# Comparison of M5 Model Tree and Artificial Neural Network for Estimating Potential Evapotranspiration in Semi-arid Climates

N. Ghahreman[a]*, M. Sameti[a]

*[a] Department of Irrigation and Reclamation Engineering, University of Tehran, Karaj, Iran*

**Abstract**

Evaporation is a fundamental parameter in the hydrological cycle. This study examines the performance of M5 model tree and artificial neural network (ANN) models in estimating potential evapotranspiration calculated by Penman- Monteith and Hargreaves- Samani equations. Daily weather data from two meteorological stations in a semi-arid climate of Iran, namely Kerman and Zahedan, were collected during 1995-2004 and included the mean, maximum and minimum air temperatures, dewpoint, relative humidity, sunshine hours, and wind speed. Results for both stations showed that the performance of the M5 model tree was more accurate (R=0.982 and 0.98 for Penman-Monteith equation and R=0.983 and 0.98 for Hargreaves-Samani equation in Kerman and Zahedan, respectively) than the ANN model (R=0.975 and 0.978 for Penman-Monteith equation and R=0.967 and 0.974 for Hargreaves-Samani equation in Kerman and Zahedan, respectively), but the models' differences were insignificant at a confidence level of 95%. It also performed better at the Zahedan station using the Penman-Monteith equation. The most significant variables affecting the potential evapotranspiration in the case of the Penman–Monteith equation were found to be mean air temperature, sunshine hours, wind speed, and relative humidity. Similarly, for the Hargreaves-Samani equation, the maximum and minimum temperatures, sunshine hours, and wind speed were determined to be the most significant variables. Further studies in other climates are recommended for further analysis.

*Keywords:* ANN, Machine learning, Penman-Monteith, Hargreaves-Samani

## 1. Introduction

Evapotranspiration, evaporation from soil and transpiration from vegetation, is an important component of the hydrologic cycle. In most arid regions, irrigation consumes the majority of developed water resources; moreover, water scarcity and misuse are substantial threats to the sustainability of agricultural production. Therefore, determining agricultural water demand is an important factor for developing a fundamental infrastructure and managing the allocation of water. The precise

quantification of crop evapotranspiration ($ET_c$) in irrigated agriculture is used to schedule irrigation and water resource management. Evapotranspiration from a reference surface in a standard condition (without stress) is called reference evapotranspiration ($ET_o$). Potential evapotranspiration can be measured or estimated by different methods. Using lysimeters is the most accurate method, but it is time-consuming and costly, and lysimeters are not widely available to all researchers, especially those in developing countries. Therefore, using alternate empirical models and equations might be preferable. These methods use physical equations (like the Penman-Monteith equation (P-M)) or of the empirical type using simple relationships of ET and meteorological variables (like the Hargreaves-

---
* Corresponding author. Tel.: +98 26 32241119,
  Fax: +98 26 32241119.
  *E-mail address:* nghahreman@ut.ac.ir

Samani equation). Many equations have been suggested for the estimation of potential evapotranspiration which, to some extent, are site- and climate-specific and should not be used in other climates without prior examination. Among these methods, the modified P-M equation proposed by FAO, the FAO56 P-M, has been widely accepted for use in different climates (Allen *et al*., 1998). Several ETo calculation packages, e.g., CROPWAT and SIMETAW, have been evaluated in different climates of Iran (Malekian *et al*., 2009; Ebrahimpour *et al*., 2014). The adequate, high quality data required for running the P-M equation is not available at all weather stations, so it is more preferable to apply equations such as the Hargreaves-Samani equation that use routine, readily available data to estimate $ET_o$ and require only daily mean, maximum and minimum temperatures, and extraterrestrial radiation (Hargreaves and Samani, 1985). Extraterrestrial radiation can be calculated theoretically (Drooger and Allen, 2002; Ghahreman and Bakhtiari, 2009). Recently, other approaches based on soft computing, like data mining, have been introduced for modeling, data classification, and clustering. These data-based approaches try to find unknown and hidden interrelations by performing certain iterative searches in a long series of data and finally provide the desired output. In other words, data mining is the identification of interesting structures in data, where structure designates patterns, statistical or predictive models of the data, and relationships among parts of the data (Fayyad and Uthurusamy, 2002). The skill of these techniques has also been proven in hydrology.

