i$=forecast water level, $Y_i$=observed water level, $\bar{Y}$=average observed flow, and
479    $N$=number of observations. RMSE provides a quantitative indication of the model absolute
480    error in terms of the units of the variable, with the characteristic that larger errors receive
481    greater attention than smaller ones. This characteristic can help eliminate approaches with
482    significant errors. However, some studies (*Kachroo and Natale, 1992; Legates and McCabe,*
483    *1999*) have reported that the index $E$ is a rather crude index, being overly sensitive to
484    extreme values, because of the square differences in the definition, while being insensitive to
485    additive and proportional differences between model predictions and observations. This
486    feature will lead to the increasing influence of large floods on the calibrated parameter values
487    and thereby enhance the forecast accuracy of the larger floods. In the present study, however,
488    parameter calibration is not based on $E$, but rather on RMSE.

489    **Results and Discussion**

490    ***Results from NNM***

491        The nearest-neighbor method belongs to typical pattern prediction. A good prediction
492    can be achieved when testing or validation patterns are as similar as possible to those of the
493    training data. In other words, a salient limitation of the NNM is that in no case can a value
494    higher than the historical discharges be predicted. This is a deficiency which would severely
495    limit the generality or even the plausibility of the NNM when used in real time forecasting
496    (*Karlsson and Yakowitz, 1987*). However, for daily management purpose in which the

497 interest is not centered on extreme values, it is viable. Therefore, it is viable for daily water
498 level prediction in the present study.
499       According to the principle of NNM, a key step is to find the optimal $k$ (the number
500 of the nearest neighbors) based on the training data. An optimization process on $k$ is graphed
501 in Fig. 8. The optimal $k$ is 7 with RMSE_tst of 0.234m and RMSE_vali of 0.242m.

502                                    *Fig. 8 should be put here*

503                                    *Fig. 9 should be put here*

504       The upper pane in Fig. 8 displays 362 validation samples and comparison of absolute
505 errors between LR and NNW prediction models is exhibited in the lower pane of Fig. 9.  As a
506 whole, error curves from LR and NNW show the same trend. However, compared with LR,
507 the NNM exhibits larger error amplification at some particular points where local extremum
508 appear on the water level curve. Obviously, the performance of NNM is slightly poorer than
509 that of LR, which seems to be discrepant with the recognized fact that NNM can be superior
510 to some linear models. Two potential aspects can contribute to the present phenomenon: first,
511 the prediction series are highly linear; second, training data is not enough for NNM which
512 make it not be able to efficiently capture these patterns reflecting local extremum points.
513 **Results from ANN-GA**
514       In the present study, the ANN-GA model played dual roles both as a counterpart
515 model and as helping determine inputs for all models. Table 3 shows the process determining
516 optimal architecture of ANN based on a three-layer network assumption. So the main task of
517 this experiment was to find the optimal number of hidden nodes and number of training
518 epochs. Here, a testing set was employed to avoid overfitting of the training set based on the
519 early stop method. These values highlighted by bold and italic typeface in 'Test' column
520 exhibit optimal training epochs for different hidden nodes. Configuration of ANN
521 corresponding to the minimum of them may be relatively optimal. Obviously, the minimum
522 is 0.2285 corresponding to $M$ =3 and epoch=7000.   Further, based on the selected
523 parameters $M$ and epoch, inputs analysis can be performed as shown in the previous section
524 of input selection. Finally, the determined architecture of ANN for the present case is 3-3-1
525 with optimal training epoch of 7000. Corresponding RMSE for training, testing and
526 validation set are 0.213m, 0.223m, and 0.237m as shown in Table 6.

527                                    *Table 3 should be put here*

528                                    *Fig. 10 should be put here*

529       Similar to Fig. 9, Fig. 10 describes prediction error processes from LR and ANN-GA.
530 While ANN-GA does not exhibit a good capturing capacity for local extremum points on the
531 curve of validation samples, it seems to exhibit a better capacity for capturing other points
532 than the LR model. Other than the small size of training samples, an unsteady prediction
533 result can contribute to the poor performance due to the unstable parameter optimization
534 method inherent in ANN although GA can lead to a relatively stable initial weights and
535 biases. In other words, the present ANN may still not an optimal ANN for this case.
536 **Results from D-SVR and Conventional SVR**
537       According to previous partition of original data set, samples in training, testing and
538 validation sets are, respectively, 724, 362, and 362. Experiment showed that computation

539 time may vary from about a couple of seconds to nearly half an hour when the number of
540 samples ranges from 50 up to about 300. The optimization process for $C, \varepsilon, \sigma$ based on GA
541 will have to run hundreds of times, which is extremely time-consuming for large-size training
542 samples. Therefore, the present training data was partitioned into eight subsets with an
543 average size of 181 (724/4=181) samples due to the overlapping between two nearest subsets.
544 When adding testing data to the training set, the sample number employed in using GA to
545 optimize parameters ($C, \varepsilon, \sigma$) for D-SVRs is 2172 in all, i.e., two times as the number of
546 training and testing samples (2172=2×(724+362)). On the other hand, for conventional SVR
547 model, GA is also employed to find optimal triplets ($C, \varepsilon, \sigma$) for training set with the help of
548 testing set to control overfitting. Table 4 displays clustering centers and the size of training
549 and testing data associated with each subset for D-SVR model.

