

# Support-Vector-Machine-Based Models for Modeling Daily Reference Evapotranspiration With Limited Climatic Data in Extreme Arid Regions

Xiaohu Wen<sup>1</sup> · Jianhua Si<sup>1</sup> · Zhibin He<sup>1</sup> · Jun Wu<sup>2</sup> ·  
Hongbo Shao<sup>3,4</sup> · Haijiao Yu<sup>1</sup>

Received: 23 December 2014 / Accepted: 24 March 2015 /  
Published online: 11 April 2015  
© Springer Science+Business Media Dordrecht 2015

**Abstract** Evapotranspiration is a major factor that controls hydrological process and its accurate estimation provides valuable information for water resources planning and management, particularly in extremely arid regions. The objective of this research was to evaluate the use of a support vector machine (SVM) to model daily reference evapotranspiration ( $ET_0$ ) using limited climatic data. For the SVM, four combinations of maximum air temperature ( $T_{max}$ ), minimum air temperature ( $T_{min}$ ), wind speed ( $U_2$ ) and daily solar radiation ( $R_s$ ) in the extremely arid region of Ejina basin, China, were used as inputs with  $T_{max}$  and  $T_{min}$  as the base data set. The results of SVM models were evaluated by comparing the output with the  $ET_0$  calculated using Penman–Monteith FAO 56 equation (PMF-56). We found that the  $ET_0$  estimated using SVM with limited climatic data was in good agreement with those obtained using the conventional PMF-56 equation employing the full complement of meteorological data. In particular, three climatic parameters,  $T_{max}$ ,  $T_{min}$ , and  $R_s$  were enough to predict the daily  $ET_0$  satisfactorily. Moreover, the performance of SVM method was also compared with that of artificial neural network (ANN) and three empirical models including Priestley-Taylor, Hargreaves, and Ritchie. The results showed that the performance of SVM method was the best among these models. This offers significant potential for more accurate estimation of the  $ET_0$  with scarce data in extreme arid regions.

---

✉ Xiaohu Wen  
xhwen@lzb.ac.cn

✉ Hongbo Shao  
shaohongbochu@126.com

<sup>1</sup> Cold and Arid Regions Environmental and Engineering Research Institute, Chinese Academy of Sciences, No. 320 Donggang West Road, Lanzhou 730000 Gansu Province, China

<sup>2</sup> Next Fuel Inc., 122 North Main Street, Sheridan, WY 82801, USA

<sup>3</sup> Key Laboratory of Coastal Biology & Bioresources Utilization, Yantai Institute of Coastal Zone Research, Chinese Academy of Sciences, Yantai 264003, People's Republic of China

<sup>4</sup> Institute of Biotechnology, Jiangsu Academy of Agricultural Sciences, Nanjing 210014, China

**Keywords** Support vector machine · Reference evapotranspiration modeling · Limited climatic data · Extreme arid regions

## 1 Introduction

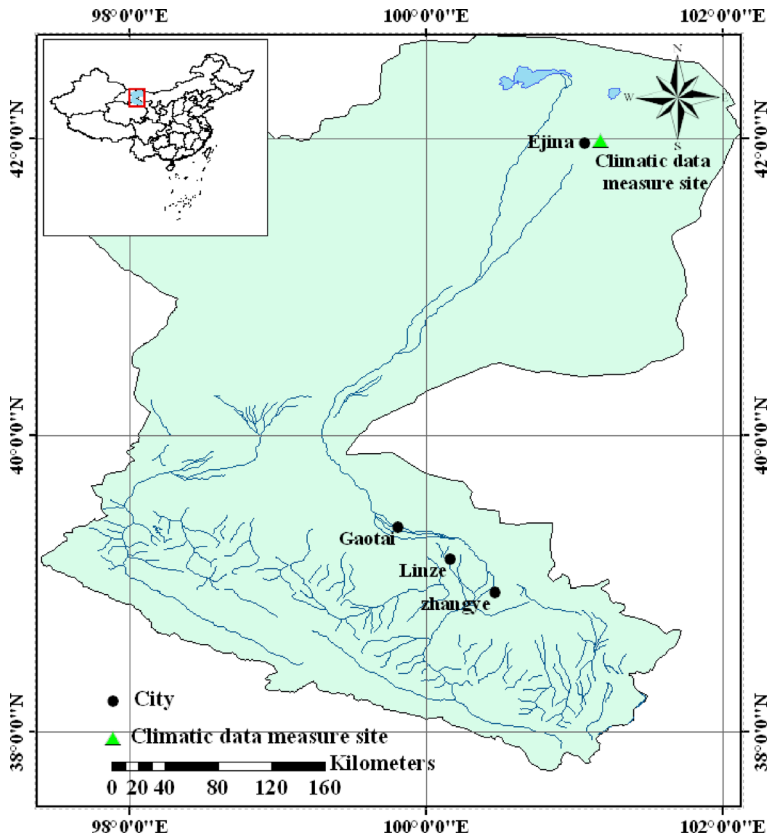
Evapotranspiration (ET) mainly controls several hydrological processes and its accurate estimation provides valuable information for water resources planning and management (Tabari et al. 2012), particularly in the arid area (Laaboudi et al. 2012).

The ET quantification, however, must be preceded by the determination of reference evapotranspiration ( $ET_0$ ) (López-Urrea et al. 2006). A great number of empirical equations have been developed for estimating  $ET_0$  using meteorological data. The Penman-Monteith FAO-56 combination equation (PMF-56) has been recommended by the Food and Agriculture Organization of the United Nations (FAO) as the standard equation for estimating  $ET_0$ . The PMF-56 equation is a physically based method, which requires a number of climatic parameters such as daily maximum temperature and minimum temperature, solar radiation, relative humidity, and wind speed. However, records for such weather variables are often incomplete or not always available for many locations so that the application of the PMF-56 model is limited (Cobaner 2011).

