

Original papers

Estimating the soil water retention curve: Comparison of multiple nonlinear regression approach and random forest data mining technique



M. Rastgou^a, H. Bayat^{b,*}, M. Mansoorizadeh^c, Andrew S. Gregory^d

^a Department of Soil Science, Faculty of Agriculture, Bu-Ali Sina University, Hamedan, Iran

^b Department of Soil Science, Faculty of Agriculture, Bu-Ali Sina University, Hamedan, Iran

^c Department of Computer Science, Faculty of Engineering, Bu-Ali Sina University, Hamedan, Iran

^d Sustainable Agriculture Sciences Department, Rothamsted Research, Harpenden, Hertfordshire AL5 2JQ, UK

ARTICLE INFO

Keywords:

Pedotransfer functions
Soil water retention curve
Soil texture
Soil structure
van Genuchten

ABSTRACT

This study evaluates the performance of the random forest (RF) method on the prediction of the soil water retention curve (SWRC) and compares its performance with those of nonlinear regression (NLR) and Rosetta-based pedotransfer functions (PTFs), which has not been reported so far. Fifteen RF and NLR-based PTFs were constructed using readily-available soil properties for 223 soil samples from Iran. The general performance of RF and NLR-based PTFs was quantified by the integral root mean square error (IRMSE), Akaike's information criterion (AIC) and coefficient of determination (R^2). The results showed that the accuracy of the RF-based PTFs was significantly ($P < 0.05$) better than the NLR-based PTFs, and that the reliability of the NLR-based PTFs was significantly ($P < 0.01$) better than the RF-based PTFs and all of the Rosetta-based PTFs. The average values of the IRMSE, AIC and R^2 of the RF method were $0.041 \text{ cm}^3 \text{ cm}^{-3}$, -16997.7 , and 0.987 , and $0.053 \text{ cm}^3 \text{ cm}^{-3}$, -15547.5 , and 0.981 for the training and testing steps of all PTFs, respectively, whereas the values for the NLR method were $0.046 \text{ cm}^3 \text{ cm}^{-3}$, -16616.4 , and 0.984 , and $0.048 \text{ cm}^3 \text{ cm}^{-3}$, -16355.6 , and 0.983 for the training and testing steps, respectively. The PTF5 of the RF and NLR methods, with inputs of sand and clay contents, bulk density, and the water content at field capacity and permanent wilting point, had the greatest R^2 values (0.987 and 0.989 , respectively), and the lowest IRMSE values (0.039 and $0.032 \text{ cm}^3 \text{ cm}^{-3}$, respectively) compared to other PTFs for the testing step. Overall, the RF method had less reliability for the prediction of the SWRC compared to the NLR method due to overprediction, uncertainty of determination of forest scale and instability in the testing step. These findings could provide the scientific basis for further research on the RF method.

1. Introduction

Soil hydraulic properties are principle factors that control the movement of water and solutes in the soil. Determination of the soil hydraulic properties is required for many distinct applications linked with irrigation, land use planning, drainage and drought risk assessment (Dobarco et al., 2019). The soil water retention curve (SWRC) is one of the most important soil hydraulic properties. It defines the relationship between soil matric potential and soil water content (Hillel, 1998). The SWRC is a crucial parameter in soil and water management for sustainable and improved agricultural production (Shwetha and Varija, 2015). The SWRC depends principally on texture, structure and bulk density (BD) of soils (Wassar et al., 2016). Many methods have been introduced for the direct measurement of the SWRC in the

laboratory (e.g., the hanging water column and pressure plate methods) (Klute, 1986) and in the field (e.g., tensiometric) (Bruce and Luxmoore, 1986). Measurements of the SWRC at several matric potentials can be expensive, difficult and time-consuming, hence it is common to predict it by modelling (Dobarco et al., 2019). Modelling of soil water is an essential tool in evaluating the effects of different managements on crop yield and environmental quality (Verhagen, 1997).

Pedotransfer functions (PTFs) translate easy-to-measure data that we have (e.g., texture class, particle size distribution (PSD) and BD) into difficult-to-measure data that we need (soil hydraulic data) (Bouma, 1989). Estimates of the SWRC by PTFs are valuable in many studies, such as hydrology, soil mapping and hydrogeology (Børgesen and Schaap, 2005). The point- and parametric-based PTFs are generally developed to predict water content at specific matric potential values

* Corresponding author.