550 *Table 4 should be put here*
551        Based on the two-step GA search approach, the optimal values of triplet parameter
552 ($C, \varepsilon, \sigma$) for each subset are obtained as shown in Table 5. The composite training error
553 (RMSE) is 0.21m with a training time of about 2hrs, and support vectors are 68.5%. Further,
554 the testing error and validation error are 0.209m and 0.211m, as shown in Table 6. As a
555 comparison, the training, testing and validation errors from conventional SVR are
556 respectively 0.213m, 0.216m, and 0.236m, which are larger than those from D-SVR, in
557 particular for the validation error. Meanwhile, the training time in conventional SVR is far
558 larger than that in D-SVR, which is unaccepted for the current one-day-ahead prediction.
559        In addition, Fig. 11 displays the comparison of absolute errors between LR and D-
560 SVR models. Their error curves exhibit similar trend, but D-SVR shows evident better
561 prediction capacity than LR in terms of absolute errors although predictions on local
562 extremum points are still not very good, which may be due to the property of the local
563 approximation performed by D-SVR model.

564 *Table 5 should be put here*

565 *Fig. 11 should be put here*
566 **Comparison among Models and Discussion**
567        Table 6 summarizes performance of different models from RMSE, E of validation
568 data, and training time. In view of its unacceptable training time, conventional SVR model
569 will be put aside in the discussion. Three nonlinear models, NNM, ANN-GA, and D-SVM,
570 show a better performance than that of LR in terms of RMSE of training and testing.
571 However, only D-SVR exhibits a better generalization than LR in terms of RMSE of
572 validation data. The value of E also proves that D-SVR's efficiency is the best. A drawback
573 of D-SVR is computationally time-consuming due to hundreds of times parameters
574 optimization via GA.
575        In order to display the performance from nonlinear models, absolute error curves of
576 them were graphed in Fig. 12. Errors from these curves are with a very similar trend that
577 predictions are underestimated at some points whereas predictions are overestimated at other
578 points such as from 230 to 290 at the X-axis.
579        Although NNM, ANN-GA, and D-SVM are all nonlinear models, they are different
580 in essence. NNM and ANN are generally called nonlinear and non-parameter models unlike
581 LR with its fixed formula form. Therefore, their performance is related to many aspects

582 including raw data quality, suitable data preprocessing, and even the ability of modelers, in
583 particular for ANN. However, different from NNM and ANN-GA, D-SVR does not depend
584 on pattern identification to carry out prediction.  To certain extent, it may be called a
585 parameter model or semi-parameter model which can be uniquely achieved under the SRM
586 principle when the triplet parameters are selected. On the other hand, a fixed prediction result
587 is never expected for ANN model due to the random starting conditions. Moreover, the
588 principle of ERM tends to make ANN and NNM be weak in the aspect of generalization.
589     The D-SVR model performed a nonlinear approximation for each subset. Obviously,
590 a local nonlinear fitting from D-SVR should be better than an empirically global fitting from
591 LR. Therefore, if over-fitting is carefully avoided, it is inevitable that the D-SVR achieves a
592 better prediction in comparison with LR.

593                                 *Table 6 should be put here*

594                                  *Fig. 12 should be put here*

595

## 596 Conclusions and Recommendations

597     As one of nonstructural flood protection measures, the future water level at a
598 downstream station was predicted by the water level series at an upstream station. Equipped
599 with LR model as a benchmark and ANN-GA and NNM as counterparts, a novel D-SVR
600 model was established to carry out the forecast using data collected from water level series at
601 the upstream Luo-Shan station and downstream Han-Kou station. ANN-GA and LR models
602 were also used to help determine input variables. A two-step GA algorithm was employed to
603 optimize the triplet parameters ($C, \varepsilon, \sigma$) for D-SVR model. The validation results revealed
604 that proposed D-SVR model can predict the water level better in comparison with the other
605 models, which may be because it implements a local approximation method and the principle
606 of SRM. However, compared with LR model, NNM and ANN-GA did not exhibit a powerful
607 mapping ability in the present case.
608     Certainly, the conclusion should not be hastily drawn that NNM and ANN are worse.
609 As a matter of fact, studies in the literature have reported that NNM and ANN are very
610 powerful in terms of nonlinear mapping. Associated with small-size training data, the present
611 case is characterized by a highly linear mapping relation, which restricts the power of NNM
612 and ANN. A complicated mapping between rainfall and runoff may be expected to really
613 expose their capabilities, which will be presented in a future study.

## 614 References

615 Abraharty, J. R., and See, L. (2000). Comparing neural network and autoregressive moving
616 average techniques for the provision of continuous river flow forecasts in two contrasting
617 catchments. *Hydrol. Process.*, 14, 2157-2172.
618 ASCE Task Committee on Application of the Artificial Neural Networks in Hydrology.
619 (2000). Artificial neural networks in hydrology I: preliminary concepts. *J. Hydrol. Engng,*
620 *ASCE*, 5(2): 115–123.