Evapotranspiration is an open, nonlinear, dynamic and complex system; therefore, it is difficult to derive an accurate formula to represent all the physical processes involved. As an alternative to traditional techniques, artificial neural networks (ANN) are highly appropriate for the modeling of non-linear processes. Many researchers have applied ANN to estimate  $ET_0$  (Chauhan and Shrivastava 2008; Laaboudi et al. 2012; Citakoglu et al. 2013; Kisi and Cengiz 2013; El-Shafie et al. 2014; Rahimikhoob 2014). These studies revealed that ANN models were superior in estimating  $ET_0$  than conventional methods such as regression and empirical equations. However, ANN have some disadvantages such as training slowly, requiring a large amount of training data, and easily getting stuck in a local minimum (Principe et al. 2000). Support vector machine (SVM), which is a novel learning machine based on statistical learning theory and a structural risk minimization principle, can be used for nonlinear system modeling (Vapnik 1995). Compared with ANN, SVM provides more reliable and better performance under the same training conditions (He et al. 2014). In last decade, SVM models have been extended to a wide range of hydrological problems (Raghavendra. N and Deka 2014).

Recently, some scientists began to use SVM for  $ET_0$  modeling. Kisi and Cimen (2009) studied the potential of SVM in modeling  $ET_0$  in central California, USA. Kisi (2012) examined the performances of least square support vector machine (LSSVM) in the modeling of  $ET_0$ . Tabari et al. (2012) examined the potential SVM for estimating  $ET_0$  in a semi-arid highland environment in Iran. Lin et al. (2013) developed SVM models for daily pan evaporation estimation and compared it with ANN models. These studies showed that SVM could be used to estimate  $ET_0$ , with relatively superior performance to ANN and empirical equations in modeling  $ET_0$ . Although SVM has excellent features, there are still limited studies using SVM in modeling  $ET_0$  research, particularly in the extremely arid regions with limited daily climatic data.

Ejjina basin, located in the lower reach of Heihe River, northwestern China (Fig. 1), is one of the most arid regions in the world. Water resource is a main controlling factor in economic development and ecological environment protection. However, the region is limited in water



**Fig. 1** Location study area and the climatic data measured site

resources with a mean annual precipitation of 42 mm. The Heihe River is the only runoff flow through the area. In the 1950s, the annual discharge of the Heihe River into the Ejina Basin was about  $12 \times 10^8 \text{ m}^3$ ; however, it was less than  $7 \times 10^8 \text{ m}^3$  in the 1990s. The accurate determination of  $ET_0$  is helpful to understand water balance in the extremely arid region and to determine the actual ecological water demand of ecosystem in the Ejina basin to serve as a reference for future water needs (Hou et al. 2010). Hence, a well performed model to improve daily  $ET_0$  estimation is always an important task to determine the actual ecological water demand and improve water use efficiency in the area (Hou et al. 2010). Generally, as a developing area, it is more difficult to collect sufficient daily meteorological data in such extreme regions for  $ET_0$  estimation.

The main objective of this study was to investigate the accuracy of SVM models for estimating daily  $ET_0$  using various combinations of daily meteorological data including maximum air temperature ( $T_{max}$ ), minimum air temperature ( $T_{min}$ ), wind speed ( $U_2$ ) and solar radiation ( $R_s$ ) in extremely arid environment of Ejina basin, northwestern China. In addition, the performances of the SVM models were compared with those of the ANN and three empirical models including Priestley–Taylor, Hargreaves and Ritchie equations to further test the SVM performance.

## 2 Materials and Methods

### 2.1 Penman-Monteith FAO 56 (PMF-56) Equation

In this paper, PMF-56 was used to provide the SVM targets to train and test the SVM models. As the sole standard method for the computation of  $ET_0$  when no measured lysimeter data are available, PMF-56 method is described by Allen et al. (1998):

$$ET_{0-PMF-56} = \frac{0.408\Delta(R_n - G) + \gamma \frac{900}{T_{mean} + 273} U_2 (e_s - e_a)}{\Delta + \gamma(1 + 0.34U_2)} \quad (1)$$

where  $ET_{0-PMF-56}$  is the reference evapotranspiration ( $\text{mm day}^{-1}$ );  $R_n$  is the net radiation ( $\text{MJ m}^{-2} \text{day}^{-1}$ );  $G$  is the soil heat flux ( $\text{MJ m}^{-2} \text{day}^{-1}$ );  $\gamma$  is the psychrometric constant ( $\text{kPa } ^\circ\text{C}^{-1}$ );  $e_s$  is the saturation vapor pressure ( $\text{kPa}$ );  $e_a$  is the actual vapor pressure ( $\text{kPa}$ );  $\Delta$  is the slope of the saturation vapor pressure-temperature curve ( $\text{kPa } ^\circ\text{C}^{-1}$ );  $T_{mean}$  is the average daily air temperature ( $^\circ\text{C}$ ); and  $U_2$  is the mean daily wind speed at 2 m ( $\text{m s}^{-1}$ ). The computation of all data required for calculating  $ET_0$  followed the method and procedure given in Chapter 3 of FAO-56 (Allen et al. 1998).

### 2.2 Hargreaves Equation

Hargreaves and Samani (1985) presented a formula for the estimation of reference evapotranspiration when daily weather data is limited or missing. The equation has the form:

$$ET_{0-Hargreaves} = 0.0023R_a \left( \frac{T_{max} + T_{min}}{2} + 17.8 \right) \sqrt{T_{max} - T_{min}} \quad (2)$$

where  $ET_{0-Hargreaves}$  is the reference evapotranspiration ( $\text{mm day}^{-1}$ );  $R_a$  is the water equivalent of the extraterrestrial radiation ( $\text{mm day}^{-1}$ ) computed according to Allen et al. (1998).

### 2.3 Ritchie Equation

Ritchie equation was described by Jones and Ritchie (1990):

$$ET_{0-Ritchie} = \alpha_1 [3.87 \times 10^{-3} R_s (0.6T_{max} + 0.4T_{min} + 29)] \quad (3)$$

where  $ET_{0-Ritchie}$  is the reference evapotranspiration ( $\text{mm d}^{-1}$ );  $R_s$  is the solar radiation ( $\text{MJ m}^{-2} \text{d}^{-1}$ ); and  $\alpha_1$  is defined as follows:

$$5 < T_{max} \leq 35 \text{ } ^\circ\text{C} \quad \alpha_1 = 1.1$$

$$T_{max} > 35 \text{ } ^\circ\text{C} \quad \alpha_1 = 1.1 + 0.05(T_{max} - 35)$$

$$T_{max} < 5 \text{ } ^\circ\text{C} \quad \alpha_1 = 0.01 \exp [0.18(T_{max} + 20)]$$

### 2.4 Priestley and Taylor Equation

Priestley and Taylor equation (Priestley and Taylor 1972) for computing  $ET_0$  value is expressed as:

$$ET_{0-Priestley-Taylor} = \frac{\alpha}{\lambda} \frac{\Delta}{\Delta + \gamma} (R_n - G) \alpha = 1.26 \tag{4}$$

Where  $ET_{0-Priestley-Taylor}$  is the reference evapotranspiration ( $\text{mm day}^{-1}$ );  $\alpha$  is empirical coefficient; and  $\lambda$  is latent heat of the evaporation ( $\text{MJ/Kg}$ ).

