**RESEARCH**

# Random Forest and Multilayer Perceptron hybrid models integrated with the genetic algorithm for predicting pan evaporation of target site using a limited set of neighboring reference station data

Sadra Shadkani[1] · Sajjad Hashemi[1] · Amirreza Pak[1] · Alireza Barzgari Lahijan[2]

## Abstract

This study explores the application of machine learning algorithms for the prediction of pan evaporation (Ep), which is a critical factor in water resource management for the assessment of water demand and usage. Specifically, this research evaluates the effectiveness of two base models: Random Forest (RF) and Multi-Layer Perceptron (MLP) and their optimized counterparts using a Genetic Algorithm (GA), designated as GA-RF and GA-MLP, for modeling Ep at a target station using data from adjacent stations. The datasets were split into a training set (70%) and a testing set (30%). The models' performances were judged using three statistical measures: Correlation Coefficient (CC), Scattered Index (SI), and Willmott's Index of agreement (WI). The enhanced models, particularly GA-MLP-5, showed superior performance with a CC of 0.8704, SI of 0.2539, and WI of 0.9212, indicating the potent ability of GA to refine RF and MLP models for predictive accuracy. Additionally, sensitivity analysis via the GA-RF indicates the varying influence of Ep from neighboring stations on the target station, shedding light on key predictors for effective water management. Conclusively, this study demonstrates that the hybrid models have significant potential in accurate Ep estimation and can be expanded to predict other meteorological variables, offering valuable tools for water resource management strategies.

**Keywords** Pan evaporation · Random forest · Scattered Index · Willmott's Index · Machine learning

## Introduction

Evaporation is an important element in the hydrological cycle and is how the water changes into the vapor then get into the air. This process needs energy absorption to change in vapor pressure (Wu et al. 2020; Sebbar et al. 2019).

✉ Sadra Shadkani
  Sshadkani@gmail.com

  Sajjad Hashemi
  hashemisajjad2009@gmail.com

  Amirreza Pak
  amirrezapakap@gmail.com

  Alireza Barzgari Lahijan
  alireza.barzgarilahijan@gmail.com

[1] Department of Water Engineering, Faculty of Agriculture, University of Tabriz, Tabriz, Iran

[2] Department of Civil Engineering, Faculty of Civil Engineering, University of Tabriz, Tabriz, Iran

Evaporation is one of the most intricate factors in the hydrologic cycle to be analyzed because of the complicated connection among water, land, and atmospheric systems. Therefore, the estimation of evaporation is a crucial part of agriculture and water resources management, and it attracts the attention of scholars all around the world (Feng et al. 2018; Fan et al. 2016; Gundalia and Dholakia 2013; Adnan et al. 2019; Yaseen et al. 2020). Several methods exist to predict the evaporations, such as Penman method, mass transfer method, energy, water balance technique, and pan evaporation (Ep). The pan evaporation technique is utilized worldwide because it is affordable and has an uncomplicated operation (Keshtegar et al. 2016). However, different factors like debris in the water, pan size, water depth in the pan, materials used to build the pan, animal activity in and around the pan can create errors in the measurement of Ep (Piri et al. 2009).

Furthermore, wind speed, the temperature of soil and air, sunshine, relative humidity, the pressure of the atmosphere, and vapor have an indispensable influence on evaporation. For instance, when the radiation and temperature

increase, the rate of evaporation will rise. Likewise, wind speed contributes to evaporation by removing water from water bodies. However, the effect of the parameters, as mentioned earlier, is not well recognized on Ep estimation. Hence, the estimation of evaporation is considered of great essence. Direct and indirect techniques are utilized for predicting evaporation (Kisi 2015; Allawi and El-Shafie 2016). However, the usage of evaporation pans can reduce Ep measurement accuracy due to instruments' restrictions in various circumstances (heavy rainfall or strong wind) and practical difficulties (Wang et al. 2017; Ghaemi et al. 2019). In the indirect methods, several climatic variables such as solar radiation, wind speed, and air temperature are required to estimate evaporation, but, in the specific regions, these elements are not attainable. Moreover, developing a mathematical relationship that encompasses all factors that affect evaporation, owing to the nonlinear and complex system of predicting evaporation, is almost impossible. Therefore, numerous hydrologic models are suggested by researchers to predict evaporation (Wang et al. 2017).

Various machine learning algorithms, including support vector machine (SVM), random forests (RF), adaptive neuro-fuzzy inference system (ANFIS) (Haddadi et al. 2022), multivariate adaptive regression splines (MARS), gene expression programming (GEP) (Birbal et al. 2021; Chaplot 2021), artificial neural network (ANN) (Mohammadi 2023; Jayathilake et al. 2023), and gradient boosted decision tree (GBDT) are used to estimate Ep and other hydrological parameters (Rahimikhoob 2009; Abghari et al. 2012; Keshtegar et al. 2016; Kisi et al. 2016; Yu et al. 2016; Yang et al. 2017; Xu et al. 2017; Chen et al. 2017; Wang et al. 2017; Behrooz et al. 2019; Zhu et al. 2009, 2019a, b). These models created better outcomes in comparison with the empirical approaches. Kisi (2009) used available climatic data to model the evaporation procedure by using three various ANN methods and revealed that the radial basis neural network (RBNN) and MLP methods could be utilized effectively for this purpose. Piri et al. (2009) used the ANN model, which is optimized by autoregressive external input (ARX) and assessed the model to estimate the Ep parameter at Southeast of Iran. Based on the result of the study, the optimized NNARX model performed better than ANN and Marciano models. Furthermore, the model with vapor and wind pressure inputs has higher precision than dew point and temperature. Lin et al. (2013) assessed the precision of two methods, including the SVM approach and a back-propagation network (BPN), in predicting the monthly evaporation and found that the SVM method produced more accurate results than other similar methods. In order to predict the daily Ep, the cascade correlation neural networks (CCNN) and the MLP model are utilized by Kim et al. (2014). The results

showed the superiority of the CCNN model when compared with MLP for both heterogeneous and homogeneous meteorological stations. Malik et al. (2017) evaluated the performance of the self-organizing map neural network, CANFIS, MLP, and RBFNN models. These models were specifically applied to estimating the monthly Evapotranspiration (Ep) parameter in India, with studies conducted at Pantnagar and Ranichauri stations. Using six input meteorological elements, the MLP and CANFIS models had the least error compared to other models. Feng et al. (2018) developed three machine learning models to predict the Ep and evaluate their performance. The results showed that the ELM method has superior performance compared with PSO-ANN (ANN optimized by particle swarm algorithm) and GA-ANN (ANN optimized by the genetic algorithm). Lu et al. (2018), in addition to using three methods of RF, GBDT, and M5, implemented four empirical models for estimating Ep in the lake of Poyang. They discovered that machine learning techniques had a superior performance than the empirical models. Among the mentioned techniques, the GBDT created accurate results than others. Majhi et al. (2019) investigated the efficiency of deep neural networks (DNN) in estimating the daily Ep in India, which showed that this method had better performance than empirical models and MLP. In order to estimate monthly, Ep, Kisi and Heddam (2019) used MARS and M5 with temperature data as inputs and utilized different splitting tactics for each model. The results disclosed that the performance of the MARS method was higher than M5, and increasing the amount of data increases its accuracy. Tree-based (M5 tree, random forests) machine learning techniques can be used for evaporation prediction due to their simple and strong nature (Alipour et al. 2014; Hassan et al. 2017). These models are popular because of their ability to compute large datasets (Hassan et al. 2017). Combining various methods to construct the hybrid models obtained the attention of researchers in the field of hydrology. Owing to the specific characteristics of each method, the hybrid models can increase the precision of models. For instance, Nourani et al. (2019) used the hybrid Wavelet-M5 for modeling the Suspended Sediment Load. The outcomes indicated that the performance of the hybrid model is better than the individual ANN and M5 models. A hybrid wavelet-linear genetic programming (WLGP), artificial neural network (ANN), Multi Linear Regression (MLR), LGP, and a hybrid wavelet-ANN (WANN) models are implemented to streamflow prediction by Ravansalar et al. (2017) in two stations of Pataveh and Shahmokhtar at Beshar River. They compared the outcomes of different models and based on the obtained results. The WLGP model increased monthly streamflow estimation precision in both Beshar River stations (Iran). Other studies have shown the superiority of combined methods over individual methods. For example,

Chaudhary et al. (2020) demonstrated the success of using the Ensemble Particle Swarm Optimization (EnsPSO) method for vegetable crop disease recognition. Chaudhary et al. (2016a) introduced an improved-RFC (Random Forest Classifier) technique for multiclass crop disease classification problems and showed the assumed model's optimal performance. Chaudhary et al. (2016b) integrated supervised instance filter- Resample and Gain Ratio feature ranking techniques and then combined the calculations of Naïve Bayes and Logistic Regression utilizing ensemble-Vote for oilseed disease classification. The results proved the success of the new hybrid model.

The previous research gaps and the current research objectives are as follows:

1- Limited studies on the integration of RF and MLP models with GA for such specific predictions.
2- Inadequate research using limited reference station data to predict pan evaporation for a different target site.
3- Lack of comparative analysis to understand the benefits of using hybrid models over traditional or singular model approaches.
4- Integration of Hybrid Models: Current research lacks in-depth investigations into the integration of RF and MLP models within a hybrid framework for the prediction of pan evaporation. This hybrid approach may leverage the strengths of both individual models and improve predictive accuracy.
5- Use of Genetic Algorithms: There is a scarcity of studies that have systematically integrated GAs with hybrid RF and MLP models for optimizing predictive performance in hydrological simulations.
6- Data Limitation Challenges: Previous studies have predominantly relied on extensive and complete datasets from the same geographical location as the target prediction site. This research contributes to the field by utilizing a limited set of neighboring reference station data to predict pan evaporation at the target site, an approach not widely addressed in existing models.
7- Generalization across Different Climates: Few studies have explored the adaptability and accuracy of hybrid models across different climatic conditions, especially when using data from reference stations with diverse characteristics.
8- Generalization across Different Climates: Few studies have explored the adaptability and accuracy of hybrid models across different climatic conditions, especially when using data from reference stations with diverse characteristics.

We aim to bridge these gaps by providing a comprehensive analysis of the model's performance and demonstrating its potential advantages over conventional methods, thereby offering new insights and directions for future research in hydrological predictions.