621     Bowden, G.J., Dandy, G.C., and Maier, H.R. (2005). Input determination for neural network
622     models in water resources applications: Part 1—background and methodology. *Journal of*
623     *Hydrology*, 301, 75–92.
624     Bray, M. and Han, D. (2004). Identification of support vector machines for runoff modelling.
625     *J. Hydroinf.*, 6(4), 265–280.
626     Burges, C.J.C. (1998). A tutorial on support vector machines for pattern recognition. *Data*
627     *mining and knowledge Discovery*, 2(2).
628     Cannon, A. J. and Whitfield, P. H. (2002) Downscaling recent streamflow conditions in
629     British Columbia, Canada, using ensemble neural network models. *Journal of Hydrology*,
630     259, 136-151.
631     Castellano-Me´ndeza, M., Gonza´lez-Manteigaa, W., Febrero-Bande, M., Manuel Prada-
632     Sa´ncheza, J., and Lozano-Caldero´n, R. (2004) Modeling of the monthly and daily behavior
633     of the runoff of the Xallas river using Box–Jenkins and neural networks methods. *Journal of*
634     *Hydrology*, 296, 38–58
635     Chau, K. W., Wu, C. L., and Li, Y. S. (2005). Comparison of Several Flood Forecasting
636     Models in Yangtze River. *Journal of Hydrologic Engineering*, 10(6), 485-491.
637     Cheng, C.T., Chau, K.W., Sun, Y.G. and Lin, J.Y. (2005). Long-term prediction of
638     discharges in Manwan Reservoir using artificial neural network models. *Lecture Notes in*
639     *Computer Science*, 3498, 1040-1045.
640     Cheng, C.T., Zhao, M.Y., Chau, K.W. and Wu, X.Y., (2006a). Using genetic algorithm and
641     TOPSIS for Xinanjiang model calibration with a single procedure. *Journal of Hydrology*,
642     316(1-4), 129-140.
643     Cheng, J., Qian, J.S., and Guo, Y.N. (2006b). A Distributed Support Vector Machines
644     Architecture for Chaotic Time Series Prediction. *Lecture Notes in Computer Science*, 4232,
645     892-899.
646     Cherkassky, V. and Ma, Y. (2004) Practical selection of SVM parameters and noise
647     estimation for SVM regression. *Neural Networks*, 17(1), 113-126.
648     Dibike, Y. B., Velickov, S., Solomatine, D., and Abbott, M. B. (2001). Model induction with
649     support vector machines: introduction and applications. *Journal of Computing in Civil*
650     *Engineering*, 15(3), 208-216.
651     Elshorbagy, A. and Simonovic, S.P. (2000). Performance evaluation of artificial neural
652     networks for runoff prediction. *J. Hydrologic Engineering*, 5(4), 424-427.
653     Galeati, G. (1990). A comparison of parametric and non-parametric methods for runoff
654     forecasting. *Hydrological Science Journal*, 35(1), 79-84.
655     García-Pedrajas, N., Ortiz-Boyer, D., and Hervás, C. (2006). An alternative approach for
656     neural network evolution with a genetic algorithm: Crossover by combinatorial optimization.
657     *Neural Network*, 19, 514-528.
658     Goswami, M., O'Connor, K.M., Bhattarai, K.P., and Shamseldin, A.Y. (2005). Assessing the
659     performance of eight real-time updating models and procedures for the Brosna River.
660     *Hydrology and Earth System Sciences*, 9(4), 394-411.
661     Gunn, S.R. (1998). Support vector machines for classification and regression. *Image, Speech*
662     *and Intelligent Systems Tech. Rep.*, University of Southampton, U.K.
663     Hsu, C. W., Chang, C.C. and Lin, C. J. (2003). A practical guide to support vector
664     classification. Available at http://www.csie.ntu.edu.tw/cjlin/papers/guide/guide.pdf.
665     [Accessed on 20/6/07].

666  Hsu, K.L., Gupta, H.V., and Sorooshian, S. (1995). Artificial neural network modeling of the
667  rainfall-runoff process. *Water Resources Research*, 31(10), 2517-2530.
668  Huang W.R., Xu B., and Hilton, A.(2004). Forecasting Flows in Apalachicola River Using
669  Neural Networks. *Hydrological Processes*, 18, 2545-2564.
670  Jain, S. K., Das, A., and Drivastava, D. K. (1999). Application of ANN for reservoir inflow
671  prediction and operation. *J. Water. Resour. Plann. Manage*., 125(5), 263–271.
672  Kachroo, R.K., and Natale, L.(1992) Non-liner modelling of the rainfall-runoff
673  transformation. *Journal of Hydrology*, 135, 241-369.
674  Karlsson, M., and Yakowitz, S. (1987). Nearest-neighbour methods for non parametric
675  rainfall runoff forecasting. *Wat. Resour. Res.,* 23(7), 1330-1308.
676  Kecman, V. (2001). Learning and soft computing: support vector machines, neural networks,
677  and fuzzy logic models. **MIT press**, Cambridge, Massachusetts.
678  Kişi, O. (2003). River flow modeling using artificial neural networks. *Journal of Hydrologic
679  Engineering*, 9(1), 60-63.
680  Legates, D. R., and McCabe, Jr, G. J. (1999). Evaluating the use of goodness-of-fit measures
681  in hydrologic and hydroclimatic model validation, *Water Resour. Res.*, 35(1), 233– 241.
682  Li, Y., and Gu, R.R. (2003). Modeling Flow and Sediment Transport in a River System
683  Using an Artificial Neural Network. *Environmental Management*, 31(1), 122-134.
684  Lin, J.Y., Cheng, C.T. and Chau, K.W.  (2006) Using support vector machines for long-term
685  discharge prediction. *Hydrological Sciences–Journal,* 51(4), 599-611.
686  Liong, S.Y., Lim, W.H., Paudyal, G.N.(2000). River Stage Forecasting in Bangladesh:
687  Neural Network Approach. *Journal of Computing in Civil Engineering*, ASCE 14(1), 1-8.
688  Liong, S.Y., and Sivapragasam, C. (2002). Flood stage forecasting with support vector
689  machines**. *Journal of American Water Resour*. 38(1):173 -186.
690  Luchetta, A., and Manetti, S. (2003). A real time hydrological forecasting system using a
691  fuzzy clustering approach. *Computers & Geosciences*, 29(9), 1111-1117.
692  Maier, H. R., and Dandy, G. C. (2000). Neural networks for the prediction and forecasting of
693  water resources variables: a review of modeling issues and applications. *Environmental
694  Modeling and Software*, 15(1): 101-123.
695  Muttil, N. and K.W., Chau.(2006). Neural network and genetic programming for modelling
696  coastal algal blooms. *Int. J. Environment and Pollution*, Vol. 28, Nos. 3/4, pp.223–238.
697  Nash, J. E. and Sutcliffe, J. V. (1970), River flow forecasting through conceptual models part
698  I — A discussion of principles, *Journal of Hydrology*, 10 (3), 282–290.
699  Nord, L.I., and Jacobsson, S.P. (1998). A novel method for examination of the variable
700  contribution to computational neural network models. *Chemometrics and Intelligent
701  Laboratory Systems*, 44(1), 153–160.
702  Prochazka, A. (1997). Neural networks and seasonal time-series prediction. *Artificial Neural
703  Networks, Fifth International Conference on* (Conf. Publ. No. 440), Vol., Iss., pp:36-41.
704  Qin, G.H., Ding, J., and Liu, G.D. (2002) River flow prediction using artif icial neural
705  networks: self-adaptive error back-propagation algorithm. *Advances in Water Science* (in
706  Chinese), 13 (1), 37-41.
707  Raman, H., and Sunilkumar, N. (1995). Multivariate modeling of water resources time series
708  using artificial neural networks. *Hydrological Sciences Journal*, 40(2):145-163.
709  Salas, J.D., Markus, M., and Tokar, A.S. (2000). Stream flow forecasting based on artificial
710  neural networks. In *artificial neural networks in hydrology*, (eds.) Govindaraju, R.S., and
711  Rao, A.R. Water Science and Technology Library.