Empirical equations are usually developed using local-related data, Allen et al. (1994) recommended that empirical equations should be calibrated using PMF-56 method. Calibrated  $ET_0$  is calculated as

$$ET_0 = a + b \times ET_{method} \tag{5}$$

where  $ET_0$  is the reference evapotranspiration defined by PMF-56 method,  $ET_{method}$  represents the evapotranspiration estimated by the evaluated empirical models, and  $a$  and  $b$  are the regression constants.

### 2.5 Support Vector Machine (SVM)

Support vector machine (SVM), which is a supervised learning model based on statistical learning theory introduced by Vapnik (1995). Generally, support vector regression (SVR) is used to describe regression with SVM. Here, we only show a brief introduction of SVR, while detailed principles and algorithms of SVM can be found in Müller et al. (1997).

In SVM, the basic idea is to map the data  $x$  into a high dimensional feature space via a nonlinear mapping  $\pi$  and to do linear regression in this space (Boser et al. 1992; Vapnik 1995).

The regression estimation with SVR is to estimate a function according to a given data set  $\{(x_i, y_i)\}_i^n$ , where  $x_i$  denotes the input vector;  $y_i$  denotes the output value and  $n$  is the total number of data sets.

In SVM, the regression function is approximated by the following function:

$$f(x) = \omega \cdot \varphi(x) + b \tag{6}$$

where  $\omega$  is a weight vector, and  $b$  is a bias.  $\varphi(x)$  denotes a nonlinear transfer function that maps the input vectors into a high-dimensional feature space in which theoretically a simple linear regression can cope with the complex nonlinear regression of the input space.

The coefficients  $\omega$  and  $b$  can be estimated by minimizing the following regularized risk function:

$$R_{reg}(f) = C \frac{1}{n} \sum_{i=1}^N L_\varepsilon(f(x_i), y_i) + \frac{1}{2} \|\omega\|^2 \tag{7}$$

$$L(f(x), y) = \begin{cases} |f(x) - y| - \varepsilon & \text{for } |f(x) - y| \geq \varepsilon \\ 0 & \text{otherwise} \end{cases} \tag{8}$$

where  $C$  is a positive constant named penalty parameter,  $L_\varepsilon(f(x_i), y_i)$  is called  $\varepsilon$ -insensitive loss

function that measures the empirical risk of the training data;  $(1/2)\|\omega\|^2$  is the regularization term;  $\varepsilon$  is the tube size of SVM.

Finally, a nonlinear regression function is obtained using the following expression

$$f(x) = \sum_{i=1}^l (\alpha_i - \alpha_i^*)k(x_i, x) + b \tag{9}$$

where  $\alpha_i$  and  $\alpha_i^*$  are the introduced Lagrange multipliers. With the utilization of the Karush-Kuhn-Tucker (KKT) conditions, only a limited number of coefficients will not be zero among  $\alpha_i$  and  $\alpha_i^*$ . The related data points could be referred to the support vectors.  $k(x_i, x)$  refers to kernel function describes the inner product in the D-dimension feature space.

$$k(x_i, x) = \sum_{j=1}^D \varphi_j(x_i)\varphi_j(x) \tag{10}$$

It can be shown that any symmetric kernel function  $k$  satisfying Mercer’s condition corresponds to a dot product in some feature space (Boser et al. 1992). In this paper, radius basis function (RBF) is selected as the kernel function. The RBF is defined as following:

$$k(x_i, x) = \exp\left(-\gamma\|x_i\|^2\right), \lambda > 0 \tag{11}$$

There are three parameters while using RBF kernels such as penalty parameter  $C$ , error exceeding  $\varepsilon$  and kernel function’s parameter  $\gamma$ . The general performance of SVM models depends on a proper setting of these parameters. In this study,  $C$ ,  $\varepsilon$  and  $\gamma$  were determined through grid-search algorithm with cross-validation as described by Hus et al. (2010), SVM algorithms were developed using Matlab libsvm Toolbox (Chang and Lin 2011).

### 2.6 Artificial Neural Network (ANN)

ANN is a massively parallel distributed information processing system that has certain performance characteristics resembling biological neural networks of the human brain (Haykin 1999). A neural network is characterized by its architecture that represents the pattern of connection between nodes, its method of determining the connection weights and the activation function. The most commonly used neural network structure is the feed forward hierarchical architecture. A typical three-layered feed-forward neural network is comprised of a multiple elements also called nodes, and connection pathways that link them (Haykin 1999). The nodes are processing elements of the network and are normally known as neurons, reflecting the fact the neural network method model is based on the biological neural network of the human brain. A neuron receives an input signal, processes it, and transmits an output signal to other interconnected neurons.

In the hidden and output layers, the net input to unit  $i$  is of the form

$$Z = \sum_{j=1}^k w_{ji}y_j + \theta_i \tag{12}$$

where  $w_{ji}$  is the weight vector of unit  $i$  and  $k$  is the number of neurons in the layer above the layer that includes unit  $i$ .  $y_j$  is the output from unit  $j$ , and  $\theta_i$  is the bias of unit  $i$ . This weighted

sum  $Z_i$  which is called the incoming signal of unit  $i$ , is then passed through a transfer function  $f$  to yield the estimates  $\hat{y}_i$  for unit  $i$ . The sigmoid function is continuous, differentiable everywhere, and monotonically increasing. The sigmoid transfer function,  $f_i$ , of unit  $i$ , is of the form

$$\hat{y}_i = \frac{1}{1 + e^{-z}} \quad (13)$$

A training algorithm is needed to solve a neural network problem. Since there are so many types of algorithms available for training a network, selection of an algorithm that provides the best fit to the data is required. In the current research, the ANN models were trained using the Levenberg–Marquardt training algorithm. The sigmoid and linear activation functions were used for the hidden and output node(s), respectively.