Data mining tools perform data analyses and may detect important data patterns contributing greatly to strategic decisions, knowledge bases, and relevant research including hydrology and soil-water management (Nagesh Kumar and Dhanya, 2009). Time series data mining, which combines chaos theory and data mining, can be effectively used in predicting river flood (Chaitanya and Yaclin, 2007). A cluster-based neural network model is effective in capturing nonlinear relationships among many hydrological processes (Parasuraman *et al*., 2007).

These relationships can be obtained through different algorithms, such as artificial neural networks (ANN), decision trees, regression trees, and model trees. ANN has been widely used in meteorological and hydrological studies such as drought prediction (Dastorani and

Afkhami, 2011) and wind speed prediction (Bakhtiari *et al*., 2013).

Evapotranspiration is dependent on several meteorological variables. A literature review revealed that recently, many studies have been carried out using data mining algorithms to model daily evaporation and evapotranspiration and relevant variables. Pal (2006) classified groundcover with an M5 model tree using data of Landsat 7 ($ETM^+$) and concluded that the accuracy of this model is higher than that of a decision tree. Terzi *et al*. (2005) used data from the Lake Egirdir region in Turkey to model daily pan evaporation and suggested the Kstar model as the best model among M5Rules and decision tree models. Using a genetic algorithm (GA), they showed that air temperature, water temperature, and relative humidity are the most significant variables. Terzi (2007) used an M5 model tree, artificial neural network, linear regression, and SMO regression to model daily pan evaporation and suggested the M5 model tree as the preferable model. Pal and Deswal (2009) used an M5 model tree to estimate evapotranspiration using the daily data from an automated weather station located at Davis, California during 1995-2005. The obtained values were compared with those calculated by FAO56 Penman-Monteith and Hargreaves-Samani equations. Based on their findings, the M5 model tree had the best performance. Moreover, sensitivity analysis revealed that solar radiation, mean temperature, humidity, and wind speed are the most significant variables in modeling potential evapotranspiration.

To date, few studies have been done on the application of data mining techniques in estimating potential evapotranspiration nationwide; therefore, the goal of this study was to use data mining algorithms to estimate potential evapotranspiration, make comparisons with other methods, and determine the most significant meteorological variables governing potential evapotranspiration.

## 2. Materials and Methods

Daily weather data from two study stations, Kerman (latitude 56º 58' E, longitude 30º 15' N, and an altitude of 1753.8 meters above MSL) and Zahedan station (latitude 60º 53' E, longitude 29º 28' N, and an altitude of 1370 meters above MSL) during 1995-2004 were used. Based on the extended-De Martonne classification (Khalili 1977), their climates are classified as arid cold and hyper-arid moderate, respectively. Daily recorded data of five

variables, including maximum temperature ($T_{max}$), minimum temperature ($T_{min}$), dewpoint temperature ($T_{dew}$), vapor pressure, and wind speed ($U_2$), covering 1995 to 2004, obtained from the Islamic Republic of Iran Meteorology Organization, were used to calculate potential evapotranspiration with the FAO-56 Penman-Monteith (FAO 56 P-M hereafter) and Hargreaves-Samani models. Reconstruction of data and filling gaps, where required, was implemented using a multivariate regression approach. The homogeneity of the data set was checked by a Run test.

For both applied models, 70% of the data were used for model building/training, and the remaining 30% was kept for model evaluation. Finally, the obtained estimations of $ET_o$ from both models were compared with corresponding values of Penman-Monteith and Hargreaves-Samani equations using statistical indices including root mean square error (RMSE), mean absolute error (MAE), and correlation coefficients. WEKA 3.6.4 and SPSS Clementine 12 packages were used to run the M5 model tree and the artificial neural network model, respectively.