Evaporation is one of the main processes in nature's water cycle and one of the most important factors in agricultural, hydrological and meteorological studies, reservoir operation, irrigation scheduling, and water resources management. Therefore, accurate estimation of evaporation has great importance in hydrological studies. According to the researches, the proper performance of machine learning methods in estimating hydrological parameters is quite evident. On the other hand, hybrid models have improved the performance of single models. So, to summarize and evaluate the previous research, it can be said that most of the researches did not have valid results. Most meteorological parameters dependent on evaporation were used, and hybrid methods were rarely used to predict evaporation. In this study, the Ep data of adjacent stations were also used to predict the target station Ep, which can be used without Ep-dependent data. Due to the simpler structure and better performance compared to experimental equations and other known machine learning methods for predicting Ep and the coupling with GA to upgrade the models, MLP and RF methods were used. This paper aims to provide high-precision modeling and investigate the applicability of the RF and MLP models for predicting pan evaporation. Moreover, to increase the accuracy of the stated models, RF and MLP are hybridized with a Genetic Algorithm (GA) to create new RF-GA and MLP-GA models for estimating Ep values in seven stations of Iran. The hybrid RF-GA models have not been utilized for evaporation estimation to the best of our knowledge.

## Materials and methods

### Study area

With an area of about 17,800 square kilometers, Ardabil province occupies about 1.1 percent of the total area of Iran. This province is located in northwestern Iran, which shares borders with the Republic of Azerbaijan from the north and Gilan, East Azerbaijan and Zanjan provinces from the east, west and south, respectively. The geographical coordinates in this province are in the range of 37° 45′ to 39° 42′ North latitude, and 47° 30′ to 48° 55′ East longitude, and its average altitude is 2,400 meters above sea level. So that its lowest point with a height of 100 meters is in Parsabad and Bileh Savar cities, and its highest point is Sabalan Mountain with a height of 4,811 meters. Ardabil province in the axis of longitude (with an expansion of 1 degree and 35 minutes), along with the height factor of its plains and mountains with a matched combination, is adjacent to the Caspian Sea

and a large extent in the north-south direction in latitude (2 degrees and 31 minutes) has given a lot of climate diversity to Ardabil province. About two-thirds of it has a mountainous texture with a large height difference, and the rest are flat and low areas. So that the north of this province with a low altitude has a relatively warm climate, and the central and southern regions have a cold mountainous climate. Also, the special geographical and topographic features of the province, such as mountain ranges with a height of more than 4,000 meters and vast plains, have caused this province to be in a better position than other regions of the country in terms of precipitation so that the province's precipitation index varies between 250 and 600 millimeters. In the present study, pan evaporation data from 7 meteorological stations of Ardabil province, including Ardabil, Sarein, Nir, Bileh Savar, Meshgin Shahr, Parsabad and Khalkhal have been used in 10 years (2008–2018) on a daily scale. Figure 1 shows the study area. Also, the statistical characteristics

of the evaporation data (Table 1) is given below. According to Table 1, among the studied cities and in terms of average evaporation value, Ardabil city with 6.3 mm/day has the highest value and Sarein city with 5.1 mm/day has the lowest average evaporation rate. Furthermore, the maximum amount of evaporation has been recorded for Parsabad station with 46 mm/day value in the studied time period.

## Random Forest (RF)

Breiman (2001) suggested the random forest method, a non-parametric statistical technique to make a prediction, and this technique uses a large number of classifications in the estimation process. The RF is part of an ensemble learning system that forms and merges multiple learners to reach the best generalization abilities (Prasad et al. 2006). In this method, the CART, one of the decision tree algorithms, is used as the base learner. Among various rule

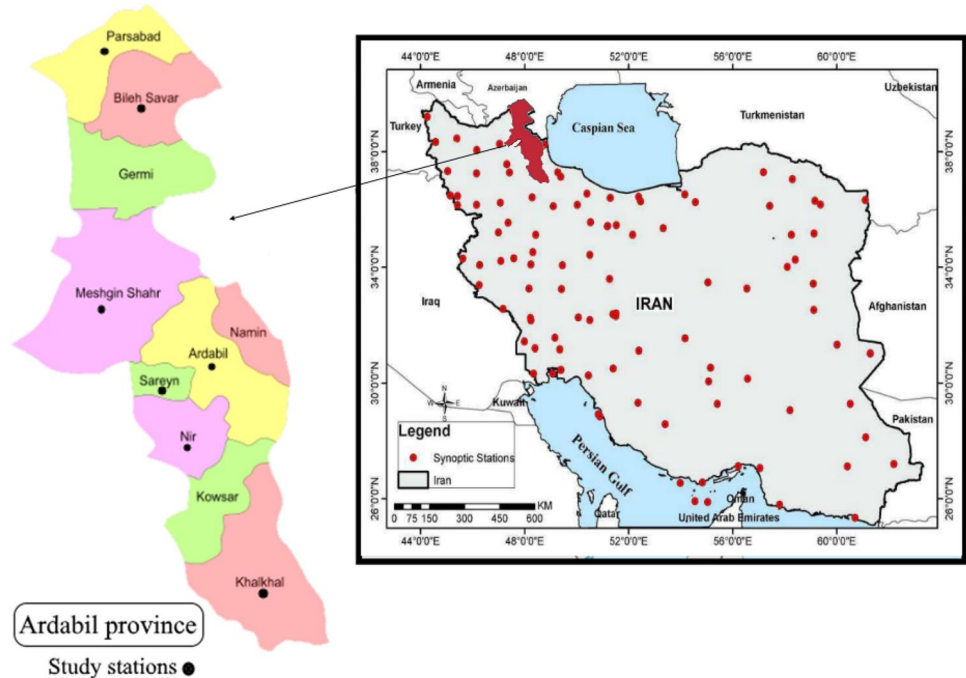**Fig. 1** Locations of studied area stations in the Iran



**Table 1** Statistical characteristics of the evaporation data

| Stations | Mean | Minimum | Maximum | Standard deviation | Coefficient of variation | Skewness |
|---|---|---|---|---|---|---|
| Ardabil | 6.3 | 0.0 | 21.1 | 3.259 | 0.513 | 0.405 |
| Sarein | 5.1 | 0.0 | 15.0 | 2.908 | 0.563 | 0.105 |
| Nir | 5.5 | 0.0 | 15.4 | 2.652 | 0.478 | 0.028 |
| Bileh Savar | 5.4 | 0.0 | 14.4 | 3.139 | 0.586 | 0.309 |
| Meshgin Shahr | 5.2 | 0.0 | 13.8 | 2.839 | 0.547 | 0.161 |
| Parsabad | 5.8 | 0.0 | 46.0 | 3.187 | 0.549 | 1.133 |
| Khalkhal | 6.198 | 0.0 | 16.0 | 2.749 | 0.444 | 0.128 |

generation techniques, the RF-based model has efficient and high performance, and it is a more vigorous way compared to other decision tree ensembles. Moreover, these rules can reveal the whole decision-making procedure. The RF working method is as follows: random subsets from original data are selected to build each decision tree consisting of two-thirds of all data (training process). The residual datasets, known as out-of-bag data (OOB), are utilized for specialized problems.

Furthermore, the lowest Gini index is implemented to choose the best split for variables at each node. By and large, the averaged aggregations indicated the RF model results. The datasets are split frequently until the predefined situation is completed. In the RF modeling, three factors should be defined, including the minimum number of nodes, the number of variables utilized to mature any tree (m try), and the number of trees in the forest (n tree) (Nawar and Mouazen 2017). The mtry factor appoints the connection between trees and the resistance of every individual tree. By reducing relations between trees and rising the tree's solidity, the RF model's proficiency can be enhanced (Ließ et al. 2012).

## Genetic Algorithm (GA)

Genetic Algorithms (GA), developed by Holland (1992) and Goldberg and Holland (1989), are presumptive algorithms (Schwefel 1993). This method emulates Darwin's evolution regulation based on natural selection and is an efficient way to solve difficult problems. In this algorithm, a population of individuals, based on their strength, ability, and desirability, remain alive and can proliferate more than other individuals. After several generations, individuals with the best efficiency will be born. Moreover, any individual is stated as a chromosome (chromosomes are completed during various generations) and the chromosomes are nominated for the problem's solution. In addition, chromosomes have a fixed number of genes and the binary coding is a popular method for demonstrating them. The specific population is comprised of a series of chromosomes and the genetics operators can affect the population and form the new population with an identical number of chromosomes. A fitness function is used to assess the quality of the solution for every chromosome. The parents for the next population are recognized by the fitness function. To attain the combination of genes that maximize or minimize the fitness function, the GA utilizes natural operators such as selection, crossover, and mutation (Holland 1992). Crossover is exerted on two chromosomes from the parents. In this process, the first part of one chromosome is utilized as the second part in the other one (binary string is recognized accidentally and cutting each of chromosomes into two parts then exchange them with crossing action). After the crossover section, to produce the randomness in the solution region the mutation is implemented. Consequently, appropriate parents are selected to create a novel population. Then, this procedure happened again for various sets of populations. In the context of a study involving the use of a Genetic Algorithm (GA) integrated with Random Forest (RF) and Multilayer Perceptron (MLP) models for predicting pan evaporation, GA parameters would play a critical role in the optimization process. Choosing appropriate boundaries for these parameters is essential for the GA to effectively search the solution space and converge to an optimal or near-optimal solution. Here's an explanation of how the GA parameters could be determined in this study:

Population Size: The number of potential solutions (individuals) in each generation. A larger population may offer a more diverse genetic pool, but it also requires more computational resources. The lower and upper boundaries are determined based on the complexity of the problem and available computational power.

Crossover Rate: This parameter determines how often a crossover (mixing of genetic material from two parents) will occur to create a new offspring. The rate should be high enough to allow for sufficient exploration of the solution space but not too high to prevent premature convergence. Typically, the range is between 60% and 90%.

Mutation Rate: The mutation rate dictates how often a mutation will alter a given gene. It ensures genetic diversity and helps prevent the algorithm from becoming stuck in local optima. A common range is between 0.5% and 1%, however, this can be adjusted based on preliminary runs to avoid excessive randomness which can lead to a loss of good solutions.

Selection Method: While not a parameter with a numerical value, the selection method determines which individuals will be selected for reproduction. Methods like tournament selection, roulette wheel selection, or elite selection can be chosen based on which performs best during initial experimentation.