### 3 Case Study

#### 3.1 Observation Data and Statistical Analysis

The climatic data in the site located near Ejina City (101°09'17.69"E, 41°58'53.95"N, altitude 927.32 m) were observed during the *Phragmites communis*' growing season of May 9th to October 1th, 2004 (Fig. 1), with the total numbers of growing days of about 146 days. An automatic weather measurement system was installed in a flat field with *Phragmites* stand used to measure the primary climatic parameters including net radiation, soil heat flux, air temperature, water vapor pressure, humidity, wind speed and direction, dew point temperature and solar radiation. The detailed measurement system and methods can be found in Si et al. (2005).

The daily climatic data employed in this study were composed of  $T_{max}$ ,  $T_{min}$ ,  $U_2$ , and  $R_s$ . The data from May 9th to August 18th, the first 102 records (about 70 % of total data) were used for training the models, and the remaining 44 records from August 19th to October 1th (about 30 %) were used for testing. The statistical parameters of daily climatic data were shown in Table 1.  $U_2$  shows a skewed distribution. According to the statistical properties of those data sets, no statistically significant differences between the divisions of the data were observed. Obviously, training data contain sufficient information about the system behavior to qualify as a system model.

In order to eliminate dimension difference, all the climatic data were scaled to [0, 1] before input the SVM model. The formula is defined as following:

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (14)$$

where  $x_{new}$  is the normalization data;  $x_{min}$  is the minimum data;  $x_{max}$  is the maximum data.

#### 3.2 Model Development

Selecting appropriate input variables is important for SVM and ANN models development since it provides the basic information about the system being modeled. Temperature is the most predominant physical factor in the evaporation process. So,  $T_{max}$  and  $T_{min}$  were selected as an input. Some studies reported that  $R_s$  and  $U_2$  are more effective variables for estimating  $ET_0$  in arid and semiarid zone (Cobaner 2011; Tabari et al. 2012). In current study, the performance of SVM and of ANN  $ET_0$  was compared with daily PMF-56  $ET_0$ . To achieve

**Table 1** Statistical parameters of climatic data and PMF-56 ET<sub>0</sub> in each data set

Climatic data and the PMF-56 ET <sub>0</sub>		Maximum	Minimum	Mean	Std.	SK	CV
<i>T</i> <sub>max</sub> (°C)	All	38.32	10.90	29.35	5.31	-0.88	0.18
	Training	38.32	14.44	31.12	4.38	-1.04	0.14
	test	32.81	10.90	25.25	5.05	-0.84	0.2
<i>T</i> <sub>min</sub> (°C)	All	23.47	-8.72	12.91	6.20	-0.94	0.48
	Training	23.47	4.82	15.46	4.10	-0.47	0.26
	Test	18.14	-8.72	7.00	6.23	-0.49	0.89
<i>R</i> <sub>s</sub> (MJ m <sup>-2</sup> day <sup>-1</sup> )	All	30.06	5.96	21.88	5.19	-0.7	0.24
	Training	30.06	5.95	23.4	5.12	-1.39	0.22
	Test	23.48	8.74	18.33	3.29	-0.80	0.18
<i>U</i> <sub>2</sub> (m s <sup>-1</sup> )	All	3.86	0.18	1.20	0.60	1.27	0.50
	Training	3.86	0.18	1.26	0.60	1.29	0.48
	Test	2.83	0.36	1.08	0.59	1.38	0.55
PMF-56 ET <sub>0</sub> (mm day <sup>-1</sup> )	All	4.86	0.29	2.92	1.08	-0.21	0.37
	Training	4.86	0.51	3.26	1.01	-0.55	0.31
	Test	3.48	0.29	2.14	0.79	-0.26	0.37

*Std.* is the standard deviation, *SK* is the skewness, *CV* is the coefficient of variation

this, various combinations of daily climatic data including *T*<sub>max</sub>, *T*<sub>min</sub>, *U*<sub>2</sub>, and *R*<sub>s</sub> were used as inputs to SVM and ANN models to estimate ET<sub>0</sub>. The four input combinations evaluated were (1) *T*<sub>max</sub> and *T*<sub>min</sub>; (2) *T*<sub>max</sub>, *T*<sub>min</sub> and *U*<sub>2</sub>; (3) *T*<sub>max</sub>, *T*<sub>min</sub> and *R*<sub>s</sub>; (4) *T*<sub>max</sub>, *T*<sub>min</sub>, *U*<sub>2</sub> and *R*<sub>s</sub>.

### 3.3 Models Performance Criteria

The performances of the models developed in this research were assessed using various standard statistical performance evaluation criteria such as coefficient of correlation (*r*), root mean squared error (RMSE), and mean absolute error (MAE). *r* measures the degree to which two variables are linearly related. RMSE and MAE provide different types of information about the predictive capabilities of the model. The RMSE measures the goodness-of-fit relevant to high ET<sub>0</sub> values whereas the MAE yields a more balanced perspective of the goodness-of-fit at moderate value distribution of the estimation errors.

The following equations were used for the computation of the above parameters:

$$r = \frac{\sum_{i=1}^n (ET_{0i}^p - \overline{ET_0^p})(ET_{0i}^o - \overline{ET_0^o})}{\sqrt{\sum_{i=1}^n (ET_{0i}^p - \overline{ET_0^p})^2 (ET_{0i}^o - \overline{ET_0^o})^2}} \tag{15}$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (ET_{0i}^p - ET_{0i}^o)^2}{n}} \tag{16}$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |ET_{0i}^p - ET_{0i}^o| \tag{17}$$



where  $ET_{0i}^p$  and  $ET_{0i}^o$  are the  $i$ th estimated and PMF-56  $ET_0$  values, respectively;  $\overline{ET_0^p}$  and  $\overline{ET_0^o}$  are the average of  $ET_{0i}^p$  and  $ET_{0i}^o$ ; and  $n$  is the total numbers of data. The best fit between observed and calculated values would have  $r=1$ ,  $RMSE=0$  and  $MAE=0$ , respectively.