## 2.1. Penman-Monteith FAO-56 equation

This method is the standard procedure when there is no measured lysimeter data (Allen *et al*., 1998; Alexandris, Kerkides, and Liakatas, 2006; Georgiou and Papamichail, 2008). Although in practice the best way to test the performance of the above-mentioned methods would be to compare their performances versus the lysimeter-measured data, this data set was not available for the study area. According to (Allen *et al*., 1998), the P-M method is summarized by the following equation:

$$ET_o = \frac{0.408\Delta(R_n - G) + \gamma[890/(T+273)]U_2(e_a - e_d)}{\Delta + \gamma(1+0.34U_2)} \quad [4]$$

Where $ET_o$ is the reference evapotranspiration (mmd$^{-1}$), $R_n$ is the daily net radiation (Mjm$^{-2}$d$^{-1}$), G is the soil heat flux (Mjm$^{-2}$d$^{-1}$), T is the mean daily air temperature at a 2 meters height (ºC), $U_2$ is the daily mean of the wind speed at a height of 2 m (ms$^{-1}$), $e_s$ is the saturation vapor pressure (kPa), $e_a$ is the actual vapor pressure (kPa), $\Delta$ is the slope of the saturation vapor pressure (KPaºC$^{-1}$), and $\gamma$ is the psychrometric constant (KPaºC$^{-1}$). All variables were calculated using procedures suggested by Allen *et al*. (1998). The soil heat flux (G) was neglected for the daily time scale (Allen *et al*., 2008). The Hargreaves equation (Hargreaves

and Samani, 1985; Hargreaves and Allen, 2007) can be written as:

$$ET_o = 0.0135(K_T).R_a.TD^{0.5}(T+17.8) \quad [5]$$

$$K_T = 0.00185(TD)^2 - 0.0433TD + 0.4023 \quad [6]$$

where $ET_o$ is reference evapotranspiration calculated by the HG method (mmd$^{-1}$), $K_T$ is an empirical coefficient, TD is the range of daily temperatures (ºC), Ra is the depth of water equal to the extraterrestrial radiation (mmd$^{-1}$), $T_{max}$ and $T_{min}$ are the daily maximum and minimum air temperatures (ºC), Tmean is the mean air temperature (ºC) computed as the average of $T_{max}$ and $T_{min}$.

## 2.2. Artificial neural network model

Neural networks are composed of simple elements operating in parallel. These elements are inspired by biological nervous systems. As in nature, the network function is determined largely by the connections between elements. A neural network can be trained to perform a specific function by adjusting the values of the connection weights between the elements. Commonly, neural networks are adjusted, or trained, so that a particular input leads to a certain output. The network is adjusted, based on a comparison of the output and the target, till the sum of square differences between the target and output values reaches a minimum value. Typically, many such input/target output pairs are used to train a network. Incremental training changes the weights and biases of a network as needed after presentation of each individual input vector. Neural networks have been used in various fields of application in different branches of sciences including pattern recognition, identification, classification, speech, vision, and control systems (Demuth and Beale, 2001).

## 2.3. M5 model tree

Model trees generalize the concepts of regression trees, which have constant values at their leaves (Witten and Frank, 2005). They are analogous to piece-wise linear functions (and hence nonlinear). An M5 model tree is a binary decision tree which has linear regression functions at the terminal (leaf) nodes which can predict continuous numerical attributes (Quinlan, 1992). Tree-based models are constructed by a divide-and-conquer method. Model tree generation requires two stages. The first stage involves using a splitting criterion to create a decision tree. The splitting criterion for the M5 model tree algorithm is based on

treating the standard deviation of the class values that reach a node as a measure of the error at that node and calculating the expected reduction in this error as a result of testing each attribute at that node. The formula to compute the standard deviation reduction (SDR) is:

$$SDR = sd(T) - \sum (\frac{|T_i|}{|T|} + sd(T_i)) \qquad [7]$$

where T represents a set of examples that reach the node, Ti represents the subset of examples that have the i[th] outcome of the potential set, and *sd* represents the standard deviation. Due to the splitting process, the data in child nodes have less standard deviation than those at the parent node and thus are more pure. After checking all the possible splits, M5 selects the one that maximizes the expected error reduction. This division often produces a large tree-like structure which may cause overfitting. To overcome the problem of overfitting, the tree must be cut back, for example, by replacing a subtree with a leaf. Thus, the second stage in the design of the model tree involves pruning the overgrown tree and replacing the sub-trees with linear regression functions. This technique of generating the model tree splits the parameter space into areas (subspaces) and builds in each of them a linear regression model.