Elitism: This involves carrying over a certain number of the best individuals to the next generation without alteration. Deciding on the number of elite individuals typically involves testing to see how it affects the balance between exploration and exploitation.

Number of Generations: The total number of iterations the GA will perform. More generations can lead to a better solution but also a longer computation time. The range is often set based on when incremental improvements plateau, which can be determined through exploratory runs.

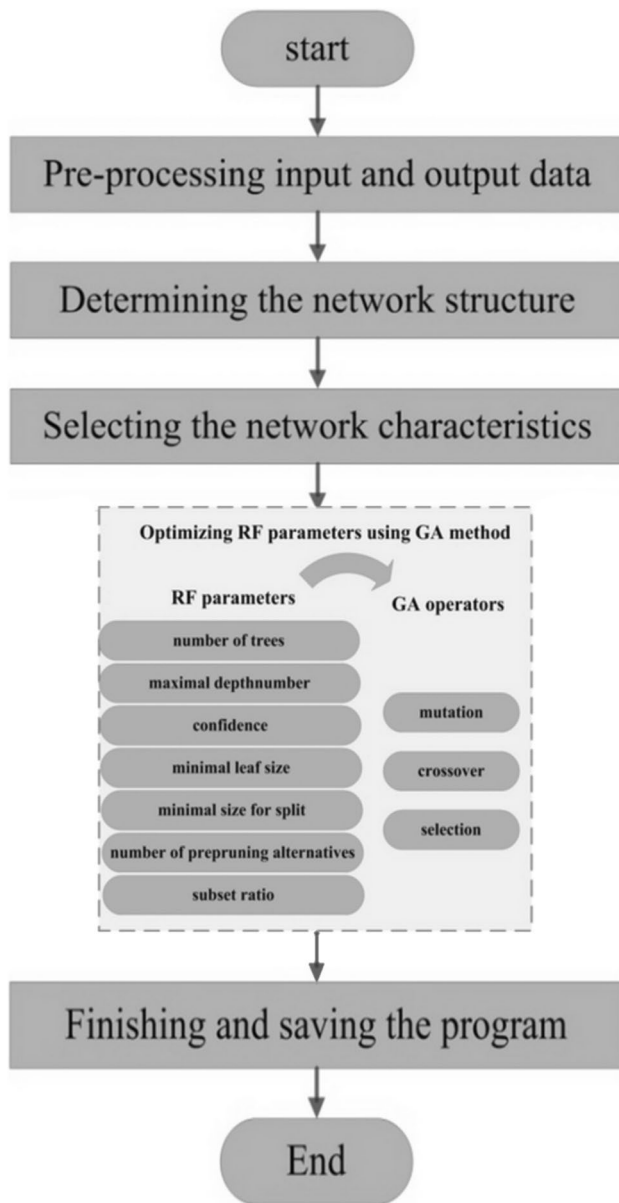The lower and upper boundaries for each parameter were determined based on a comprehensive approach

**Fig. 2** The flowchart of GA-RF model

involving preliminary testing, literature benchmarking against similar studies, analysis of the prediction task complexity, and an assessment of computational resources. For example, the mutation rate was constrained between 0.5% and 1% to maintain genetic diversity without inducing excessive randomness in the population. These values, along with those for the population size, crossover rate, and number of generations, were iteratively adjusted in exploratory runs until the GA demonstrated consistent convergence towards optimal or near-optimal solutions. The chosen parameters provided a robust search capability within the solution space, balancing exploration and exploitation aptly for the problem at hand.
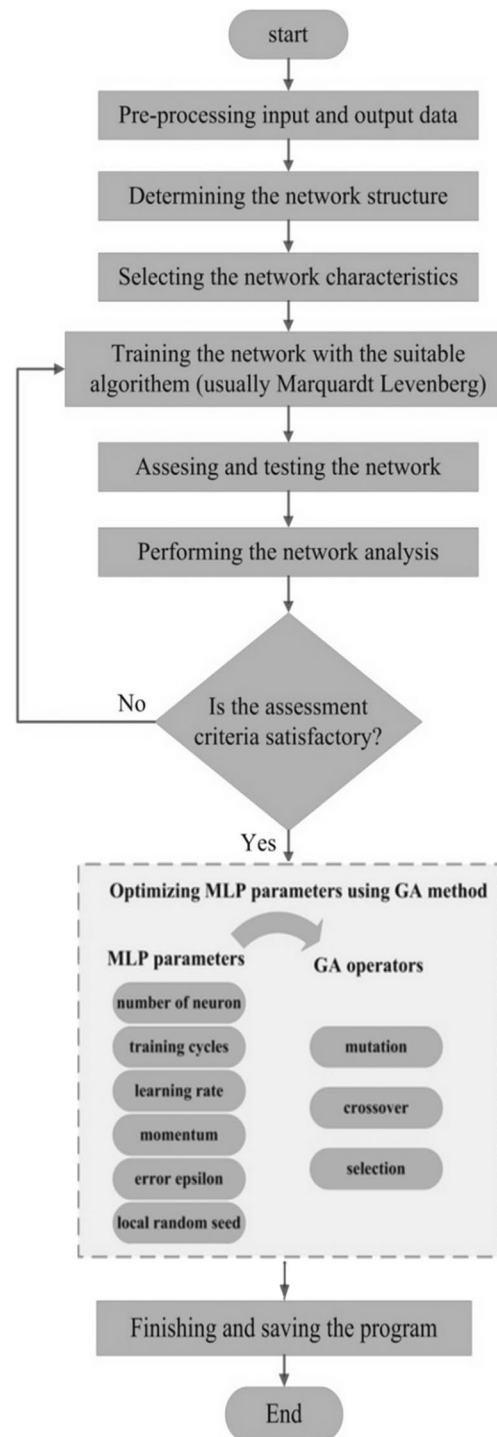


**Fig. 3** The flowchart of GA-MLP model

## Hybrid GA-RF

In the random forest method, two parameters of ntrees and mtry which are defined earlier are the most significant factors in its performance. Evident is that optimizing these parameters can increase its operating accuracy

(Rodriguez-Galiano et al. 2012). In this study, the GA is utilized as an optimizer which can minimize the score achieved by the fitness function. Hence, it leads to select the most appropriate subsets for ntrees and mtry. Furthermore, in the traditional RF method, trees have various portions in its precision. For instance, maybe some trees fortify inaccurate estimations (Adnan and Islam 2016). Based on some researches, a greedy algorithm, as an optimizer, is suggested to ameliorate the RF model but in some cases, this algorithm gives rise to becoming trapped at local optima. Likewise, a small number of high-quality individual learners are selected to create a better performance (Zhou et al. 2002). Consequently, the genetic algorithm is performed better in improving the random forest model's accuracy. The flowchart of the GA-RF model is shown in Fig. 2.

## Multi-layer Perceptron (MLP)

Multi-layer perceptron are the most common neural networks. These networks are part of the feed-forward neural networks that can be selected by appropriate number of layers and neurons; perform a nonlinear mapping with desired

**Table 2** Name of the models

| Target station | Reference stations | RF | GA-RF | MLP | GA-MLP |
|---|---|---|---|---|---|
| Ardabil | Sarein, Nir, Bileh Savar, Meshgin Shahr, Parsabad, Khalkhal | RF-1 | GA-RF-1 | MLP-1 | GA-MLP-1 |
| Sarein | Ardabil, Nir, Bileh Savar, Meshgin Shahr, Parsabad, Khalkhal | RF-2 | GA-RF-2 | MLP-2 | GA-MLP-2 |
| Nir | Ardabil, Sarein, Bileh Savar, Meshgin Shahr, Parsabad, Khalkhal | RF-3 | GA-RF-3 | MLP-3 | GA-MLP-3 |
| Bileh Savar | Ardabil, Sarein, Nir, Meshgin Shahr, Parsabad, Khalkhal | RF-4 | GA-RF-4 | MLP-4 | GA-MLP-4 |
| Meshgin Shahr | Ardabil, Sarein, Nir, Bileh Savar, Parsabad, Khalkhal | RF-5 | GA-RF-5 | MLP-5 | GA-MLP-5 |
| Parsabad | Ardabil, Sarein, Nir, Bileh Savar, Meshgin Shahr, Khalkhal | RF-6 | GA-RF-6 | MLP-6 | GA-MLP-6 |
| Khalkhal | Ardabil, Sarein, Nir, Bileh Savar, Meshgin Shahr, Parsabad | RF-7 | GA-RF-7 | MLP-7 | GA-MLP-7 |

**Table 3** Correlation coefficients of E values between different stations

| Stations | Ardabil | Sarein | Nir | Bileh Savar | Meshgin Shahr | Parsabad | Khalkhal |
|---|---|---|---|---|---|---|---|
| Ardabil | 1.000 | | | | | | |
| Sarein | 0.666 | 1.000 | | | | | |
| Nir | 0.732 | 0.738 | 1.000 | | | | |
| Bileh Savar | 0.604 | 0.607 | 0.641 | 1.000 | | | |
| Meshgin Shahr | 0.686 | 0.699 | 0.745 | 0.738 | 1.000 | | |
| Parsabad | 0.516 | 0.519 | 0.560 | 0.763 | 0.652 | 1.000 | |
| Khalkhal | 0.650 | 0.612 | 0.690 | 0.734 | 0.690 | 0.647 | 1.000 |

**Table 4** Parameters of the RF and GA-RF models

| Models | Parameter | | | | | | |
|---|---|---|---|---|---|---|---|
| | A | B | C | D | E | F | G |
| RF-1, RF-2, RF-3, RF-4, RF-5, RF-6, RF-7 | 100 | 10 | 0.100 | 2 | 4 | 3 | 0.200 |
| GA-RF -1 | 81 | 5 | 0.362 | 41 | 20 | 55 | 0.150 |
| GA-RF -2 | 81 | 5 | 0.346 | 41 | 20 | 55 | 0.142 |
| GA-RF -3 | 94 | 5 | 0.331 | 41 | 20 | 3 | 0.189 |
| GA-RF -4 | 94 | 5 | 0.366 | 41 | 20 | 3 | 0.152 |
| GA-RF -5 | 94 | 5 | 0.369 | 41 | 20 | 3 | 0.167 |
| GA-RF -6 | 81 | 5 | 0.349 | 41 | 20 | 55 | 0.158 |
| GA-RF -7 | 81 | 36 | 0.191 | 1 | 80 | 55 | 0.293 |