Shamseldin, A.Y., and O'Connor, K.M. (1996). A nearest neighbour linear perturbation model for river flow forecasting. *Journal of Hydrology*, 179, 353-375.

Sheta, A.F., and El-Sherif, M.S. (1999). Optimal prediction of the Nile River flow using neural networks. *Neural Networks*, **1999**. IJCNN '99. International Joint Conference on, Vol.5,3438-3441.

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

Sivapragasam, C., and Liong, S.Y. (2004). Identifying optimal training data set – a new approach. In: Liong, S.Y., Phoon, K.K., Babovic, V. (Eds.), *Proceedings of the Sixth International Conference on Hydroinformatics*, Singapore, 21–24 June 2004. World Scientific Publishing Co., Singapore.

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

Thirumalaiah, K, and Deo, M.C. (1998). River stage forecasting using artificial neural networks. *Journal of Hydrologic Engineering*, 3 (1), 26 -32.

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

Vapnik, V. (1995). *The nature of statistical learning theory*, Springer, New York.

Vapnik, V. (1998). *Statistical learning theory*, Wiley, New York.

Vapnik, V.N. (1999). An overview of statistical learning theory. **IEEE Transactions on Neural Networks** 10 (5), 988–999.

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

Yakowitz, S. (1987). Nearest-neighbour methods for time series analysis. *Journal of time series analysis*, 8(2), 235-247.

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

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

742

743 **Table 1.** Statistical Parameters for Training, Testing, and Validation Sets

| Model variables and data sets | Statistical parameters | | | | |
|---|---|---|---|---|---|
| | Mean | Standard deviation | Minimum | Maximum | Range |
| $X_{t-4}$ (m) | | | | | |
|    Training set | 23.46 | 3.71 | 17.37 | 30.96 | 13.59 |
|    Testing set | 23.46 | 3.71 | 17.35 | 30.93 | 13.58 |
|    Validation set | 23.46 | 3.71 | 17.37 | 31.04 | 13.67 |
| $X_{t-2}$ (m) | | | | | |
|    Training set | 23.46 | 3.71 | 17.35 | 31.04 | 13.69 |
|    Testing set | 23.46 | 3.71 | 17.39 | 30.96 | 13.57 |
|    Validation set | 23.47 | 3.71 | 17.37 | 30.80 | 13.43 |
| $X_t$ (m) | | | | | |
|    Training set | 23.47 | 3.71 | 17.37 | 30.96 | 13.59 |
|    Testing set | 23.46 | 3.71 | 17.35 | 30.93 | 13.58 |
|    Validation set | 23.46 | 3.71 | 17.37 | 31.04 | 13.67 |
| $Y_{t+1}$ (m) | | | | | |
|    Training set | 18.64 | 3.75 | 12.20 | 25.71 | 13.51 |
|    Testing set | 18.64 | 3.75 | 12.26 | 25.70 | 13.44 |
|    Validation set | 18.64 | 3.75 | 12.21 | 25.69 | 13.48 |

744

745 Table 2 Akaike's Information Criterion (AIC) for Models

| Model | RMSE_trn (m) | RMSE_tst (m) | RMSE_vali (m) | AIC |
|---|---|---|---|---|
| $LR(X_6,X_8,X_{10})$ | 0.2396 | 0.240 | 0.237 | 1.630 |
| $LR(X_6,X_8,X_{10},X_9)$ | 0.2395 | 0.242 | 0.238 | 1.632 |
| $LR(X_6,X_8,X_{10},X_3)$ | 0.2394 | 0.240 | 0.237 | 1.634 |
| $ANN\text{-}GA(X_6,X_8,X_{10})$ | 0.213* | 0.223* | 0.237* | 1.601* |
| $ANN\text{-}GA\ (X_6,X_8,X_{10},X_9)$ | 0.210* | 0.229* | 0.245* | 1.778* |
| $ANN\text{-}GA\ (X_6,X_8,X_{10},X_3)$ | 0.210* | 0.229* | 0.242* | 1.778* |