In order to test the robustness of the developed model, it is important to test the model using some other performance evaluation criteria such as relative error ( $RE$ ) and threshold statistics ( $TS$ ) (Jain and Indurthy 2003). The  $TS$  for a level of  $x\%$  is a measure of the consistency in modeling errors from a particular model. The  $TS$  are represented as  $TS_x$  and expressed as a percentage. This criterion can be expressed for different levels of relative error ( $RE$ ) from the model.

$$RE = \frac{|ET_{0i}^p - ET_{0i}^o|}{ET_{0i}^o} \quad (18)$$

$$TS_x = \frac{n_x}{n} \times 100 \quad (19)$$

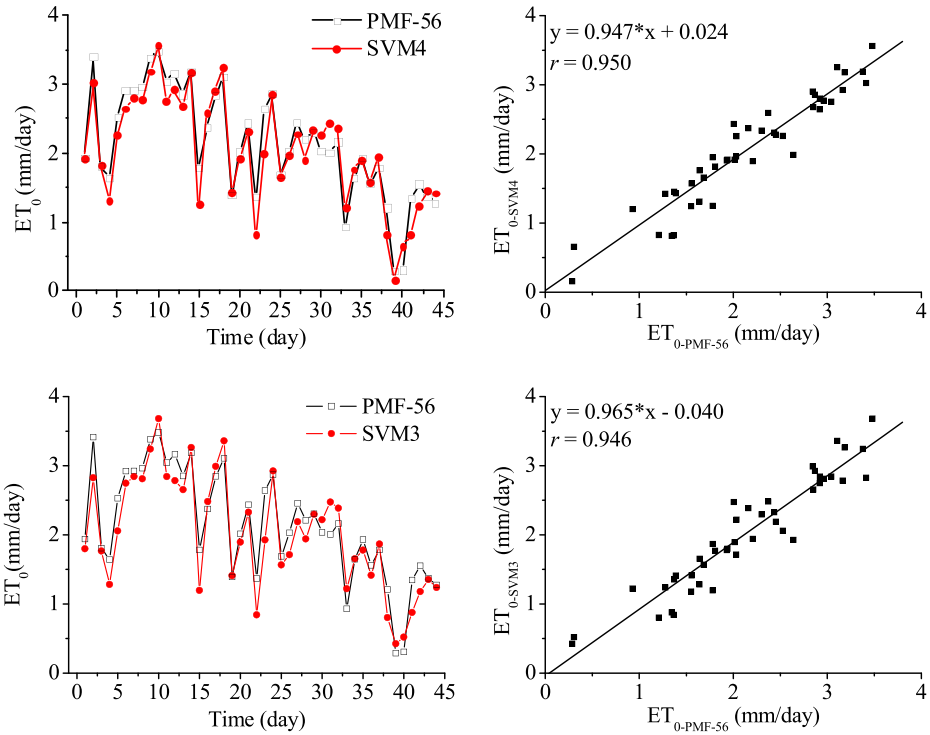
where,  $n_x$  is the number of data points for which the  $RE$  is less than  $x\%$ ;  $n$  is the total number of data points computed. Clearly, higher  $n_x$  and  $TS_x$  values would indicate better model performance.

### 3.4 Results and Discussion

The performance of SVM models for PMF-56  $ET_0$  and the parameters  $C$ ,  $\varepsilon$ ,  $\gamma$  of the optimum SVM model were given in Table 2. It is apparent that all of the models performed similarly in training periods and testing periods, as the values of  $RMSE$  and  $MAE$  don't vary significantly, and all  $r$  are also very close to unity. In testing periods, it is apparent that SVM4 and SVM3 models were better than SVM1 and SVM2 models for PMF-56  $ET_0$  estimation (Table 2). Based on the performance statistics, SVM4 whose inputs combinations were  $T_{max}$ ,  $T_{min}$ ,  $U_2$  and  $R_s$  had the smallest value of the  $RMSE$  (0.262 mm/day),  $MAE$  (0.207 mm/day) and higher value of  $r$  (0.950) than other model in the testing periods. Therefore, it was selected as the best-fit model for estimating the PMF-56  $ET_0$ . SVM3 model whose inputs include  $T_{max}$ ,  $T_{min}$  and  $R_s$  with  $RMSE$  of 0.282 mm/day,  $MAE$  of 0.228 mm/day and  $r$  of 0.946 provided the secondly best PMF-56  $ET_0$  estimation. Comparative analysis of the performance statistics showed that, SVM4 and SVM3 models performed similarly. Moreover,  $r$  values were also very close to

**Table 2** Optimal SVM parameters and the performance statistics of SVM models during training and testing periods

Models	Input	Parameter			Training periods			Testing periods		
		C	$\gamma$	$\varepsilon$	$r$	RMSE mm/day	MAE mm/day	$r$	RMSE mm/day	MAE mm/day
SVM1	$T_{min}$ , $T_{max}$	5.278	1.741	0.079	0.818	0.581	0.461	0.772	0.539	0.446
SVM2	$T_{min}$ , $T_{max}$ , $U_2$	5.728	1.000	0.004	0.826	0.569	0.411	0.773	0.504	0.418
SVM3	$T_{min}$ , $T_{max}$ , $R_s$	5.278	1.000	0.064	0.927	0.383	0.244	0.946	0.286	0.228
SVM4	$T_{min}$ , $T_{max}$ , $U_2$ , $R_s$	0.330	0.574	0.002	0.921	0.395	0.269	0.950	0.262	0.207



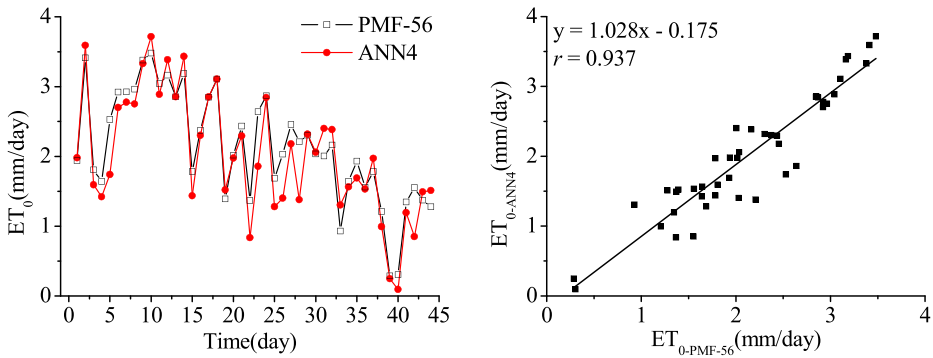
**Fig. 2** Comparison of the  $ET_0$  values estimated by SVM4, SVM3 and PMF-56 equation during testing periods

unity. For practical applications, SVM4 and SVM3 had good accuracy in PMF-56  $ET_0$  modeling and the selection of one model over the other should be dependent upon the available meteorological data. Furthermore, SVM3, in which  $T_{max}$ ,  $T_{min}$  and  $R_s$  are needed, performed well in PMF-56  $ET_0$  modeling and could be used in the developing areas with limited weather data.