A: Random Forest.number_of_trees, B: Random Forest.maximal_depth, C: Random Forest.confidence, D: Random, Forest.minimal_leaf_size, E: Random Forest.minimal_size_for_split, F: Random Forest.number_of_prepruning_alternatives, G: Random Forest.subset_ratio

**Table 5** Parameters of the MLP and GA-MLP models

| Model | Parameter | | | | |
|---|---|---|---|---|---|
| | A | B | C | D | E |
| MLP-1, MLP-2, MLP-3, MLP-4, MLP-5, MLP-6, MLP-7 | 200 | 0.0100 | 0.9000 | 0.0001 | 1,992 |
| GA-MLP-1 | 77 | 0.1392 | 0.5319 | Infinity | 29 |
| GA-MLP-2 | 77 | 0.3913 | 0.5447 | Infinity | 77 |
| GA-MLP-3 | 7 | 0.1335 | 0.2405 | Infinity | 29 |
| GA-MLP-4 | 7 | 0.1423 | 0.2405 | Infinity | 29 |
| GA-MLP-5 | 7 | 0.1657 | 0.2502 | Infinity | 29 |
| GA-MLP-6 | 7 | 0.1530 | 0.2458 | Infinity | 29 |
| GA-MLP-7 | 4 | 0.1372 | 0.0714 | Infinity | 29 |

A: Neural Net.training_cycles, B: Neural Net.learning_rate, C: Neural Net.momentum, D: Neural Net.error_epsilon, E: Neural Net.local_random_seed

accuracy (Du and Swamy 2006). MLP networks have several layers: input layer, output layer and hidden layer or layers where the output of the first layer is the input vector of the second layer. In the same way, the output of the second layer is the input vector of the third layer. The output of the second layer shows the actual network response. The neurons in the upper layer are related to the neurons in the lower layer. The role of each neuron is to calculate the sum of the given inputs and then pass this sum through a function called the transfer function. The transfer function can be a linear or nonlinear function. Two common functions in multi-layer perceptron networks are sigmoid function and sigmoid tangent. The multi-layer perceptron works in such a way that a pattern is supplied to the network and its output is calculated. Comparing the actual output with the desired output causes the weight factor of the network to change so that a more accurate output is obtained next time (Chelani et al. 2002).
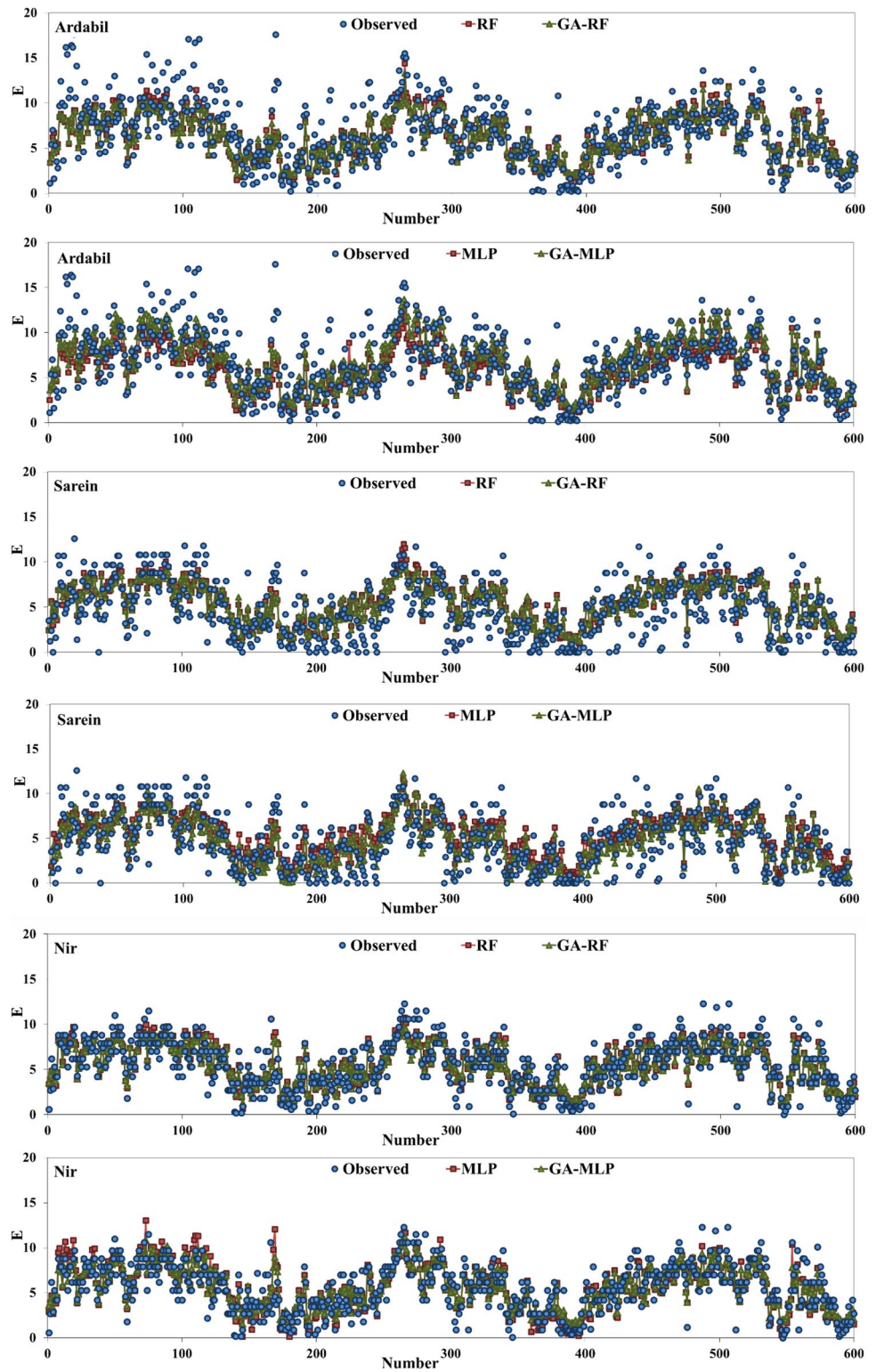
## Hybrid GA-MLP

One of the complex modeling processes in the MLP technique is specifying the number of neurons in hidden layers, local random seed, error epsilon, momentum, learning rate and training cycles. To solve this difficulty, a novel technique was developed in this research in which MLP model was integrated with Genetic Algorithm. First step of the GA-MLP model begins with selecting the population randomly. Then by considering the weight, features of each individual from the initial generation will be selected. In the second step, each individual from the population will be investigated. For this purpose, the Multi-Layer Perceptron will be implemented by defined weights, inputs and outputs for each layer and neuron. Moreover, at the end of the procedure, the difference between experiment and output models

will occur. Based on the amount of MLP errors, individuals will be selected in the third step. Then, these individuals will be rated and according to the minimum error, finest population, which contains elite persons, will be chosen. Then the finest parents will be chosen for breeding by utilizing of the GA operators. The process will replicated for the next generations and the algorithm will run in the specified number of series and the attained results are stored (Samadianfard et al. 2021). Finally, if the termination standards had satisfaction outcomes, the person with the finest function is saved. If not, this method will discover a suitable population with a different function. The algorithm utilized frequently in the training phase of this procedure is Levenberg–Marquardt, which has an accidental nature. Using GA will secure the model against this issue and elects the finest transfer function for the hidden and output layers. The flowchart of the GA-MLP model is shown in Fig. 3.

**Table 6** General results of the computations for the RF, GA-RF, MLP and GA-MLP models

| Model | Statistical parameters | | |
|---|---|---|---|
| | CC | SI | WI |
| RF-1 | 0.7670 | 0.3424 | 0.8442 |
| GA-RF-1 | 0.7743 | 0.3412 | 0.8386 |
| MLP-1 | 0.7796 | 0.3481 | 0.8407 |
| GA-MLP-1 | 0.7813 | 0.3284 | 0.8695 |
| RF-2 | 0.7797 | 0.4509 | 0.8416 |
| GA-RF-2 | 0.7858 | 0.4498 | 0.8349 |
| MLP-2 | 0.7892 | 0.4366 | 0.8484 |
| GA-MLP-2 | 0.7892 | 0.4113 | 0.8774 |
| RF-3 | 0.8428 | 0.2615 | 0.9103 |
| GA-RF-3 | 0.8486 | 0.2584 | 0.9065 |
| MLP-3 | 0.8349 | 0.2743 | 0.9114 |
| GA-MLP-3 | 0.8497 | 0.2654 | 0.9129 |
| RF-4 | 0.8502 | 0.2906 | 0.9199 |
| GA-RF-4 | 0.8511 | 0.2790 | 0.9203 |
| MLP-4 | 0.8519 | 0.3176 | 0.9077 |
| GA-MLP-4 | 0.8477 | 0.2854 | 0.9191 |
| RF-5 | 0.8617 | 0.2621 | 0.9245 |
| GA-RF-5 | 0.8673 | 0.2568 | 0.9226 |
| MLP-5 | 0.8692 | 0.2824 | 0.9158 |
| GA-MLP-5 | 0.8704 | 0.2539 | 0.9212 |
| RF-6 | 0.7091 | 0.4221 | 0.7872 |
| GA-RF-6 | 0.7162 | 0.4206 | 0.7825 |
| MLP-6 | 0.7240 | 0.4125 | 0.8031 |
| GA-MLP-6 | 0.7134 | 0.4162 | 0.7984 |
| RF-7 | 0.8123 | 0.2738 | 0.8682 |
| GA-RF-7 | 0.8105 | 0.2731 | 0.8691 |
| MLP-7 | 0.8128 | 0.2841 | 0.8624 |
| GA-MLP-7 | 0.8156 | 0.2626 | 0.8864 |