E-mail addresses: mostafa.rastgo@gmail.com (M. Rastgou), h.bayat@basu.ac.ir (H. Bayat), mansoorm@basu.ac.ir (M. Mansoorizadeh), andy.gregory@rothamsted.ac.uk (A.S. Gregory).

<https://doi.org/10.1016/j.compag.2020.105502>

Received 15 November 2019; Received in revised form 8 April 2020; Accepted 11 May 2020

0168-1699/ © 2020 Elsevier B.V. All rights reserved.

Nomenclature

| | |
|---------------|--|
| S | Sand content (%) |
| C | Clay content (%) |
| d_g | Geometric mean diameter (mm) |
| δ_g | Geometric standard deviation (-) |
| BD | Bulk density (g cm^{-3}) |
| TP | Total porosity ($\text{cm}^3 \text{cm}^{-3}$) |
| θ_{FC} | Water content at field capacity, 33 kPa ($\text{cm}^3 \text{cm}^{-3}$) |

| | |
|----------------|--|
| θ_{PWP} | Water content at 1500 kPa ($\text{cm}^3 \text{cm}^{-3}$) |
| OM | Organic matter content (%) |
| K_s | Saturated hydraulic conductivity (cm day^{-1}) |
| θ_s | Saturated water content ($\text{cm}^3 \text{cm}^{-3}$) |
| θ_r | Residual water content ($\text{cm}^3 \text{cm}^{-3}$) |
| RF | Random forest |
| NLR | Nonlinear regression |
| SWRC | Soil water retention curve |

and the entire SWRC, respectively, by multiple linear (MLR) and non-linear regression (NLR) methods (Gunarathna et al., 2019b; Merdun et al., 2006; Minasny et al., 1999; Rajkai et al., 2004; Tomasella et al., 2000). Data mining techniques including artificial neural networks (ANNs) (Bayat et al., 2013a; Bayat et al., 2013b; Gunarathna et al., 2019a; Koekkoek and Booltink, 1999; Pachepsky et al., 1996), group method of data handling (GMDH) (Bayat et al., 2011; Neyshaburi et al., 2015; Pachepsky and Rawls, 1999), nonparametric nearest neighbor technique (Botula et al., 2013; Gunarathna et al., 2019a; Haghverdi et al., 2015; Nemes et al., 2006; Nguyen et al., 2017) and support vector machine (SVM) (Khlosi et al., 2016; Lamorski et al., 2008; Lamorski et al., 2014; Twarakavi et al., 2009), have been applied successfully for PTF development.

Random forest (RF), or random decision forests, has become a popular approach as an ensemble learning method for prediction and classification (Verikas et al., 2011). The RF method has been developed by Breiman (2001) as an expansion of the classification and regression trees (CART) technique to provide better performance of prediction results (Wiesmeier et al., 2011). So far, few studies have been carried out on the application of the RF method in soil science. Tóth et al. (2014) applied the RF method to analyze the relationship between soil water content at four matric suctions (0.1, 33, and 1500 kPa, and 150 MPa) and Hungarian soil map information. They found that the importance of soil properties in the prediction of the soil water content varied according to soil type and matric suction. Recently Szabó et al. (2019) have developed PTFs based on RF and geostatistics methods to map soil hydraulic properties, such as water contents at saturation, field capacity and wilting point, for the Balaton catchment area in Hungary. Araya and Ghezzehei (2019) compared the performances of four machine-learning algorithms including the k-nearest neighbors (kNNs), support vector regression (SVR), RF, and boosted regression tree (BRT) for prediction of saturated hydraulic conductivity. They found that the BRT model outperformed the other algorithms closely followed by the RF model. Gunarathna et al. (2019a) tested three machine-learning algorithms including ANN, kNN, and RF to estimate volumetric water content at matric suctions of 10, 33 and 1500 kPa for soils in Sri Lanka. They recommended that the PTFs to be developed using the RF algorithm. Ließ et al. (2012) studied uncertainty in the spatial prediction of soil texture by comparison of the RF and regression tree techniques for 56 soil profiles and found that the former method provided a better result. Also, Wiesmeier et al. (2011) utilized the RF technique to develop digital mapping of the soil organic matter content in 120 soil profiles. They found that the prediction accuracy of the RF modeling was acceptable. A review of literatures therefore revealed that the RF data mining technique has been applied to develop PTFs to predict specific points of the SWRC, such as field capacity and permanent wilting point, or particular properties such as saturated hydraulic conductivity, but it has not been used to develop parametric-based PTFs of the van Genuchten model parameters, so far. Therefore, the objective of the present study was to develop simple parametric-PTFs to predict the SWRC with greater accuracy and reliability using a novel approach with the RF data mining technique. We compare its performance with those of the multiple NLR approach and with Rosetta software (Schaap et al., 2001) on the prediction of the SWRC through finding the best