The comparison of the  $ET_0$  values computed by the PMF-56 equation and the values estimated by SVM4 and SVM3 models were shown in Fig. 2, in the form of line graphs and scatter plots. The  $ET_0$  values estimated by the SVM models are closely to that computed using the PMF-56  $ET_0$  values and followed the same trend. The consistence revealed that the two models showed good estimation accuracy of the PMF-56  $ET_0$  (Fig. 2).

**Table 3** The structure and the performance statistics of ANN models during training and testing periods

Models	Input	Input-hidden-output nodes	Training periods			Testing periods		
			<i>r</i>	RMSE mm/day	MAE mm/day	<i>r</i>	RMSE mm/day	MAE mm/day
ANN1	$T_{min}, T_{max}$	2-6-1	0.856	0.522	0.400	0.750	0.561	0.440
ANN2	$T_{min}, T_{max}, U_2$	3-3-1	0.865	0.506	0.397	0.682	0.587	0.467
ANN3	$T_{min}, T_{max}, R_s$	3-3-1	0.935	0.359	0.254	0.923	0.337	0.268
ANN4	$T_{min}, T_{max}, U_2, R_s$	4-2-1	0.936	0.355	0.266	0.937	0.322	0.236



**Fig. 3** Comparison of the  $ET_0$  values estimated by ANN4 and PMF-56 equation model during testing periods

In order to evaluate the ability of SVM model relative to ANN model, four ANN models were developed using the same variables combinations for  $ET_0$  modeling. The optimal number of neuron in the hidden layer was identified using a trial and error procedure by varying the number of hidden neurons from 2 to 15. Furthermore, the optimal network architecture was selected based on the one with minimum of MSE. The final ANN architecture and the performance statistics of each model were shown in Table 3. According to the testing periods results, ANN4 (4-2-1) model with the input combination  $T_{max}$ ,  $T_{min}$ ,  $U_2$  and  $R_s$  had the smallest RMSE (0.322 mm/day), MAE (0.268 mm/day) and the highest  $r$  (0.937), performed best. ANN3 (3-3-1) model, whose inputs were  $T_{max}$ ,  $T_{min}$  and  $R_s$  had smaller RMSE (0.337 mm/day), MAE (0.268 mm/day) and higher  $r$  (0.923), ranked the second in  $ET_0$  estimations. However, a comparison of the performance criteria for ANN models (Table 3) with those of SVM in Table 2 showed that all the SVM models have performed better than the corresponding ANN models in modeling the PMF-56  $ET_0$ .

It is important to evaluate not only the average estimation error but also the distribution of estimation errors when assessing the performance of any model for its applicability in modeling  $ET_0$ . Comparing the best SVM model SVM4 and the best ANN model ANN4 for modeling  $ET_0$ , SVM4 gave 28 estimates lower than the 10 % relative error in the testing periods, while ANN4 had 23 estimates lower than the 10 % error. Furthermore, SVM4 had 15 estimates lower than the 5 % error, while ANN4 had 13 estimates lower than the 5 % relative error, respectively. The SVM model yielded more accurate results than the ANN model.

The comparison between best ANN model ANN4 in modeling the  $ET_0$  and PMF-56  $ET_0$  was shown in Fig. 3. Compared to Fig. 2 for SVM4 model, it further confirmed that although both the ANN and SVM had comparable performance during testing periods, the SVM models

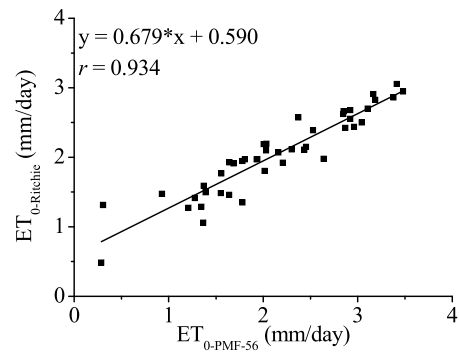
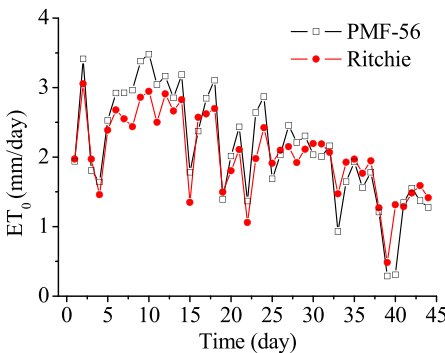
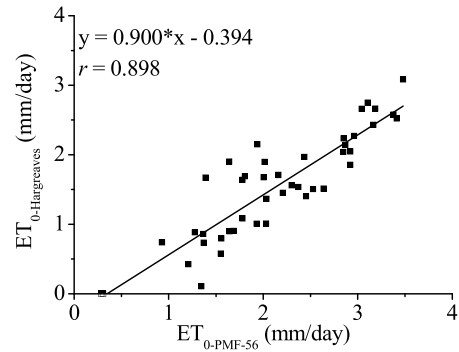
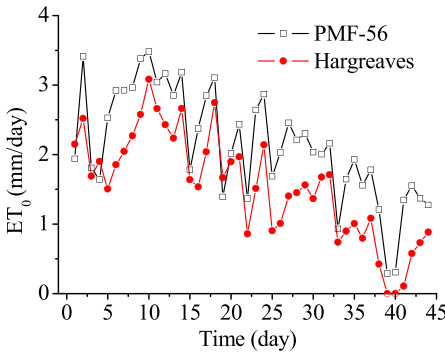
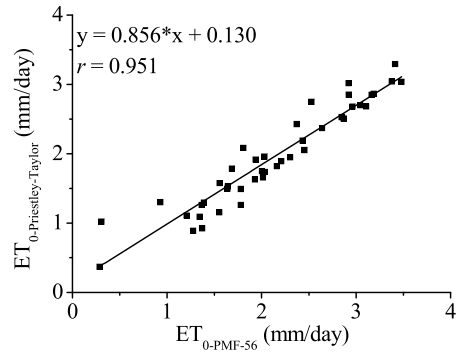
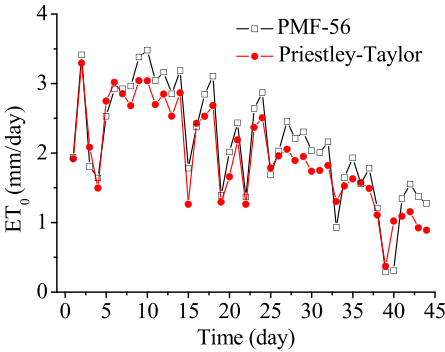
**Table 4** The calibration coefficients of the empirical models and performance statistics of the empirical models during testing periods