746 *Average over ten time tests

747

Table 3 RMSE of Train and Test Set with Changing Hidden Nodes ($M$) and Epochs

| $M$ | 2 | | 3 | | 4 | | 5 | | 6 | | 7 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Epoch | Train | Test | Train | Test | Train | Test | Train | Test | Train | Test | Train | Test |
| 1000 | 0.2075 | 0.2326 | 0.2035 | 0.2399 | 0.1992 | 0.2775 | 0.1958 | **0.2364** | 0.1952 | **0.2788** | 0.4705 | 0.4533 |
| 3000 | 0.2075 | **0.2325** | 0.2080 | 0.9157 | 0.1980 | 0.2545 | 0.1937 | 0.7127 | 0.1857 | 0.3041 | 0.1855 | **0.3210** |
| 5000 | 0.6113 | 0.5673 | 0.2035 | 0.2399 | 0.1991 | 0.2573 | 0.1936 | 0.2873 | 0.1891 | 0.4448 | 0.1866 | 0.3679 |
| 7000 | 0.2148 | 0.2362 | 0.2032 | **0.2285** | 0.2008 | **0.2554** | 0.1869 | 0.2598 | 0.1905 | 0.4036 | 0.1902 | 0.4881 |
| 9000 | 0.2075 | 0.2325 | 0.2033 | 0.3080 | 0.1991 | 0.2769 | 0.1964 | 0.3792 | 0.1913 | 0.3381 | 0.1861 | 0.3975 |
| 11000 | 0.2075 | 0.2325 | 0.2021 | 0.2387 | 0.2002 | 0.3402 | 0.1942 | 0.5286 | 0.1869 | 0.2877 | 0.1831 | 0.3218 |
| $M$ | 8 | | 9 | | 10 | | 11 | | 12 | | 13 | |
| 1000 | 0.1871 | 0.3467 | 0.1760 | 1.4003 | 0.1746 | 0.4095 | 0.1711 | 0.4803 | 0.1664 | 0.6671 | 0.1644 | 0.4986 |
| 3000 | 0.1783 | **0.3004** | 0.1701 | **0.3191** | 0.1682 | 0.6884 | 0.1666 | 0.6667 | 0.1574 | **0.3723** | 0.1564 | 0.5020 |
| 5000 | 0.1877 | 2.5142 | 0.1795 | 0.4923 | 0.1761 | 0.9704 | 0.1728 | **0.3499** | 0.1576 | 1.1967 | 1.0416 | 0.9948 |
| 7000 | 0.1823 | 1.0629 | 0.1769 | 0.5054 | 0.1776 | **0.3178** | 0.8488 | 0.8334 | 0.1596 | 0.3657 | 0.1540 | 0.5319 |
| 9000 | 0.1797 | 0.4496 | 0.1806 | 0.5123 | 0.1702 | 0.8181 | 0.1662 | 0.4935 | 0.1676 | 2.2748 | 0.1579 | **0.3514** |
| 11000 | 0.4303 | 0.4220 | 0.1785 | 0.6683 | 0.1695 | 0.5109 | 0.1688 | 0.7671 | 0.1649 | 0.5135 | 0.1577 | 0.3803 |
| $M$ | 14 | | 15 | | 16 | | 17 | | 18 | | 19 | |
| 1000 | 0.1571 | 0.4847 | 1.5085 | 1.4068 | 0.1520 | 0.7772 | 0.1465 | 0.8242 | 0.1477 | 0.8621 | 0.1411 | 1.0257 |
| 3000 | 0.1505 | 0.7819 | 0.1486 | 0.5772 | 0.1544 | 0.9932 | 0.1508 | 0.5340 | 0.1324 | 0.7250 | 0.1327 | 0.7653 |
| 5000 | 0.1587 | 0.9445 | 0.1519 | **0.3877** | 0.1464 | 1.2091 | 0.1413 | 0.5150 | 0.1403 | 0.5975 | 0.1340 | **0.6271** |
| 7000 | 0.1600 | **0.4541** | 0.1459 | 0.6086 | 0.1461 | **0.7584** | 1.0348 | 0.9819 | 0.1316 | **0.5012** | 0.1387 | 0.7769 |
| 9000 | 0.1580 | 0.6617 | 0.1486 | 0.6439 | 0.1436 | 1.0023 | 0.1400 | 1.7086 | 0.1332 | 1.0227 | 0.1417 | 0.9096 |
| 11000 | 0.1552 | 0.5632 | 0.1465 | 1.2319 | 0.1467 | 0.7592 | 0.4745 | **0.4632** | 0.1362 | 0.8169 | 0.1341 | 2.5940 |

750 Note: values in Train and Test columns correspond to their RMSE

Table 4 Characteristics of Subsets Partitioned by FCM

| Subset n | Clustering center | | | | Training & testing data number |
|---|---|---|---|---|---|
| | X(t-4) | X(t-2) | X(t) | Y(t+1) | |
| 1 | 27.9 | 27.9 | 27.9 | 23.1 | 304 |
| 2 | 24.9 | 24.9 | 24.9 | 20.1 | 258 |
| 3 | 26.7 | 26.6 | 26.6 | 21.8 | 252 |
| 4 | 23.4 | 23.4 | 23.4 | 18.7 | 307 |
| 5 | 18.5 | 18.5 | 18.5 | 13.5 | 284 |
| 6 | 21.7 | 21.6 | 21.6 | 16.9 | 240 |
| 7 | 29.5 | 29.5 | 29.5 | 24.5 | 143 |
| 8 | 19.8 | 19.8 | 19.8 | 14.9 | 384 |
| Sum | | | | | 2172 |