## 3. Results

In Table 1, the performances of the M5 model tree and artificial neural network in estimating potential evapotranspiration by both equations, are compared. At the Kerman station, the M5 model tree estimated evapotranspiration with a correlation coefficient of 0.98 and a mean absolute error of 0.270, but with the ANN method, the correlation coefficient and mean absolute error were 0.975 and 0.33, respectively. At the Zahedan station, R and MAE values with the M5 model tree were 0.9832 and 0.252, and with the ANN model they were 0.978 and 0.3, respectively. According to the results of the student's t-test ($P_{\alpha=0.05}$=0.914, t= -0.11 for the Kerman station and $P_{\alpha=0.05}$=0.973, t= -0.03 for the Zahedan station), there was no significant difference between the M5 model tree and the ANN model at a confidence level of 95%. The performance of both models in estimating potential evapotranspiration was higher at the Zahedan station. To determine the most important variables, sensitivity analyses were applied by using all input variables, and the importance of every variable was examined by removing each of them from the results. At both stations, the most common variables affecting potential evapotranspiration were average daily temperature, sunshine hours, average wind speed, and average relative humidity. At the Kerman station, dewpoint temperature and actual vapor pressure were the sensitive parameters in artificial neural networks, and at the Zahedan station, actual vapor pressure was the sensitive parameter in the M5 model tree. It can be concluded that the three parameters of mean air temperature, sunshine hours, and wind speed had the most influence on the results, because by removing them, the performance decreased significantly.

Table 1. Comparison of ANN and M5 model tree for estimating potential evapotranspiration of FAO 56 Penman Monteith method for two study stations during the period of 1995-2004.

| Stations | Model | R | MAE | Combination of effective variables |
|---|---|---|---|---|
| Kerman | M5 | 0.982 | 0.270 | T,n,w,RH |
| | ANN | 0.975 | 0.33 | T,n,w,RH,dew,e |
| Zahedan | M5 | 0.983 | 0.252 | T,n,w,RH,e |
| | ANN | 0.978 | 0.3 | T,n,w,RH |

T, mean air temperature; n, sunshine hours; w, wind speed; RH, relative humidity; dew, dewpoint temperature; e, vapor pressure.

Table 2. Sensitivity analysis of M5 model tree for Kerman station

| M5 model tree | R | RMSE | MAE |
|---|---|---|---|
| T,n,w,RH,dew,e,R | 0.983 | 0.341 | 0.253 |
| T,n,w,RH,dew,e | 0.983 | 0.341 | 0.254 |
| T,n,w,RH,dew | 0.983 | 0.343 | 0.254 |
| * T,n,w,RH,e | 0.983 | 0.339 | 0.252 |
| T,n,w,dew,e | 0.982 | 0.347 | 0.257 |
| T,n,RH,dew,e | 0.92 | 0.73 | 0.555 |
| T,n,w,RH | 0.983 | 0.344 | 0.255 |
| T,n,w,e | 0.982 | 0.348 | 0.257 |
| T,n,w,dew | 0.983 | ۰,۳۴۶ | 0.256 |
| T,n,w | 0.98 | 0.37 | 0.275 |

Table 3. Sensitivity analysis of M5 model tree for Zahedan station

| M5 model tree | R | RMSE | MAE |
|---|---|---|---|
| T,n,w,RH,dew,e,R | 0.982 | 0.363 | 0.273 |
| T,n,w,RH,dew,e | 0.982 | 0.363 | 0.272 |
| T,n,w,RH,dew | 0.982 | 0.362 | 0.271 |
| T,n,w,RH,e | 0.982 | 0.363 | 0.272 |
| T,n,w,dew,e | 0.981 | 0.364 | 0.272 |
| T,n,RH,dew,e | 0.93 | 0.71 | 0.56 |
| ∗ T,n,w,RH | 0.982 | 0.36 | 0.27 |
| T,n,w,dew | 0.981 | 0.364 | 0.272 |
| T,n,w | 0.98 | 0.373 | 0.285 |

When the Hargreaves–Samani method is used, the accuracy of the models increases using maximum and minimum air temperature instead of mean temperature. In Table 4, the performances of the M5 model tree and the artificial neural network method for estimating potential evapotranspiration of the Hargreaves–Samani method at both stations are compared. It can be seen that the M5 model tree has a higher correlation coefficient (R=0.983 and 0.98 for Kerman and Zahedan, respectively) and a smaller mean absolute error (MAE=0.496 and 0.423 for Kerman and Zahedan, respectively) compared with the artificial neural network model (R=0.967, MAE=0.671 and R=0.974 , MAE= 0.471 for Kerman and Zahedan, respectively), but the differences were not significant using the student's t-test.