**Fig. 4** Observed and estimated E values



## Evaluation of results

For evaluating the studied model's performance, various standard statistics are utilized. In this study, Taylor diagram known as one of the evaluation meters is used. The performance indexes including Correlation coefficient (*CC*), Scattered Index (*SI*), and Willmott's Index of agreement (*WI*) are calculated with the following formulas:
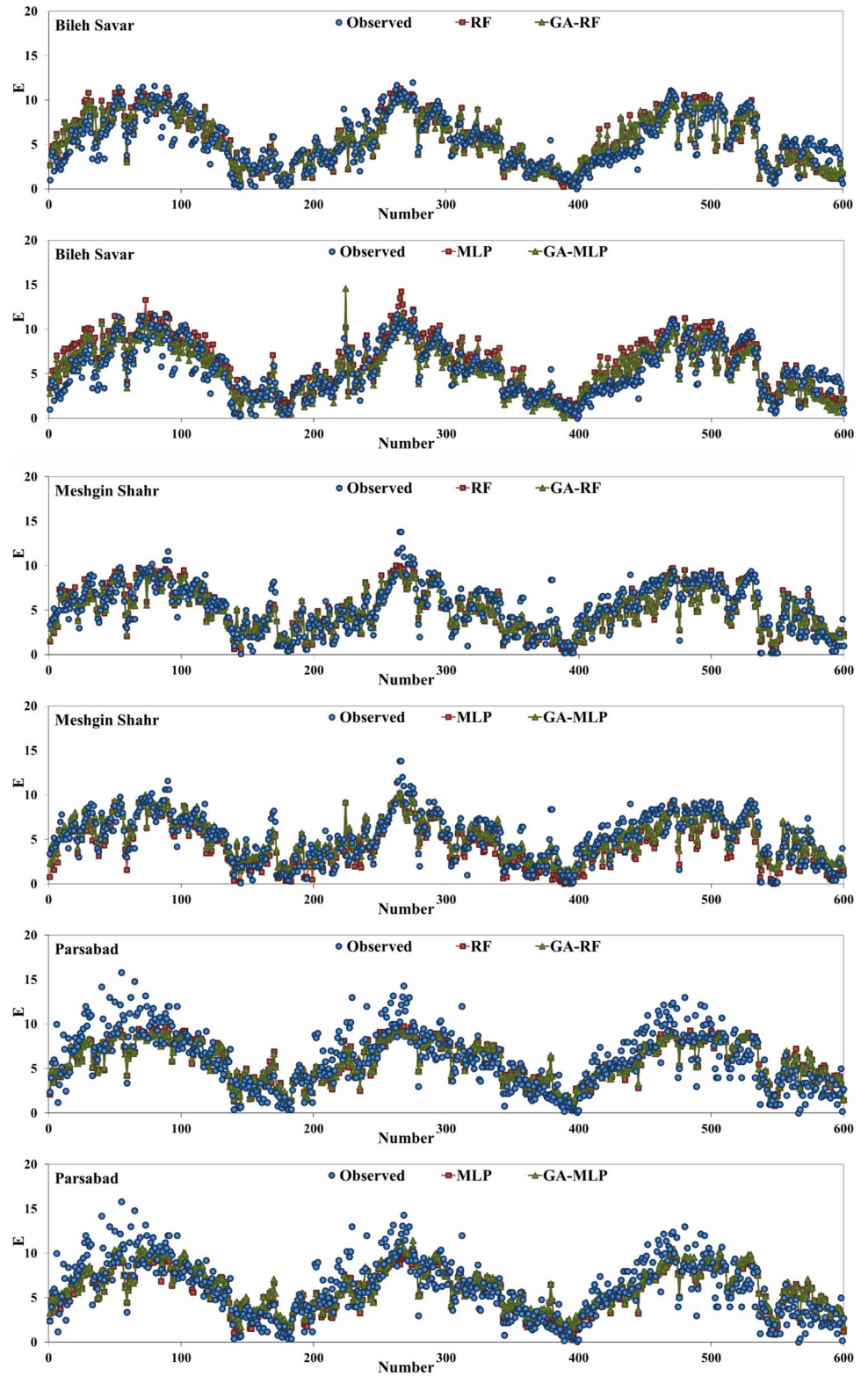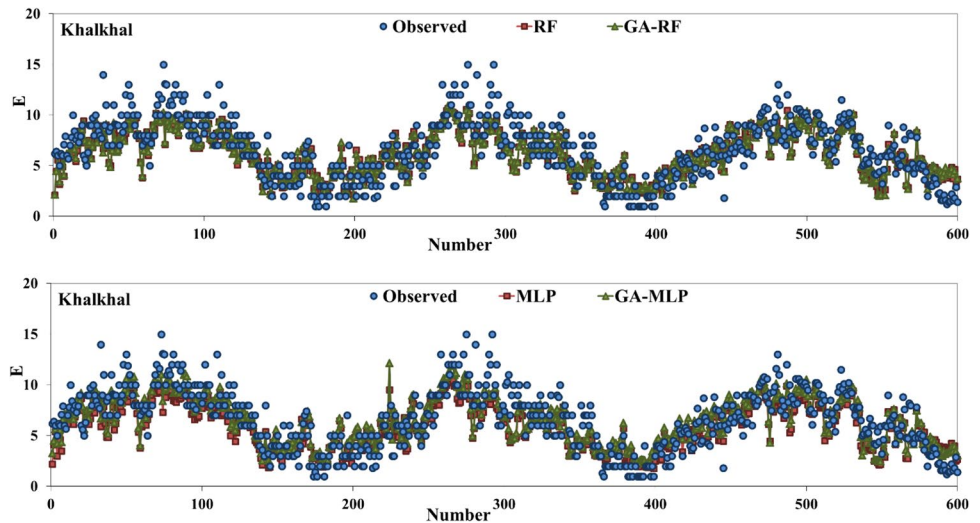
**Fig. 4** (continued)

**Fig. 4** (continued)



$$CC = \frac{\left( \sum\limits_{i=1}^{n} O_i P_i - \frac{1}{n} \sum\limits_{i=1}^{n} O_i \sum\limits_{i=1}^{n} P_i \right)}{\left( \sum\limits_{i=1}^{n} O_i^2 - \frac{1}{n} \left( \sum\limits_{i=1}^{n} O_i \right)^2 \right) \left( \sum\limits_{i=1}^{n} P_i^2 - \frac{1}{n} \left( \sum\limits_{i=1}^{n} P_i \right)^2 \right)} \quad (1)$$

$$SI = \frac{\sqrt{\frac{1}{n} \sum\limits_{i=1}^{n} \left( P_i - O_i \right)^2}}{\overline{O}} \quad (2)$$

$$WI = 1 - \left[ \frac{\sum\limits_{i=1}^{n} \left( O_i - P_i \right)^2}{\sum\limits_{i=1}^{n} \left( \left| P_i - \overline{O}_i \right| + \left| O_i - \overline{O}_i \right| \right)^2} \right] \quad (3)$$

In which n is the number of data, $P_i$ and $O_i$ are considered as the predicted and observed value for $i^{th}$ parameter of pan evaporation.

## Results and discussion

In this study, capabilities of the Random Forest (RF) and Multilayer Perceptron (MLP) models and their optimized forms with GA are investigated in estimating pan evaporation by the usage of various stations. In the current research, seven stations (Ardabil, Sarein, Nir, Bileh Savar, Meshgin Shahr, Parsabad, and Khalkhal) are considered for pan evaporation prediction. Furthermore, there is no direct method to split training and testing datasets. For example, for developing the model, Deo et al. (2018) and Samadianfard et al. (2018, 2019a, b, 2020) implemented 70% of their data, while Qasem et al. (2019) applied 67% and Zounemat-Kermani et al. (2019) used 80% of whole

data for training step. Accordingly, 70% of data is used as training, and the rest are applied for the testing part. In this research, the value of the pan evaporation in one station is considered as output, and the rest are input parameters. Using Ep of six stations, the Ep last station is estimated, and the precision of each model is examined. As shown in Table 2, the RF, MLP and the hybrid GA-RF and GA-MLP methods are used for Ep estimation in each station. Additionally, Table 3 represents the pan evaporation value correlation coefficient among all stations.

Tables 4 and 5 demonstrate the parameters of RF, GA-RF, MLP and GA-MLP models that are utilized in the development of models. In all of the RF models, Random Forest. number_of_trees (A) is 100, Random Forest.maximal_depth (B) is 10, Random Forest.confidence (C) is 0.100, Random Forest.minimal_leaf_size (D) is 2, Random Forest.minimal_ size_for_split (E) is 4, Random Forest.number_of_prepruning_alternatives (F) is 3, and Random Forest.subset_ratio (G) of 0.200 is utilized. Additionally, in all of the MLP models, Neural Net.training_cycles (A) is 200, Neural Net.learning_rate (B) is 0.0100, Neural Net.momentum (C) is 0.9000, Neural Net.error_epsilon (D) is 0.0001 and Neural Net.local_ random_seed (E) is 1,992. Furthermore, in the hybrid models, the genetic algorithm changed and improved all of the above-mentioned values to increase the accuracy of estimation.