input variables and PTFs for the SWRC.

2. Materials and methods

2.1. Sample collection and determination

In the present study 223 undisturbed and disturbed soil samples were taken from six provinces of Iran including west Azarbaijan ($35^\circ 8' - 39^\circ 46' \text{N}$, $44^\circ 3' - 47^\circ 23' \text{E}$; 60 data), Hamedan ($33^\circ 59' - 35^\circ 48' \text{N}$, $47^\circ 34' - 49^\circ 36' \text{E}$; 55 data), Kermanshah ($33^\circ 41' - 35^\circ 17' \text{N}$, $45^\circ 24' - 48^\circ 6' \text{E}$; 26 data), Kurdistan ($34^\circ 45' - 36^\circ 31' \text{N}$, $45^\circ 31' - 48^\circ 13' \text{E}$; 22 data), Mazandaran ($35^\circ 46' - 36^\circ 58' \text{N}$, $50^\circ 21' - 58^\circ 08' \text{E}$; 30 data) and Fars ($27^\circ 2' - 31^\circ 42' \text{N}$, $50^\circ 42' - 55^\circ 38' \text{E}$; 30 data). Steel cylinders, measuring 5.1 cm in diameter and 3.5 cm in height, were used to collect the undisturbed samples. Since the sampling was done from different locations of the various provinces, the topsoil and subsoil layers of soil at different locations had different depths and thicknesses. We collected samples from the center of the topsoil and subsoil layers, which represented the pedological A and B horizons, respectively. The sampling depths varied from 10 to 35 cm for topsoil (208 samples) and from 20 to 45 cm for subsoil (15 samples), reflecting the variation in the soil profiles.

Soil PSD was analyzed by the hydrometer method (Gee and Or, 2002), and the geometric mean and standard deviation of particle diameter (d_g and δ_g , respectively) were calculated by equations from Shirazi and Boersma (1984). Organic matter (OM) content was determined by the Walkley and Black (1934) method and BD by the core method (Blake and Hartge, 1986). Total porosity (TP) was calculated from BD and particle density, and the saturated hydraulic conductivity (K_s) was measured with a constant head permeameter (Klute and Dirksen, 1986). The SWRC was constructed by measuring the volumetric water content at matric suctions of 0 (saturation status of soil samples), 1, 2 and 5 kPa with a sandbox apparatus, and at 10, 25, 50, 100, 200, 500, 1000 and 1500 kPa with a pressure plate apparatus. Undisturbed samples were used for measurement of the matric suctions from 0 to 100 kPa and disturbed samples were used for matric suctions from 200 to 1500 kPa. Two key points in the SWRC are the water contents at field capacity (30 kPa suction; θ_{FC}) and permanent wilting point (1500 kPa suction; θ_{PWP}).

2.2. Soil-water retention equation

The van Genuchten–Mualem (Eq. (1)) model (Mualem, 1976; van Genuchten, 1980) was utilized to describe the SWRC data.

$$\theta = \theta_r + (\theta_s - \theta_r) \times \frac{1}{[1 + (\alpha h)^n]^{(1-\frac{1}{n})}} \quad (1)$$

where θ_r and θ_s are residual and saturated water contents ($\text{cm}^3 \text{cm}^{-3}$), respectively, and h is the soil water suction (kPa). The parameter α is related to the inverse of the air entry pressure (> 0 , kPa^{-1}) and n (> 1 , dimensionless parameter) is related to the pore size distribution of the soil (van Genuchten, 1980). In the present study, van Genuchten model parameters θ_r , θ_s , α and n were obtained using the MATLAB software (MathWorks, 2018).