$(P_{\alpha=0.05}=0.98$, t= -0.03 for the Kerman station and $P_{\alpha=0.05}=0.898$, t=0.13 for the Zahedan station) at a confidence level of 95%. A comparison of Table 1 and Table 4 indicates that values of mean absolute error in the Penman-Monteith equation are less than those in the Hargreaves-Samani equation. With the Hargreaves–Samani method, the most significant variables affecting potential evapotranspiration were mean, maximum and minimum temperatures, sunshine hours, and wind speed. At the Kerman station, relative humidity was found to be the most important variable in the artificial neural network model. Similarly, at the Zahedan station for both models, the mean actual vapor pressure was the major affecting variable.

Table 4. Comparison of ANN and M5 model trees for estimating potential evapotranspiration of Hargreaves-Samani method for two stations of Kerman and Zahedan during the study period 1995-2004

| Stations | Model | R | MAE | Combination of effective variables |
|---|---|---|---|---|
| Kerman | M5 | 0.983 | 0.496 | Tmax,Tmin,n,w |
| | ANN | 0.967 | 0.671 | Tmax,Tmin,n,w,RH |
| Zahedan | M5 | 0.980 | 0.423 | Tmax,Tmin,n,w,e |
| | ANN | 0.974 | 0.471 | Tmax,Tmin,n,w,e |

Tmax, maximum temperature; Tmin, minimum temperature.

One of the advantages of the M5 model tree is that this model allows access to a combination of several simple linear relationships that can be used for predicting potential evapotranspiration. From Figure 1 it can be concluded that the most important variable in constructing an M5 model tree is average daily temperature that branches the

tree in T>17.95 and T≤17.95. After that, according to splitting ratio in T>17.95 the wind speed and in T≤17.95 the average daily temperature are the most important variables. Branching continued until 19 linear models were obtained for the Kerman station dataset, each of which should be applied for its specific conditions.
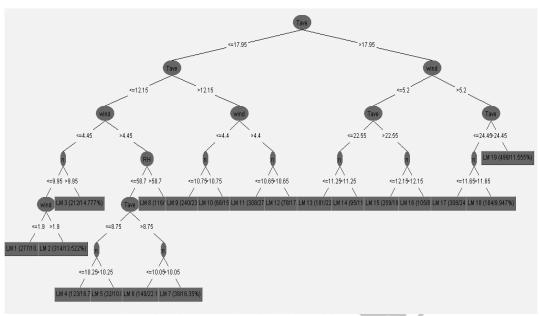
Fig. 1. Schematic representation of derived M5 model tree for Kerman station