Table 6 shows the outcomes of the RF, GA-RF, MLP and GA-MLP models in Ep estimation for each station. According to Table 6, among all of the models which are calculated by the RF method, the RF-5 with CC of 0.8617, SI of 0.2621, and WI of 0.9245 has high accuracy in Ep prediction. Moreover, among all of the models which are calculated by the MLP method, the MLP-3 with CC of 0.8349, SI of 0.2743, and WI of 0.9114 has suitable accuracy. In addition, when a genetic algorithm utilized as an optimizer, the performance of all RF and MLP models, except MLP-6, increases. In hybrid models the Willmott's Index and Correlation
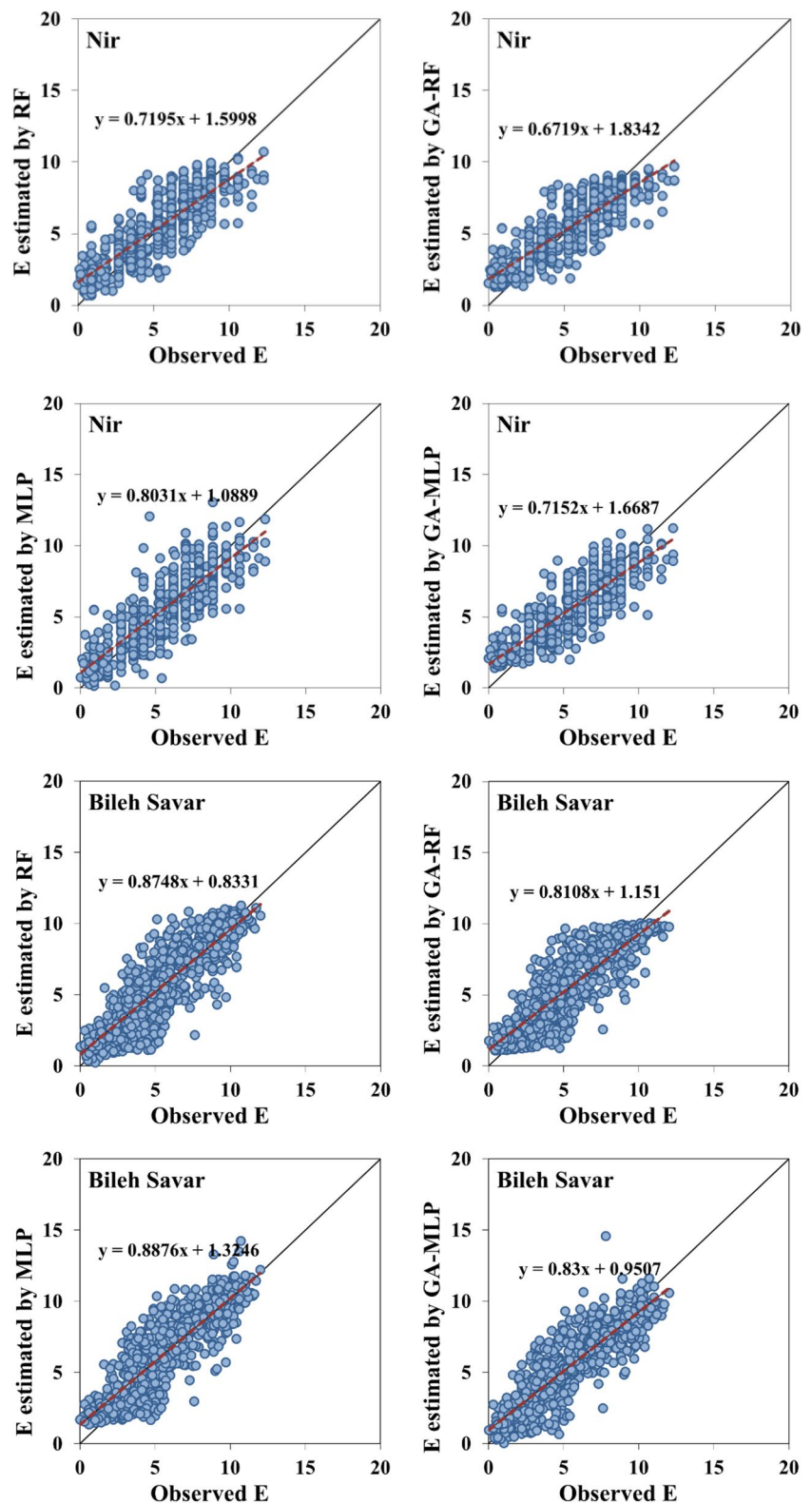
**Fig. 5** The scatter plots of observed and estimated E values



coefficient parameters are higher than standalone models. However, the GA-MLP-5 model with *CC* of 0.8704, *SI* of 0.2539, and *WI* of 0.9212 creates better results when compared with other ones. It is obvious that the genetic algorithm

as an optimizer improve the performance of models in Ep prediction. This method decreases the *SI* parameter by 10% in the best model (GA-MLP-5). Among the studied stations, in hybrid models, Bileh Savar and Meshgin Shahr had

**Fig. 5** (continued)



the best performance according to their error meters. Both GA-RF-4, GA-MLP-4 and GA-RF-5, GA-MLP-5 models provide almost same performances however, as GA-MLP-5 can better capture the higher Ep value, it slightly outperforms GA-RF-4, GA-MLP-4 and GA-RF-5 models. Additionally, closer examination of Table 5 represents that the effect of the GA optimizer at RF model for pan evaporation forecasting at Khalkhal station was minimal.
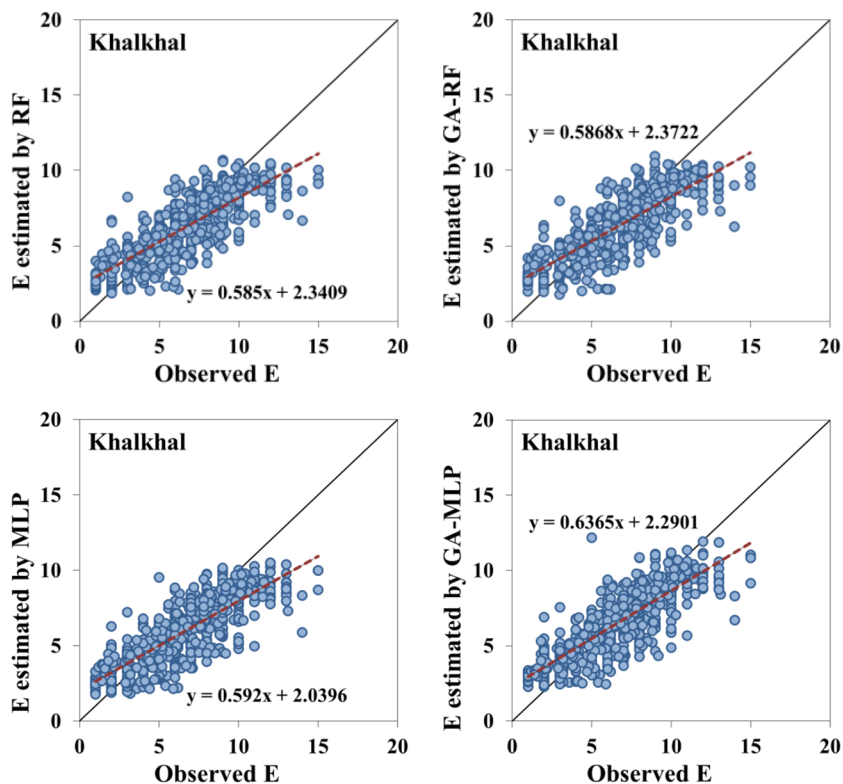
**Fig. 5** (continued)

**Fig. 5** (continued)



The performance of the standalone and hybrid models is demonstrated in Fig. 4 for seven stations. From Fig. 4, it can be inferred that the genetic algorithm ameliorated the precision of the Ep estimation. Furthermore, scatter plots of observed and estimated values of pan evaporation for each station are shown in Fig. 5 using RF, GA-RF, MLP and GA-MLP methods. Evident is that the RF-3 and MLP-5 models has high accuracy among other default models and the GA-RF-5 and GA-MLP-5 have the best performance in comparison with all models in prediction of Ep value. In Fig. 5 blue solid points is predicted E by RF, GA-RF, MLP and GA-MLP methods and estimated E, black solid line is bisector line (y = x), and red dotted line is trend line, in this figure, the higher the scatter of blue points around the black line, the higher the accuracy of the model. Based on Fig. 5, in both stations (Bileh Savar and Meshgin Shahr), hybrid GA-RF and GA-MLP models are in great agreement with the bisector line. However, almost all of the other models are below the 1:1 line and underestimate the Ep values. By and large, the GA-MLP-5 model slightly demonstrates precise results than GA-RF-5.

Using Taylor diagrams, the correlation and standard deviation values between estimated and observed pan evaporation are investigated. In Fig. 6, Taylor diagrams are presented for all RF, GA-RF, MLP and GA-MLP models. *RMSE* parameter in the diagram is defined as the distance from the reference point (green dot) to any other point. Therefore, the most accurate model is the minimum space between green and the correspondent dot (Taylor,
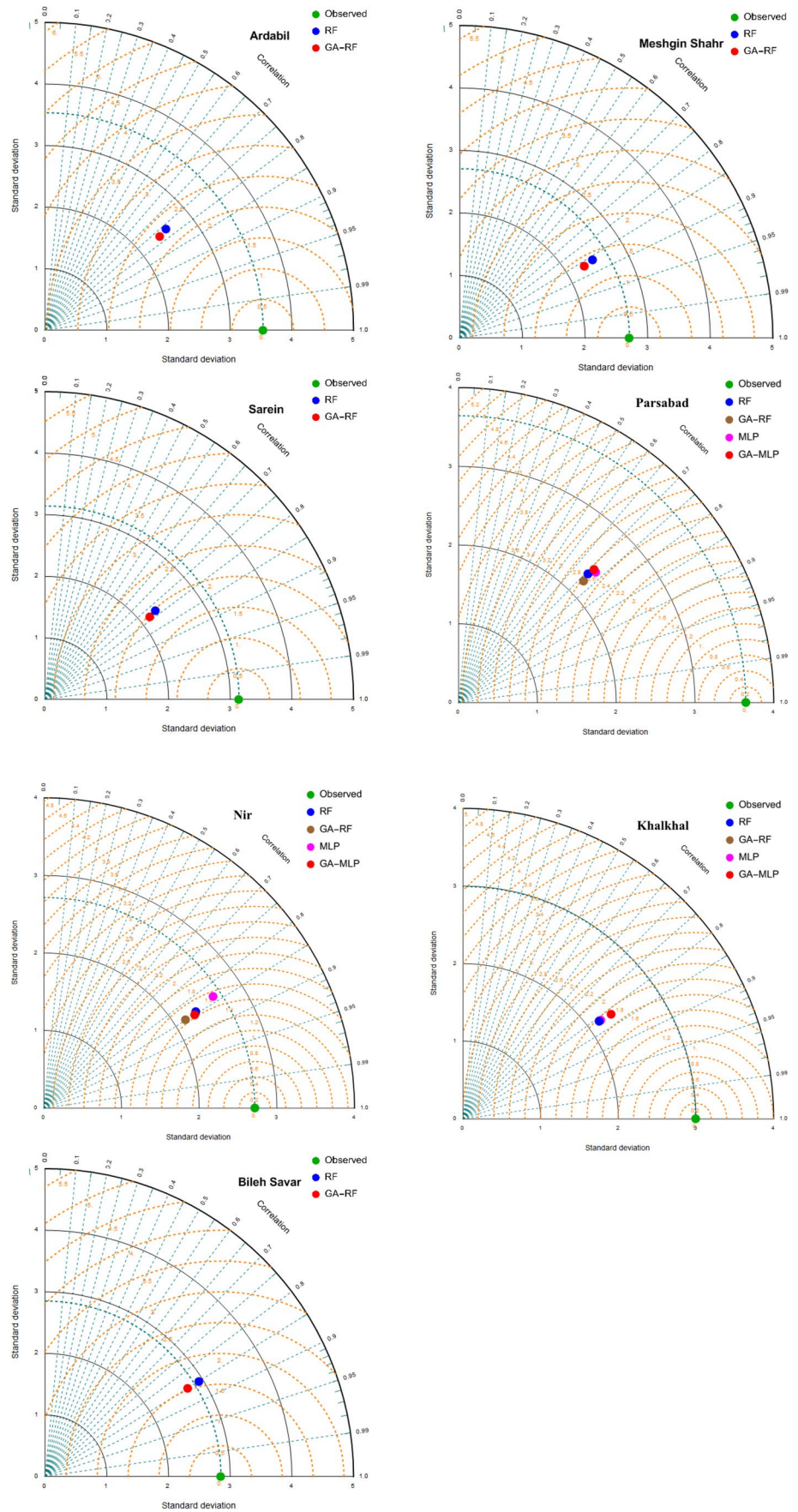
2001). According to Fig. 6, the red point (GA-MLP-5) in the Meshgin Shahr station is the closest to the green point and has a minimum distance from the reference point. Therefore, it yields accurate results with minimum error in Ep estimation. Conclusively, it can be stated that in the case of lack of measured Ep data, especially in developing countries, it may be estimated using the corresponding values in the neighboring stations with acceptable accuracy.