2.3. Data pre-processing

Data pre-processing and regression assumptions, including detection of outliers, normality test of the residuals, multicollinearity and independence of the residuals, were applied for all variables (Berry, 1993). The outliers in the data were identified by the inter-quartile range (IQR) method (Seo, 2006) and were replaced by the lower and upper threshold values (MathWorks, 2018). Before developing PTFs, all variables were evaluated by Kolmogorov-Smirnov normality and multicollinearity tests by the SPSS 24 software (IBM, 2016). The degree of multicollinearity in the PTFs was tested by the variance inflation factor ($VIF = 1/1-R_j^2$, where R_j^2 is the R^2 value obtained by regressing the j^{th} predictor on the remaining predictors) (Hocking, 2013). Also, to avoid multicollinearity between textural contents, the silt fraction was not used as a predictor. The variables clay content, sand content, d_g , δ_g , OM, K_s , α and n had non-normal distributions, therefore, transformations were applied to normalize them.

2.4. Developing PTFs

The PTF inputs were arranged in four steps (Fig. 1). The first step (PTFs 1–5) was based on basic soil properties (i.e., sand content (%), clay content (%), BD (g cm^{-3}), θ_{FC} ($\text{cm}^3 \text{cm}^{-3}$) and θ_{pwp} ($\text{cm}^3 \text{cm}^{-3}$)) according to Rosetta-based PTFs (Schaap et al., 2001) for comparison of SWRC estimates by other methods. The parameters of the van Genuchten model were predicted in all steps. In the second step (PTFs 6–9), d_g (mm) and δ_g were used as new inputs instead of sand and clay contents in the previous step to evaluate the efficiency of using statistical descriptors of PSD to predict the parameters of the van Genuchten model. To build the third step (PTFs 10–12), TP ($\text{cm}^3 \text{cm}^{-3}$) replaced

BD from PTFs 3–5 to evaluate the effect of using TP instead of BD on the prediction of the parameters of the van Genuchten model. In other words, the purpose of the second and third steps was to evaluate whether the use of another form of descriptors of soil structure (TP instead of the BD) and soil texture (d_g and δ_g instead of the sand and clay contents) would improve the accuracy of the estimates or not. In the last step, PTFs 13–15 were developed by including OM (%) and K_s (cm day^{-1}) as new variables to evaluate the efficiency of these instead of the water content at specific matric suctions on the prediction of the van Genuchten model parameters. The input variables of the 15 PTFs are shown in Fig. 1.

To compare the results of PTFs 1–5 of the RF and NLR methods with those of the Rosetta models, the parameters of the van Genuchten model (θ_r , θ_s , α and n) were estimated by the PTFs built in the Rosetta software (PTFs 1–5), using the measured values of input variables based on PTFs 1–5 as predictors in the Rosetta program. The estimated coefficients of the van Genuchten model were used to calculate the estimated water content at matric suctions from 0 to 1500 kPa (estimated SWRCs). Then curve-by-curve comparison of the measured and estimated SWRCs was performed with different evaluation statistics. Since there is no training step in the Rosetta software, the results of the Rosetta model was only compared with the results of the testing step. To evaluate the effect of using different descriptors of PSD on the prediction of the SWRC, PTFs 6, 7, 8 and 9 from the second step were compared with PTFs 2, 3, 4 and 5 from the first step, respectively (Fig. 1). In the same way, to evaluate effect of using different descriptors of soil structure on the prediction of the SWRC, PTFs 10, 11 and 12 from the third step were compared with PTFs 3, 4 and 5 from the first step, respectively. Also, the PTFs 13–15 were compared with the PTFs 4 and 5 to find out the efficiency of OM and K_s variables as