### 3.1. Regression linear models for Kerman station

LM 1: ET-PM = 0.1784 * wind + 0.0279 * Tave + 0.0041 * RH + 0.0056 * n+ 0.1576

LM 2: ET-PM = 0.1494 * wind + 0.0525 * Tave – 0.0027 * RH + 0.0189 * n + 0.2936

LM 3: ET-PM = 0.1936 * wind + 0.0636 * Tave + 0.0056 * RH + 0.8104 * n – 8.0414

LM 4: ET-PM = 0.0876 * wind + 0.0698 * Tave – 0.0083 * RH + 0.0405 * n + 0.7028

LM 5: ET-PM = 0.0957 * wind + 0.0935 * Tave – 0.0046 * RH + 0.6019 * n – 5.3738

LM 6: ET-PM = 0.1077 * wind + 0.0936 * Tave – 0.0097 * RH + 0.0603 * n + 0.4363

LM 7: ET-PM = 0.0973 * wind + 0.1128 * Tave – 0.0135 * RH + 0.3047 * n – 1.8025

LM 8: ET-PM = 0.0751 * wind + 0.0817 * Tave – 0.0054 * RH + 0.0667 * n + 0.3939

LM 9: ET-PM = 0.2447 * wind + 0.1602 * Tave + 0.0105 * RH + 0.0561 * n – 1.8327

LM 10: ET-PM = 0.2569 * wind + 0.1178 * Tave + 0.018 * RH + 0.7122 * n – 8.1259

LM 11: ET-PM = 0.1478 * wind + 0.1881 * Tave – 0.007 * RH + 0.0816 * n – 1.1996

LM 12: ET-PM = 0.1666 * wind + 0.164 * Tave – 0.0007 * RH + 0.5194 * n – 5.6033

LM 13: ET-PM = 0.2493 * wind + 0.177 * Tave + 0.0096 * RH + 0.0713 * n – 2.1323

LM 14: ET-PM = 0.2753 * wind + 0.0993 * Tave + 0.0274 * RH + 0.5917 * n – 6.5035

LM 15: ET-PM = 0.2749 * wind + 0.1622 * Tave + 0.001 * RH + 0.1197 * n – 1.9257

LM 16: ET-PM = 0.3032 * wind + 0.145 * Tave + 0.0156 * RH + 0.4873 * n – 6.1078

LM 17: ET-PM = 0.2036 * wind + 0.1763 * Tave – 0.0131 * RH + 0.0874 * n – 1.0811

LM 18: ET-PM = 0.2403 * wind + 0.1333 * Tave – 0.008 * RH + 0.3852 * n – 3.7715

LM 19: ET-PM = 0.2893 * wind + 0.1446 * Tave – 0.0102 * RH + 0.1523 * n – 1.5047

### 4. Discussion

In this study, the performances of an M5 model tree and artificial neural network modes were evaluated for estimating potential evapotranspiration using FAO-56 Penman-Monteith and Hargreaves–Samani methods for the two stations of Kerman and Zahedan. The results indicated no significant difference between the M5 model tree and the artificial neural network model in either station for the two chosen methods of Penman–Monteith and Hargreaves-Samani using the student's t-test. For the Hargreaves-Samani method, the most significant affecting variables were mean air temperature, sunshine hours, wind speed, and relative humidity. In the case of the Penman-Monteith equation, the most important variables were average maximum and minimum temperatures, sunshine hours, and wind speed. Comparing the results of the two study stations, it was found that the results of the Zahedan station are more accurate. This might be due to the strong moisture advection in the Kerman region which could lead to the underestimation

of ET values. The main advantage of an M5 model tree is that it offers several multilinear regression equations which are valid for certain climatic conditions. The first branches of the tree are identified as the most important variables cause the tree to branch. Pal and Deswal (2009) showed that solar radiation, mean air temperature, relative humidity, and wind speed are the most significant parameters in estimating potential evapotranspiration with an M5 model tree. Terzi (2007) used several algorithms to estimate evaporation and concluded that there is a close agreement between results of an M5 model tree and measured daily pan evaporation values. He used air temperature, water temperature, solar radiation, and relative humidity parameters as effective parameters of evaporation. Sattari (2013) suggested that the ANN model performs better than the M5 model tree with the dataset from Ankara, Turkey. He concluded that a combination of parameters including minimum and maximum temperatures, minimum and maximum relative humidity, wind speed, and sunshine hours is the best combination of variables for calculating ETo. Further research might be recommended using lysimetric data for more scrutiny and suggesting preferable models.

## References

Allen, R.G., L.S. Pereira, D. Rae, M. Smith, 1998. Crop evapotranspiration, Irrigation and Drainage Paper No. 56. Food and Agriculture Organization: Rome, Italy.

Allen, R.G., I.A. Walter, R.L. Elliott, T.A. Howell, D. Itenfisu, M.E. Jensen, R.L. Snyder, 2005. The ASCE standardized reference evapotranspiration equation. Task Committee on Standardization of Reference Evapotranspiration of the EWRI of the ASCE, USA.

Alexandris, S., P. Kerkides, A. Liakatas, 2006. Daily reference evapotranspiration estimates by the "Copais" approach. AGRICULTURAL WATER MANAGEMENT, 82;371-386.