Overall, GA-MLP-5 has superior results than all other models. Moreover, it should be noted that this study has some limitations. For instance, the gathered datasets in seven stations are for Ardabil province in Iran. Therefore, the studied places have almost the same climates. Hence, it would be better to use proposed techniques for various locations with diverse climates and examine the accuracy of new models. Additionally, using the same datasets, novel methods could be implemented for the extension of the study.

## Sensitivity analysis

To better understand the impact of adjacent station on prediction of Ep of target stations, SI evaluation parameter was used for various stations. To achieve this purpose, the GA-RF were utilized for sensitivity analysis to discover the most effective adjacent station on target stations The table appears to be showing the results of a sensitivity analysis for a GA-RF model. The Table 7 lists several target stations in the first

**Fig. 6** The Taylor diagrams of observed and estimated E values

column and corresponding nearby stations at the top of each column. The "All station (SI)" row shows a baseline Scatter Index (SI) value for the GA-RF model when all nearby stations are included. Below this, there are rows that represent the new SI value after eliminating data from one of the nearby stations. Here's a more detailed analysis based on the provided table:

Ardabil as Target Station: The baseline SI when all stations are included is 0.3412. After eliminating data from each of the nearby stations one at a time, we see slight increases in the SI value, indicating a minor decrease in model accuracy. The largest increase in SI occurs when Sarein is eliminated, suggesting that Sarein's data is most important for predicting pan evaporation at Ardabil.

Sarein as Target Station: The baseline SI is 0.4498. The most significant impact is observed when the data from Ardabil is eliminated, which results in an SI of 0.4556, indicating that Ardabil has the most influence on the model's accuracy for Sarein.

Nir as Target Station: The baseline SI is 0.2584. The SI increases to 0.2736 when Sarein is eliminated, which is the highest change among the stations listed. This suggests that Sarein's data is quite influential on the predictions for Nir.

Bileh Savar as Target Station: The baseline SI is 0.2790. The highest SI observed after elimination is 0.2900 when Sarein is removed, hinting that Sarein's data contributes significantly to the accuracy for Bileh Savar.

Meshgin Shahr as Target Station: The baseline SI is 0.2568. When Bileh Savar is eliminated, the SI increases the most to 0.2635, indicating that Bileh Savar's data has a substantial impact on Meshgin Shahr's predictions.

Parsabad as Target Station: The baseline SI is 0.4206. Eliminating Meshgin Shahr's data increases the SI to 0.4228, which is the largest increase, suggesting Meshgin Shahr's data is important for predictions at Parsabad.

Khalkhal as Target Station: The baseline SI is 0.2731. The highest increase in SI is to 0.2814 when Parsabad is eliminated, which means Parsabad's data has a notable influence on the model's accuracy for Khalkhal.

In summary, the table shows that for each target station, the exclusion of data from certain nearby stations causes a decrease in model accuracy, as indicated by the increase in SI values. This highlights the contribution of specific stations to the predictive power of the GA-RF model and helps identify which nearby stations' data are most valuable for accurate pan evaporation predictions at each target station. Also, the Meshgin Shahr, Nir, Sarein, Parsabad, Bileh Savar, Bileh Savar and Parsabad are the most effective adjacent station to Ardabil, Sarein, Nir, Bileh Savar, Meshgin Shahr, Parsabad and Khalkhal stations, respectively.

## Study limitations and future outlook

The study focuses on the application of Random Forest (RF) and Multilayer Perceptron (MLP) models integrated with the genetic algorithm (GA) for predicting pan evaporation at target stations using neighboring reference station data. The research demonstrates that the GA-MLP-5 model outperforms other models and highlights the importance of accurate estimation of evaporation in hydrological studies and water resource management. However, the study has limitations, such as the use of datasets from only seven stations in Ardabil province, Iran, and the need for further research in diverse climate conditions. In the future, it would be beneficial to expand the study to different locations with diverse climates to examine the accuracy of the proposed techniques. Additionally, the study could benefit from exploring novel methods for extending the research using the same datasets.

**Table 7** The effect of removing an adjacent station in predicting the target station Ep

| Target station | Ardabil | Sarein | Nir | Bileh Savar | Meshgin Shahr | Parsabad | Khalkhal |
|---|---|---|---|---|---|---|---|
| All station (SI) | 0.3412 | 0.4498 | 0.2584 | 0.2790 | 0.2568 | 0.4206 | 0.2731 |
| Eliminate | Sarein | Ardabil | Ardabil | Ardabil | Ardabil | Ardabil | Ardabil |
| SI | 0.3451 | 0.4556 | 0.2628 | 0.2811 | 0.2618 | 0.4210 | 0.2793 |
| Eliminate | Nir | Nir | Sarein | Sarein | Sarein | Sarein | Sarein |
| SI | 0.3453 | 0.4792 | 0.2736 | 0.2900 | 0.2571 | 0.4233 | 0.2745 |
| Eliminate | Bileh Savar | Bileh Savar | Bileh Savar | Nir | Nir | Nir | Nir |
| SI | 0.3430 | 0.4503 | 0.2589 | 0.2817 | 0.2614 | 0.4239 | 0.2805 |
| Eliminate | Meshgin Shahr | Meshgin Shahr | Meshgin Shahr | Meshgin Shahr | Bileh Savar | Bileh Savar | Bileh Savar |
| SI | 0.3472 | 0.4524 | 0.2652 | 0.2911 | 0.2635 | 0.4525 | 0.2782 |
| Eliminate | Parsabad | Parsabad | Parsabad | Parsabad | Parsabad | Meshgin Shahr | Meshgin Shahr |
| SI | 0.3429 | 0.4500 | 0.2609 | 0.3071 | 0.2586 | 0.4228 | 0.2778 |
| Eliminate | Khalkhal | Khalkhal | Khalkhal | Khalkhal | Khalkhal | Khalkhal | Parsabad |
| SI | 0.3459 | 0.4501 | 0.2650 | 0.2878 | 0.2572 | 0.4271 | 0.2814 |

# Conclusion

In this study, for estimating the pan evaporation in the target station by the usage of six available stations as input, the tree-based method (random Forest), multi-layer perceptron and their hybrid form which is optimized by the genetic algorithm are utilized. For comparison of the models, three statistical indicators including Correlation coefficient (CC), Scattered Index (SI), and Willmott's Index of agreement (WI) are used. Based on the obtained results, the RF-5 and MLP-3 exhibited better performance among other RF and MLP models in Ep estimation. As well as, hybrid models outperformed all default RF and MLP models, which indicates the high efficiency of GA to improve the results of machine learning methods. Furthermore, the GA-MLP-5 model demonstrated the highest prediction precision compared with all other models with the CC of 0.8704, SI of 0.2539, and WI of 0.9212. Sensitivity analysis was performed using the GA-RF method to identify the most important neighbor station on the target station, which is an accurate method of introducing the most sensitive parameters or stations in modeling. The results carried out by using hybrid models created better estimation with high potential in pan evaporation at the specific station, which can be utilized from the hybrid models used to predict other meteorological parameters at other stations and water resource management issues. It is also possible to increase the quality of modeling by integrating GA with other machine learning methods. Conclusively, in the absence of Ep-dependent data, the Ep data of adjacent stations can be used. Conclusively, the results carried out using hybrid models created better estimation with high potential in pan evaporation at the specific station.

# Declarations

# Reference

Abghari H, Ahmadi H, Besharat S, Rezaverdinejad V (2012) Prediction of daily pan evaporation using wavelet neural networks. Water Resour Manage 26(12):3639–3652

Adnan MN, Islam MZ (2016) Optimizing the number of trees in a decision forest to discover a subforest with high ensemble accuracy using a genetic algorithm. Knowl - Based Syst 110:86–97

Adnan RM, Malik A, Kumar A, Parmar KS, Kisi O (2019) Pan evaporation modeling by three different neuro-fuzzy intelligent systems using climatic inputs. Arab J Geosci 12:606

Alipour A, Yarahmadi J, Mahdavi M (2014) Comparative study of M5 model tree and artificial neural network in estimating reference evapotranspiration using MODIS products. J Climatol 2014:1-11

Allawi MF, El-Shafie A (2016) Utilizing RBF-NN and ANFIS methods for multi-lead ahead prediction model of evaporation from reservoir. Water Resour Manag 30(13):4773–4788

Behrooz K, Salim H, Abderrazek S, Shun-Peng Z, Nguyen-Thoi T (2019) SVR-RSM: a hybrid heuristic method for modeling monthly pan evaporation. Environ Sci Pollut Res 26:35807–35826

Birbal P, Azamathulla H, Leon L, Kumar V, Hosein J (2021) Predictive modelling of the stage–discharge relationship using gene-expression programming. Water Supply 21(7):3503–3514

Breiman L (2001) Random forests. Mach Learn 45(1):5–32

Chaplot B (2021) Prediction of rainfall time series using soft computing techniques. Environ Monit Assess 193(11):721

Chaudhary A, Kolhe S, Kamal R (2016a) An improved random forest classifier for multi-class classification. Inf Process Agric 3(4):215–222

Chaudhary A, Kolhe S, Kamal R (2016b) A hybrid ensemble for classification in multiclass datasets: an application to oilseed disease dataset. Comput Electron Agric 124:65–72