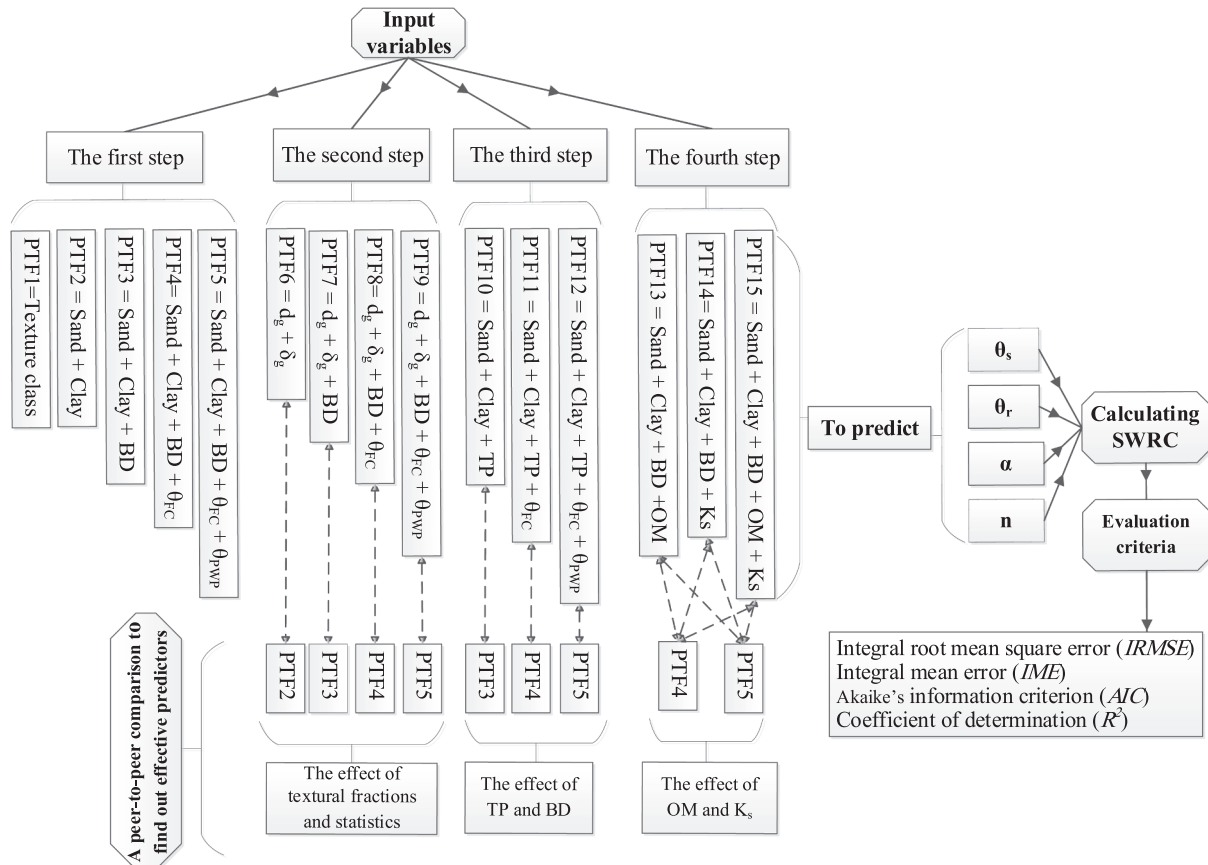


Fig. 1. Input variables of the 15 pedotransfer functions (PTFs) for predicting the van Genuchten model parameters (θ_r , θ_s , α and n) of the soil water retention curve (SWRC). A list of abbreviations is available in the notation box.

predictors (Fig. 1).

In the present study, the k-fold cross validation approach (Efron and Tibshirani, 1994) was utilized to obtain training and testing datasets for each PTF. The number of folds (i. e., k) was obtained by trial and error. To do so, some PTFs, selected randomly, were developed with 10, 15 and 20-fold cross-validation. Then, the k value which resulted in the best performance of the PTFs, was selected to develop all PTFs in this study. The results showed that 20-fold cross validation performed better than the other folds in most of the PTFs (Table 1). Therefore, 20-fold cross validation was selected to develop PTFs in this study. Based on this approach, the 223 samples were randomly divided into 20 subsets and 20 models were developed by each predicting technique for each PTF. In each model, training and testing datasets were based on a ratio of 19:1. Finally, the average of the results of 20 models was calculated for each PTF. Therefore, all data were used for the training and testing steps of the PTFs.

2.5. Description of modeling techniques

2.5.1. Multiple nonlinear regression

A NLR model based on a second-order polynomial for the prediction of the response variable y from a number of p predictors can be written as (Rawls and Brakensiek, 1985):

$$y = a + \sum_{i=1}^p (b_i x_i + c_i x_i^2) \quad (2)$$

where a is the intercept, and two regression coefficients b_i and c_i are determined for every input variable x_i .