755
756

Table 5 Calibration Results of Triplet Parameters ($C, \varepsilon, \sigma$) in D-SVR

| Model | | Triplet Parameters | | | RMSE_trn | Percentage of SVs |
|---|---|---|---|---|---|---|
| | | $C$ | $\varepsilon$ | $\sigma$ | (m) | (%) |
| D-SVR | Submodel1 | 242.83 | 0.0004 | 0.689 | | |
| | Submodel2 | 144.03 | 0.0065 | 0.422 | | |
| | Submodel3 | 80.53 | 0.0031 | 3.297 | | |
| | Submodel4 | 724.31 | 0.0033 | 0.477 | 0.210 | 68.5 |
| | Submodel5 | 239.29 | 0.0062 | 0.066 | | |
| | Submodel6 | 873.04 | 0.0231 | 0.527 | | |
| | Submodel7 | 894.02 | 0.0006 | 0.596 | | |
| | Submodel8 | 137.51 | 0.0035 | 0.922 | | |
| Conventional SVR | | 649.36 | 0.0049 | 0.515 | 0.213 | 68.6 |

757
758
759
760

Table 6 Performances for Different Models

| Model | RMSE_trn (m) | RMSE_tst (m) | RMSE_vali (m) | E_vali | Training time (s) |
|---|---|---|---|---|---|
| LR | 0.234 | 0.240 | 0.237 | 0.9960 | - |
| NNM | - | 0.234 | 0.242 | 0.9961 | 10 |
| ANN-GA | 0.213 | 0.223 | 0.237 | 0.9960 | 53 |
| Conventional SVR | 0.213 | 0.216 | 0.236 | 0.9960 | 153532 |
| D-SVR | 0.210 | 0.209 | 0.211 | 0.9968 | 7110 |

761
762

763



Fig. 1 Bound on Actual Risk Is Sum of Empirical Risk and Confidence Interval
(adapted from Vapnik, *1998*)



Fig. 2 Nonlinear SVR with Vapnik's $\varepsilon$-insensitive loss function

771



Fig. 3 Configuration of D-SVR
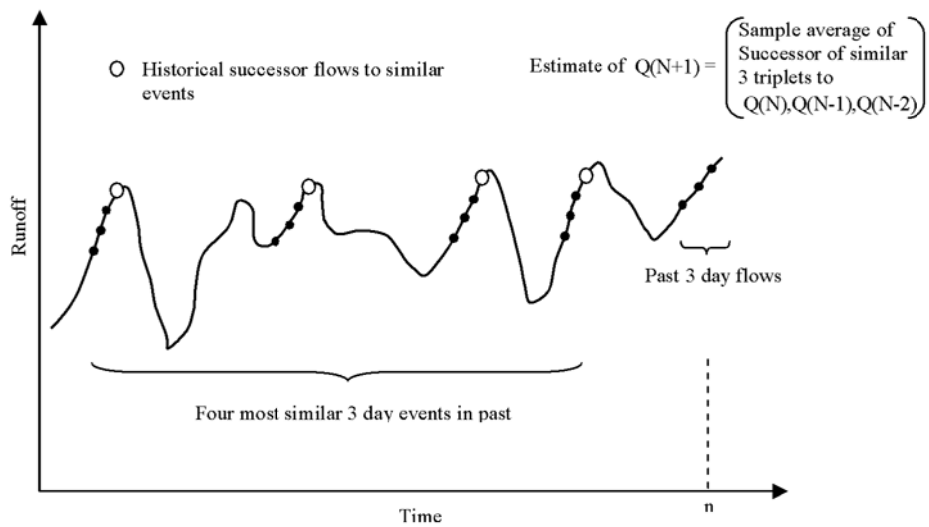
772
773
774