Bakhtiari, B., N. Ghahreman, I. Rahimi, 2013. Application of neural network for prediction of wind speed in short time scale (Case study: Jiroft Station). SOIL AND WATER RESEARCH, 44(1); 11-20 (In Farsi).

Chaitanya, D., A. Yalcin, 2007. Flood prediction using time series data mining. HYDROLOGY, 333; 305-316.

Demuth, H., M. Beale, 2001. Neural network toolbox user's guide. Vol.4: The Math Works, Inc., Natick, Massachusetts, USA.

Droogers, P., R.G. Allen, 2002. Estimating reference evapotranspiration under inaccurate data conditions. IRRIGATION AND DRAINAGE SYSTEMS, 16; 33–4.

Dastorani, M.T., H. Afkhami, 2011. Application of artificial neural networks on drought prediction in Yazd (Central Iran). DESERT, 16 (1); 39-48.

Ebrahimpour, M., N. Ghahreman, M. Orang, 2014. Assessment of climate change impacts on reference evapotranspiration and simulation of daily weather data using SIMETAW. IRRIGATION AND DRAINAGE ENGINEERING, 140(2); 10.1061/(ASCE)IR.1943-4774.0000669.

Fayyad, U.M., R. Uthurusamy, 2002. Evolving data mining in to solutions for insights Communication ACM. 45(8); 34 pp.10.

Georgiou, P.E., D.M. Papamichail, 2008. Optimization model of an irrigation reservoir for water allocation and crop planning under various weather conditions. Irrigation Science, 26; 487-504.

Ghahreman, N., B. Bakhtiari, 2009. Solar radiation estimation from rainfall and temperature data in arid and semi-arid climates of Iran. DESERT, 14(2); 141-150.

Hargreaves, G.H., R.G. Allen, 2003. History and evaluation of Hargreaves evapotranspiration equation .IRRIGATION AND DRAINAGE ENGINEERING, 129; 53-63.

Hargreaves, G.H., Z.A. Samani, 1985. Reference crop evapotranspiration from temperature. APPLIED ENGINEETING IN AGRICLUTURE, 1; 96- 99.

Khalili, A., 1997. Integrated water plan of Iran. Vol.4: Meteorologycal studies, Ministry of Power, Iran.

Malekian, A., H. Ghasemi, A. Ahmadian, 2009. Evaluation of the efficiency of CROPWAT model for determining plant water requirement in arid regions. DESERT, 14 (2); 209-215.

Nagesh Kumar, D., C.T. Dhanya, 2009. Data mining and its applications for modelling rainfall extremes, HYDRAULIC ENGINEERING, 15(1): 25-51.

Pal, M., 2006. M5 model tree for land cover classification. Int J Remote Sens. 27(4): 825– 831.

Pal, M., S. Deswal, 2009. M5 model tree based modeling of reference evapotranspiration. HYDROLOGYCAL PROCESSES, 23; 1437-1443.

Parasuraman, K., A. Elshorbagy, S. Carey, 2007. Modelling the dynamics of the evapotranspiration process using genetic programming. HYDROLOGYCAL SCIENCES, 52(3); 563-578.

Quinlan, J.R., 1992. Learning with continuous classes. Proceedings of Australian Joint Conference on Artificial Intelligence. In: Hobart, asmania.Australia. 16-18 Nov.1992. World Scientific Press, Singapore. pp. 343–348.

Sattari, M., M. Pal, K. Yürekli, A. Ünlükara, 2013. M5 Model trees and neural network of $ET_o$ in Ankara, Turkey. ENGINEERING AND ENVIROMENTAL SCIENCES. 37; 211-219.

Terzi, Ö., E.U. Kücüksille, M.E. Keskin, 2005. Modelling of daily pan evaporation using data mining. Proc International Symposium on Innovation in Intelligent Systems and Applications. Pp. 182- 185.

Terzi, Ö., 2007. Data mining approach for estimation evaporation from free water surface. APPLIED SCIENCES, 7(4); 593- 596.

Witten, I.H., E. Frank, 2005. Data Mining: Practical machine learning tools and techniques with Java implementations. Morgan Kaufmann, San Francisco, USA.