Chaudhary A, Kolhe S, Kamal R (2020) A particle swarm optimization-based ensemble for vegetable crop disease recognition. Comput Electron Agric 178:1–7

Chelani A, Chalapati Rao CV, Phadke KM, Hasan MZ (2002) Prediction of sulphur dioxide concentration using artificial neural networks. Environ Model Soft 17:161–168

Chen YY, Cheng Y, Cheng Q, Yu H, Li D (2017) Short-term prediction model for ammonia nitrogen in aquaculture pond water based on optimized LSSVM. Int Agric Eng J 26(3):416–427

Deo RC, Ghorbani MA, Samadianfard S, Maraseni T, Bilgili M, Biazar M (2018) Multi-layer perceptron hybrid model integrated with the firefly optimizer algorithm for windspeed prediction of target site using a limited set of neighboring reference station data. Renew Energy 116:309–323

Du KL, Swamy MN (2006) Neural networks in a soft computing framework. Springer Science and Business Media L, London, pp 566

Fan J, Wu L, Zhang F, Xiang Y, Zheng J (2016) Climate change effects on reference crop evapotranspiration across different climatic zones of China during 1956–2015. J Hydrol 542:923–937

Feng Y, Jia Y, Zhang Q, Gong D, Cui N (2018) National-scale assessment of pan evaporation models across different climatic zones of China. J Hydrol 564:314–328

Ghaemi A, Rezaie-Balf M, Adamowski J, Kisi O, Quilty J (2019) On the applicability of maximum overlap discrete wavelet transform integrated with MARS and M5 model tree for monthly pan evaporation prediction. Agric For Meteorol 278:107647

Goldberg DE, Holland JH (1989) Genetic algorithms and machine learning. Mach Lear 3:95–99

Gundalia MJ, Dholakia MB (2013) Estimation of pan evaporation using mean air temperature and radiation for monsoon season in Junagadh region. Int J Eng Res Appl 3(6):64–70

Haddadi F, Moazenzadeh R, Mohammadi B (2022) Estimation of actual evapotranspiration: A novel hybrid method based on remote sensing and artificial intelligence. J Hydro 609:127774

Hassan MA, Khalil A, Kaseb S, Kassem MA (2017) Exploring the potential of tree-based ensemble methods in solar radiation modeling. Appl Energ 203:897–916

Holland JH (1992) Adaptation in natural and artificial systems: an introductory analysis with applications to biology, control, and artificial intelligence. Ann Arbor 6:126–137

Jayathilake T, Gunathilake MB, Wimalasiri EM, Rathnayake U (2023) Wetland water level prediction in the context of machine learning techniques: Where do we stand? Environments 10(5–75):1–17

Keshtegar B, Piri J, Kisi O (2016) A nonlinear mathematical modeling of daily pan evaporation based on conjugate gradient method. Comput Electron Agric 127:120–130

Kim S, Singh VP, Seo Y (2014) Evaluation of pan evaporation modeling with two different neural networks and weather station data. Theor Appl Climatol 117(1–2):1–13

Kisi O (2009) Daily pan evaporation modelling using multi-layer perceptrons and radial basis neural networks. Hydrol Process 23(2):213–223

Kisi O (2015) Pan evaporation modeling using least square support vector machine, multivariate adaptive regression splines and M5 model tree. J Hydrol 528:312–320

Kisi O, Heddam S (2019) Evaporation modelling by heuristic regression approaches using only temperature data. Hydrol Sci J 64(6):653–672

Kisi O, Genc O, Dinc S, Zounemat-Kermani M (2016) Daily pan evaporation modeling using chi-squared automatic interaction detector, neural networks, classification and regression tree. Comput Electron Agric 122:112–117

Ließ M, Glaser B, Huwe B (2012) Uncertainty in the spatial prediction of soil texture: comparison of regression tree and Random Forest models. Geoder 170:70–79

Lin GF, Lin HY, Wu MC (2013) Development of a support-vector-machine-based model for daily pan evaporation estimation. Hydrol Process 27(22):3115–3127

Lu X, Ju Y, Wu L, Fan J, Zhang F, Li Z (2018) Daily pan evaporation modeling from local and cross-station data using three tree-based machine learning models. J Hydrol 566:668–684

Majhi B, Naidu D, Mishra AP, Satapathy SC (2019) Improved prediction of daily pan evaporation using Deep-LSTM model. Neural Comput Appl 32:7823–7838

Malik A, Kumar A, Kisi O (2017) Monthly pan-evaporation estimation in Indian central Himalayas using different heuristic approaches and climate-based models. Comput Electron Agric 143:302–313

Mohammadi B (2023) Modeling various drought time scales via a merged artificial neural network with a firefly algorithm. Hydrology 10(3):58

Nawar S, Mouazen AM (2017) Comparison between random forests, artificial neural networks and gradient boosted machines methods of on-line Vis-NIR spectroscopy measurements of soil total nitrogen and total carbon. Sens 17(10):2428

Nourani V, Molajou A, Tajbakhsh AD, Najafi H (2019) A wavelet-based data mining technique for suspended sediment load modeling. Water Resour Manag 33(5):1769–1784

Piri J, Amin S, Moghaddamnia A, Han D, Remesun D (2009) Daily pan evaporation modelling in a hot and dry climate. J Hydrol Eng 14:803–811

Prasad AM, Iverson LR, Liaw A (2006) Newer classification and regression tree techniques: Bagging and Random Forests for ecological prediction. Ecosystems 9(2):181–199

Qasem SN, Samadianfard S, Kheshtgar S, Jarhan S, Kisi O, Shamshirband S, Chau KW (2019) Modeling monthly pan evaporation using wavelet support vector regression and wavelet artificial neural networks in arid and humid climates. Eng Appl Computat Fluid Mech 13(1):177–187

Rahimikhoob A (2009) Estimating daily pan evaporation using artificial neural network in a semi-arid environment. Theor Appl Climatol 98(1):101–105

Ravansalar M, Rajaee T, Kisi O (2017) Wavelet-linear genetic programming: a new approach for modeling monthly streamflow. J Hydrol 549:461–475

Rodriguez-Galiano V, Ghimire B, Rogan J, Chica-Olmo M, Rigol-Sanchez J (2012) An assessment of the effectiveness of a random forest classifier for land-cover classification. ISPRS J Photogramm Remote Sens 67:93–104

Samadianfard S, Jarahan S, Sadri HN (2018) Application of support vector regression integrated with firefly optimization algorithm for predicting global solar radiation. J Energy Syst 2(4):180–189

Samadianfard S, Majnooni-Heris A, Qasem SN, Kisi O, Shamshirband S, Chau KW (2019a) Daily global solar radiation modeling using data-driven techniques and empirical equations in a semi-arid climate. Eng Appl Comput Fluid Mech 13(1):142–157

Samadianfard S, Jarhan S, Salwana E, Mosavi A, Shamshirband S, Akib S (2019b) Support Vector regression integrated with fruit fly optimization algorithm for river flow forecasting in Lake Urmia basin. Water 11:19–34

Samadianfard S, Hashemi S, Kargar K, Izadyar M, Mostafaeipour A, Mosavi A, Shamshirband S (2020) Wind speed prediction using a hybrid model of the multi-layer perceptron and whale optimization algorithm. Ener Rep 6:1147–1159

Samadianfard S, Kargar K, Shadkani S, Abbaspour A, SadeghSafar MJ (2021) Hybrid models for suspended sediment prediction: optimized random forest and multi-layer perceptron through genetic algorithm and stochastic gradient descent methods. Neural Comput Appli 34:3033–3051

Schwefel HP (1993). Evolution and Optimum Seeking: The Sixth Generation. John Wiley & Sons Ltd., Hoboken

Sebbar A, Heddam S, Djemili L (2019) Predicting daily pan evaporation (Epan) from dam reservoirs in the mediterranean regions of Algeria: OPELM vs OSELM. Environ Processes 6(1):309–319

Taylor KE (2001) Summarizing multiple aspects of model performance in a single diagram. J Geophys Res Atmos 106(D7):7183–7192

Wang L, Niu Z, Kisi O, Li Ca, Yu D (2017) Pan evaporation modeling using four different heuristic approaches. Comput Electron Agr 140:203–213

Wu L, Huang G, Fan J, Ma X, Zhou H, Zeng W (2020) Hybrid extreme learning machine with meta-heuristic algorithms for monthly pan evaporation prediction. Comput Electron Agr 168:105115

Xu L, Liu S, Li D (2017) Prediction of water temperature in prawn cultures based on a mechanism model optimized by an improved artificial bee colony. Comput Electron Agric 140:397–408

Yang H, Hassan SG, Wang L, Li D (2017) Fault diagnosis method for water quality monitoring and control equipment in aquaculture

based on multiple SVM combined with DS evidence theory. Comput Electron Agric 141:96–108

Yaseen ZM, Al-Juboori AM, Beyaztas U, Al-Ansari N, Chau KW, Qi C, Shahid S (2020) Prediction of evaporation in arid and semi-arid regions: a comparative study using different machine learning models. Eng Appl Comput Fluid Mech 14(1):70–89

Yu H, Chen Y, Hassan SG, Li D (2016) Prediction of the temperature in a Chinese solar greenhouse based on LSSVM optimized by improved PSO. Comput Electron Agric 122:94–102

Zhou Z-H, Wu J, Tang W (2002) Ensembling neural networks: many could be better than all. Artif Intell 137:239–263

Zhu S, Bonacci O, Oskoruš D, Hadzima-Nyarko H, Wu S (2009) Long term variations of river temperature and the influence of air temperature and river discharge: case study of Kupa River watershed in Croatia. J Hydrol Hydromech 67(4):305–313

Zhu S, Hadzima-Nyarko M, Gao A, Wang F, Wu J, Wu S (2019a) Two hybrid data-driven models for modeling water-air temperature relationship in rivers. Environ Sci Pollut Res 26:12622–12630

Zhu S, Nyarko EK, Hadzima-Nyarko M, Heddam S, Wu S (2019b) Assessing the performance of a suite of machine learning models for daily river water temperature prediction. PeerJ 7:e7065

Zounemat-Kermani M, Seo Y, Kim S, Ghorbani MA, Samadianfard S, Naghshara S, Kim NW, Singh VP (2019) Can decomposition approaches always enhance soft computing models? Predicting the dissolved oxygen concentration in the St. Johns River, Florida. Appl Sci 9:25–34