Models	Calibration coefficients		$r$	RMSE mm/day	MAE mm/day
	a	b			
Priestley–Taylor	-0.198	0.763	0.951	0.304	0.266
Hargreaves	-1.772	0.829	0.898	0.705	0.642
Ritchie	-0.493	0.689	0.934	0.338	0.279

provided more accurate  $ET_0$  estimates than the ANN during the more important independent testing stage. Overall, the results obtained confirmed the capability of SVM models for  $ET_0$  estimates.

The performance of SVM models was further compared with three different empirical models including Priestley–Taylor, Hargreaves and Ritchie equations. These empirical models were firstly applied to calculate evapotranspiration based on the training data, and then calibrated using the PMF-56  $ET_0$  by the equation (5).

Priestley–Taylor, Hargreaves and Ritchie were calibrated by  $a$  and  $b$  coefficients. The performance statistics in testing periods of each model was given in Table 4. Priestley–Taylor



**Fig. 4** Comparison of the  $ET_0$  values estimated by empirical models and PMF-56 equation during testing periods

equation had the smallest RMSE, MAE and the highest  $r$ , with the best performance. Ritchie equation performed the second best in  $ET_0$  estimations. Hargreaves model performed the worst in the PMF-56  $ET_0$  estimation. Compared with those of SVM4 in Table 2, Priestley–Taylor equation had the highest  $r$  (0.951) that provided information for linear dependence between observations and corresponding estimates. It is not always expected that  $r$  is in agreement with performance criteria such as the RMSE. In the present study the main model performance criterion is the RMSE. The best model was selected by considering this criterion. From this viewpoint, it revealed that SVM4 model gave more accurate results than the empirical models in modeling PMF-56  $ET_0$ .

For the distribution of estimation errors, in testing periods, Priestley-Taylor, Hargreaves and Ritchie methods had 16, 3 and 19 estimates lower than the 10 % error, respectively. Furthermore, Priestley-Taylor, Hargreaves and Ritchie methods had 7, 0 and 6 estimates lower than the 5 % relative error, respectively. From the viewpoint of relative error, SVM4 model still performed better than the empirical methods.

The  $ET_0$  estimates of the empirical methods were illustrated in Fig. 4 in the form of line graphs and scatter plots. All of the empirical models underestimated the  $ET_0$  values calculated by PMF-56 model. The performance differences between the empirical equations and the SVM approaches models showed that the SVM models performed better than the empirical equations.

The estimation of total PMF-56  $ET_0$  obtained from the estimated  $ET_0$  values was also considered for comparison due to its importance in water balance calculation, water resources planning and management. The total estimated  $ET_0$  amounts in testing periods were given in Table 5. It showed that all models underestimate total PMF-56  $ET_0$  value in testing periods. SVM4 and ANN4 models whose input parameters were  $T_{max}$ ,  $T_{min}$ ,  $U_2$  and  $R_s$  estimated the total PMF-56  $ET_0$  value of 94.14 mm as 90.27 mm and 89.03 mm, with an underestimation of 4.1 % and 5.4 %, respectively. While SVM3, Priestley-Taylor, Hargreaves and Ritchie equations estimated the total PMF-56  $ET_0$  value as 89.05 mm, 86.31 mm, 67.91 mm and 89.88 mm with underestimation of 5.4 %, 8.3 %, 27.9 % and 4.5 %, respectively. The total PMF-56  $ET_0$  amount estimates of SVM4, SVM3, ANN4 and Ritchie equation were closer to the PMF-56  $ET_0$  values. Among the models, SVM4 model had the best estimate (−4.1 %) and Ritchie equation had the secondly best estimate (−4.5 %), while Hargreaves equation had the worst (−27.9 %) in terms of total estimated PMF-56  $ET_0$  values.

As a whole, the findings of this study revealed that SVM model seemed to be more adequate than ANN, Priestley-Taylor, Hargreaves and Ritchie equations for the  $ET_0$  modeling and can be employed successfully in  $ET_0$  estimation in the extreme arid regions with limited

**Table 5** Total  $ET_0$  values calculated by various models during testing period

Models	Total evapotranspiration (mm)	Relative error (%)
PMF-56	94.14	–
SVM3	89.05	−5.4
SVM4	90.26	−4.1
ANN4	89.03	−5.4
Priestley-Taylor	86.31	−8.3
Hargreaves	67.91	−27.9
Ritchie	89.88	−4.5

climatic data. SVM3 model which only needed  $T_{max}$ ,  $T_{min}$  and  $R_s$  can be considered as simple model that offers a significant potential for accurate estimation of daily  $ET_0$ . SVM4 model with  $T_{max}$ ,  $T_{min}$ ,  $U_2$  and  $R_s$  as input variables exhibited good daily  $ET_0$  estimation ability and produced better results. They are the recommended models by a lack of appropriate meteorological data for the application of the PMF-56 equation in extreme arid regions.

Generally, there are few limitations when SVM model is applied in practice. For many data-driven techniques, the amount of data size used to develop the model usually does limit their performance. To realize reliable forecasts, long-term weather data are required, but our computations show that SVM model is still reliable in  $ET_0$  modeling even short-term weather records are used. The other limitation is that the model has been developed using data from a single site. However, this should not be seen as constituting a major problem since the analysis can easily be widened if more data from other stations become available for analysis. More data from different sources would allow the model to capture the patterns of data from a wider range of scenarios, thus increasing the geographical scope of its validity (Adeloye et al. 2012).