Fig.4 NN rule for runoff case (adapted from *Karlsson and Yakowitz, 1987*)
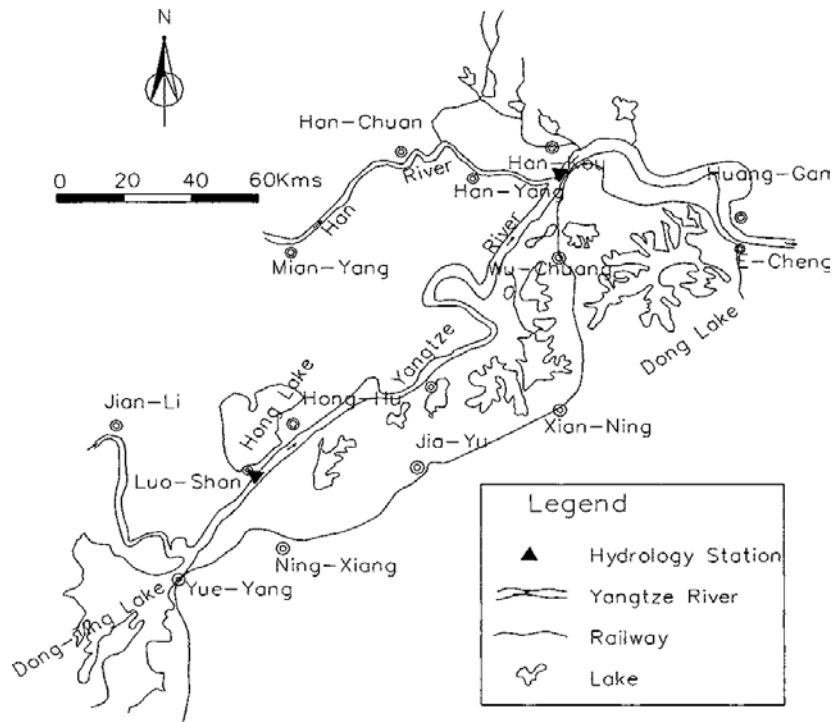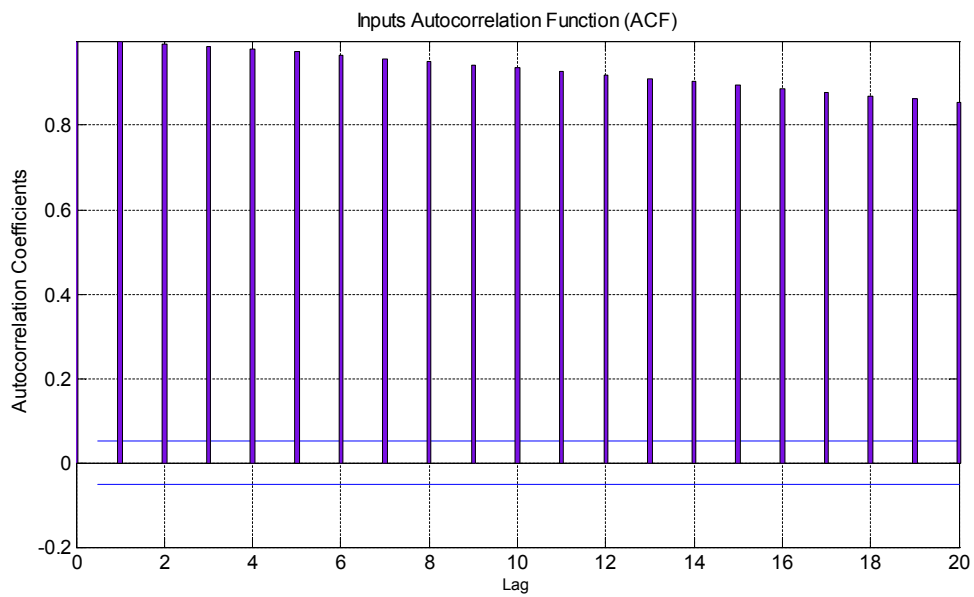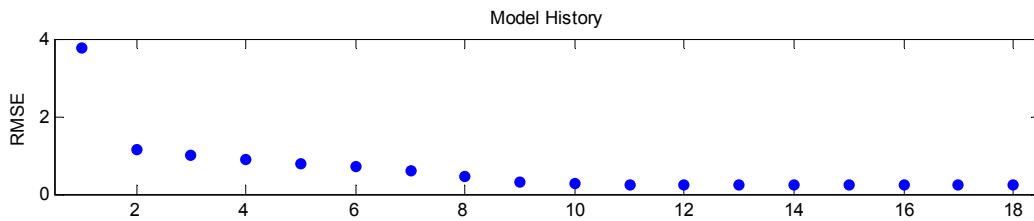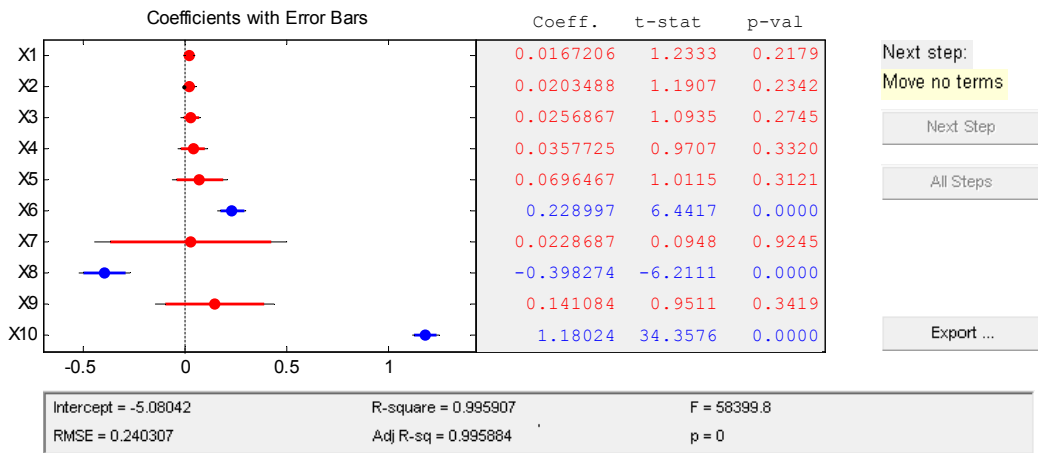
775
776
777

778



779
780 Fig. 5 Study Area
781



782
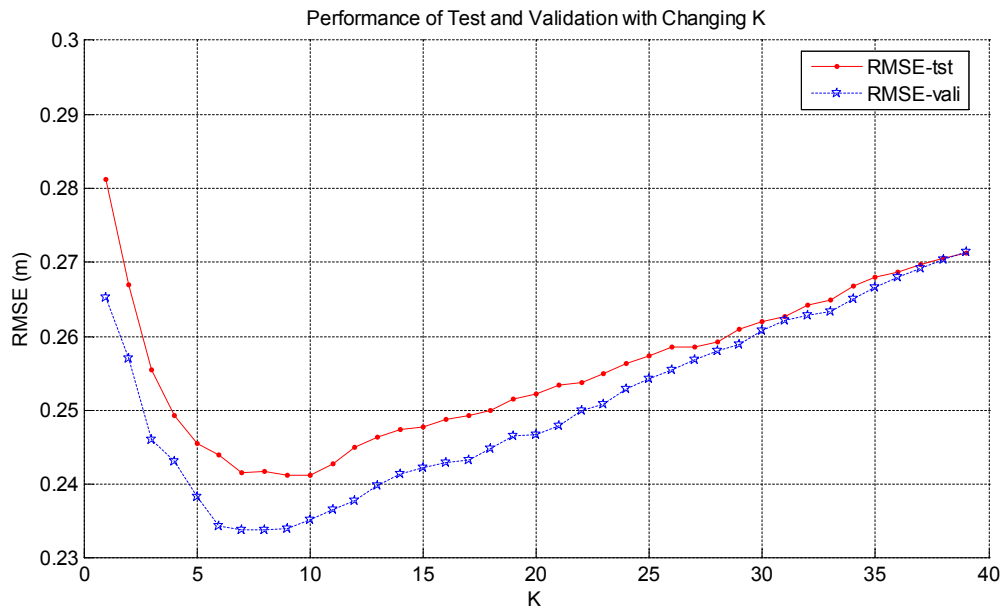783 Fig. 6 Auto-correlation of Water Levels at Lou-Shan Station
784