## 4 Conclusions

The accurate estimation of evapotranspiration is one of the most important issues in the management of water resources. This work investigated the applicability of SVM for daily  $ET_0$  modeling using limited climatic data in the extremely arid regions of Ejina basin, northwestern China. Four models were developed using different combinations of four daily climatic data including  $T_{max}$ ,  $T_{min}$ ,  $R_s$  and  $U_2$ . The developed SVM models were tested using the  $ET_0$  calculated by PMF-56. The results demonstrated that SVM could be applied successfully to establish accurate and reliable PMF-56  $ET_0$  modeling. Particularly, SVM model whose inputs included  $T_{max}$ ,  $T_{min}$  and  $R_s$  provided good  $ET_0$  estimate, this is especially true in the developing areas where reliable weather data sets are limited.

Based on the comparison of SVM models with ANN and empirical models such as Priestley–Taylor, Hargreaves, Ritchie equations, the SVM gave more accurate results than the ANN and empirical models in the estimation of PMF-56  $ET_0$ . SVM can be successfully used for modeling daily PMF-56  $ET_0$  when there are limited climatic data in extreme arid regions.

**Acknowledgments** This work was funded by the National Basic Research Program of China (2013CB429906), the authors also wish to thank anonymous reviewers for their reading of the manuscript and for their suggestions and critical comments.

## References

- Adeloye AJ, Rustum R, Kariyama ID (2012) Neural computing modeling of the reference crop evapotranspiration. *Environ Model Softw* 29(1):61–73
- Allen RG, Smith M, Pereira LS (1994) An update for the definition of reference evapotranspiration. *ICID Bull* 43:1–34
- Allen RG, Pereira LS, Raes D, Smith M (1998) Crop evapotranspiration-guidelines for computing crop water requirements. *FAO irrigation and drainage. paper no. 56*. FAO, Rome
- Boser, B.E., Guyon, I.M., Vapnik, V.N., (1992) A training algorithm for optimal margin classifiers. *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*. ACM Press, pp.144–152
- Chang CC, Lin CJ (2011) LIBSVM: a library for support vector machines. *ACM Trans Intell Syst Technol* 2:1–27, Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>

- Chauhan S, Shrivastava RK (2008) Performance evaluation of reference evapotranspiration estimation using climate based methods and artificial neural networks. *Water Resour Manag* 23(5):825–837
- Citakoglu H, Cobaner M, Haktanir T, Kisi O (2013) Estimation of monthly mean reference evapotranspiration in Turkey. *Water Resour Manag* 28(1):99–113
- Cobaner M (2011) Evapotranspiration estimation by two different neuro-fuzzy inference systems. *J Hydrol* 398: 292–302
- El-Shafie A, Najah A, Alsulami HM, Jahanbani H (2014) Optimized neural network prediction model for potential evapotranspiration utilizing ensemble procedure. *Water Resour Manag* 28(4):947–967
- Hargreaves GH, Samani ZA (1985) Reference crop evapotranspiration from temperature. *Appl Eng Agric* 1(2): 96–99
- Haykin S (1999) *Neural network—a comprehensive foundation*. Prentice-Hall, Englewood Cliffs
- He Z, Wen X, Liu H, Du J (2014) A comparative study of artificial neural network, adaptive neuro fuzzy inference system and support vector machine for forecasting river flow in the semiarid mountain region. *J Hydrol* 509:379–386
- Hou LG, Xiao HL, Si JH, Xiao SC, Zhou MX, Yang YG (2010) Evapotranspiration and crop coefficient of *Populus euphratica* Oliv forest during the growing season in the extreme arid region northwest China. *Agric Water Manag* 97(2):351–356
- Hsu, C.W., Chang, C.C., Lin, C.J., (2010). A practical guide to support vector classification. URL <http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>.
- Jain A, Indurthy SKVP (2003) Comparative analysis of eventbased rainfall-runoff modeling techniques—deterministic, statistical, and artificial neural networks. *J Hydrol Eng ASCE* 8(2):93–98
- Jones, J.W., Ritchie, J.T., (1990) Crop growth models. Management of farm irrigation systems. In: Hoffman, G.J., Howel, T.A., Solomon, K.H. (Eds.), *ASAE Monograph No. 9*. ASAE, St. Joseph, Mich., pp.63–89
- Kisi O (2012) Least squares support vector machine for modeling daily reference evapotranspiration. *Irrig Sci*. doi:10.1007/s00271-012-0336-2
- Kisi O, Cengiz TM (2013) Fuzzy genetic approach for estimating reference evapotranspiration of Turkey: Mediterranean Region. *Water Resour Manag* 27(10):3541–3553
- Kisi O, Cimen M (2009) Evapotranspiration modeling using support vector machines. *Hydrol Sci J* 54(5):918–928
- Laaboudi A, Mouhouche B, Draoui B (2012) Neural network approach to reference evapotranspiration modeling from limited climatic data in arid regions. *Int J Biometeorol* 56(5):831–841
- Lin GF, Lin HY, Wu MC (2013) Development of a support-vector-machine-based model for daily pan evaporation estimation. *Hydrol Process* 22:3115–3127
- López-Urrea R, de Santa M, Olalla F, Fabeiro C, Moratalla A (2006) Testing evapotranspiration equations using lysimeter observations in a semi-arid climate. *Agric Water Manag* 85:15–26
- Müller K, Smola A, Rätsch G, Schölkopf B, Kohlmorgen J, Vapnik VN (1997) Predicting time series with support vector machines. *Artif Neural Networks—ICANN* 97(1327):999–1004
- Priestley CHB, Taylor RJ (1972) On the assessment of surface heat flux and evaporation using large scale parameters. *Mon Weather Rev* 100:81–92
- Principe JC, Euliano NR, Lefebvre CW (2000) *Neural and adaptive systems: fundamentals through simulations*. Wiley, New York
- Raghavendra NS, Deka PC (2014) Support vector machine applications in the field of hydrology: a review. *Appl Soft Comput* 19:372–386
- Rahimikhoob A (2014) Comparison between M5 model tree and neural networks for estimating reference evapotranspiration in an arid environment. *Water Resour Manag* 28(3):657–669
- Si JH, Feng Q, Zhang XY, Liu W, Su YH, Zhang YW (2005) Growing season evapotranspiration from *Tamarix ramosissima* stands under extreme arid conditions in northwest China. *Environ Geol* 48(7):861–870
- Tabari H, Kisi O, Ezani A, Hosseinzadeh Talaei P (2012) SVM, ANFIS, regression and climate based models for reference evapotranspiration modeling using limited climatic data in a semi-arid highland environment. *J Hydrol* 444–445:78–89
- Vapnik VN (1995) *The nature of statistical learning theory*. Springer, New York