| | Coeff. | t-stat | p-val |
|---|---|---|---|
| X1 | 0.0167206 | 1.2333 | 0.2179 |
| X2 | 0.0203488 | 1.1907 | 0.2342 |
| X3 | 0.0256867 | 1.0935 | 0.2745 |
| X4 | 0.0357725 | 0.9707 | 0.3320 |
| X5 | 0.0696467 | 1.0115 | 0.3121 |
| X6 | 0.228997 | 6.4417 | 0.0000 |
| X7 | 0.0228687 | 0.0948 | 0.9245 |
| X8 | -0.398274 | -6.2111 | 0.0000 |
| X9 | 0.141084 | 0.9511 | 0.3419 |
| X10 | 1.18024 | 34.3576 | 0.0000 |

Intercept = -5.08042   R-square = 0.995907   F = 58399.8
RMSE = 0.240307   Adj R-sq = 0.995884   p = 0

Fig. 7 Stepwise Linear Model Process
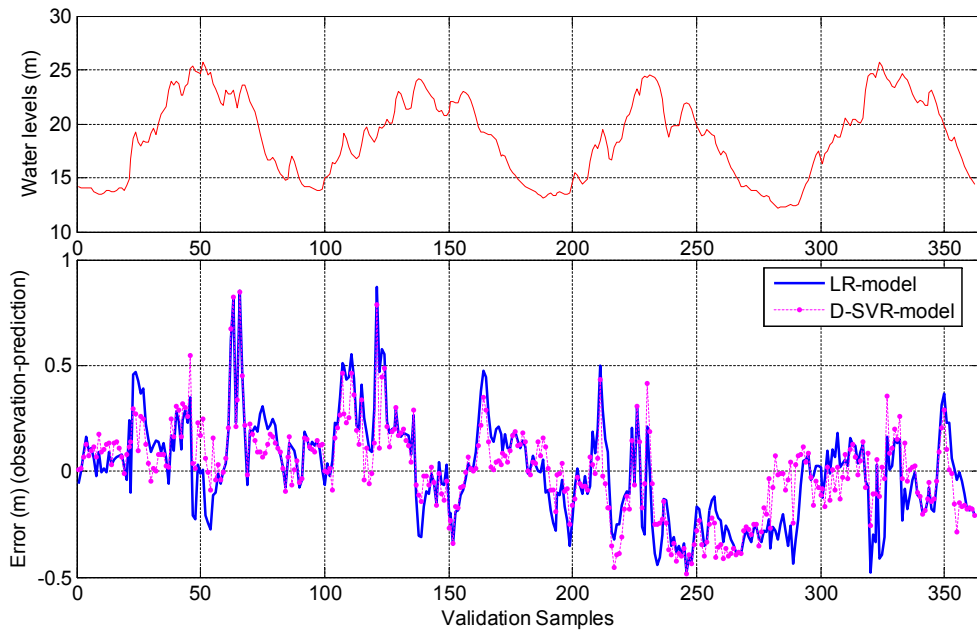
Fig. 8 Determining Optimal $k$ for NNM

791



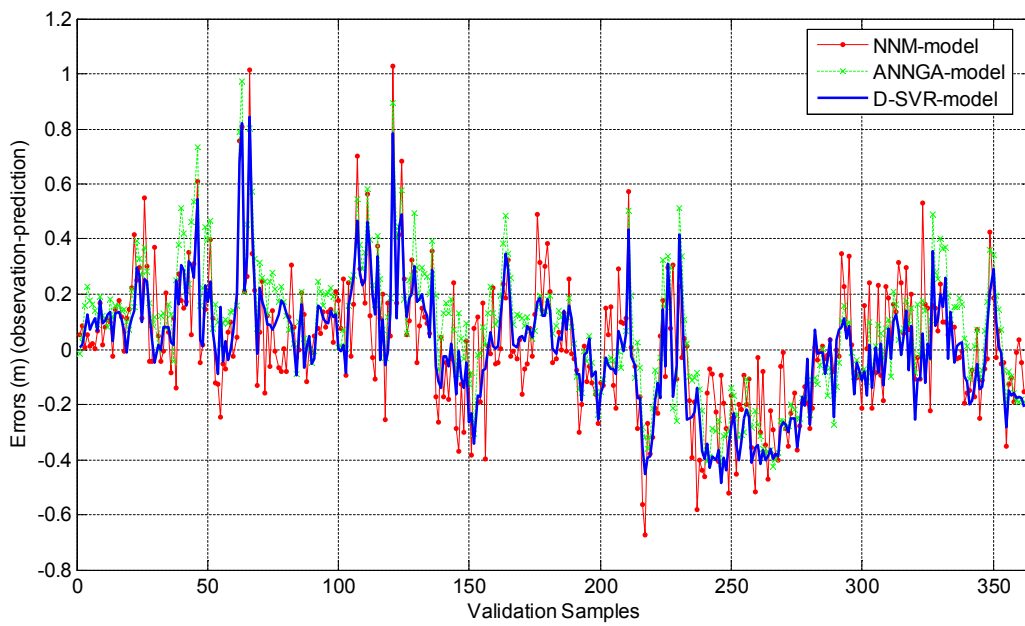Fig. 9 Comparison of absolute errors between LR and NNM



Fig. 10 Comparison of absolute errors between LR and ANN-GA

797



798
799
Fig. 11 Comparison of absolute errors between LR and D-SVR



800
801
Fig. 12 Comparison of absolute errors among NNM, ANN-GA, and D-SVR
802