2.5.2. Random forest: An ensemble of regression trees

RF has become a popular tool for regression and classification problems. The RF is an ensemble method based on the regression tree methodology (i.e., CART) that was introduced for better performance (Breiman, 2001). The model building process in the RF is the same as that in the CART method but without pruning (Breiman, 1984). Also, whereas a regression tree only grows by a single tree the RF grows by forest of trees. In other words, unlike a regression tree, in the RF for each tree only a subset of the input variables is applied. The number of inputs in each tree and also the number of trees in the forest can be distinct and it depends on the dataset. Least-squares boosting (LSBoost) fits regression ensembles. At every step, the ensemble fits a new learner

to the difference between the observed response and the aggregated prediction of all learners grown previously. The ensemble fits to minimize the mean-squared error (MathWorks, 2018). The number of trees used here was 16 which was established by trial and error. An architecture of the RF algorithm is shown in Fig. 2 where input matrix X consists of N samples and M input variables (sample set $S = [(x_i, y_i), i = 1, 2, \dots, N]$, $(X, Y) \in R^M \times R$). The bootstrap method is utilized to construct n tree sample sets from the sample set S . At each bootstrap sample, about one-third of the dataset S was utilized as out of the bootstrap data or out-of-bag (OOB) data; whereas the rest is called in-bag data (Ibrahim and Khatib, 2017) (Fig. 2). Modeling of the regression tree is done for each sample set. In the RF algorithm, all individual trees give a predictive result. The final prediction value is calculated based on an average result of all individual trees (Wiesmeier et al., 2011). The prediction error is defined as follows (Liaw and Wiener, 2002):

$$MSE_{OOB} = \frac{\sum_{i=1}^{n_{tree}} (y_i - \hat{y}_i^{OOB})^2}{n_{tree}} \quad (3)$$

where MSE_{OOB} is the mean square error of the OOB data prediction, n_{tree} is the number of trees, and y_i and \hat{y}_i^{OOB} are the actual value of the OOB data and the average of all OOB predictions, respectively. Among all the ensemble methods, the RF method has high capability in solving classification and regression problems, because the RF method combines several simple regression trees to better optimize prediction (Zaklouta and Stanculescu, 2012). The RF method increases differences for each single tree through random selection of the training samples and different variables at each splitting node. In the present study, the NLR and RF algorithms were implemented by fitnlm and fitensemble functions in the MATLAB software, respectively. (MathWorks, 2018).

2.6. Evaluation criteria

The estimated water content was computed by estimated parameters of the van Genuchten model for each PTF at matric suctions from 0 to 1500 kPa. For curve-by-curve comparison of the measured and predicted SWRCs, different evaluation statistics were used. Various statistical criteria including integral root mean square error (IRMSE), integral mean error (IME) (Tietje and Tapkenhinrichs, 1993), Akaike's information criterion (AIC) (Akaike, 1974) and coefficient of

Table 1

The results of 10, 15 and 20-fold cross-validation (k) for van Genuchten model parameters of the soil water retention curve derived from nonlinear regression (NLR) and random forest (RF) techniques based on root mean square error (RMSE) for pedotransfer functions PTF 3, 5 and 11 in the train and test datasets.

| | | | θ_r | | | θ_s | | | α | | | n | | |
|-------|--------|-----|------------|-------|-------|------------|-------|-------|----------|-------|-------|-------|-------|-------|
| | | | RMSE | | | RMSE | | | RMSE | | | RMSE | | |
| | | | Train | Test | Mean | Train | Test | Mean | Train | Test | Mean | Train | Test | Mean |
| PTF3 | k = 10 | NLR | 0.058 | 0.060 | 0.059 | 0.063 | 0.065 | 0.064 | 1.017 | 1.037 | 1.027 | 0.426 | 0.436 | 0.431 |
| | | RF | 0.052 | 0.061 | 0.056 | 0.058 | 0.073 | 0.066 | 0.893 | 1.084 | 0.989 | 0.374 | 0.442 | 0.408 |
| | k = 15 | NLR | 0.058 | 0.060 | 0.059 | 0.064 | 0.064 | 0.064 | 1.017 | 1.030 | 1.024 | 0.426 | 0.434 | 0.430 |
| | | RF | 0.052 | 0.061 | 0.057 | 0.058 | 0.070 | 0.064 | 0.894 | 1.033 | 0.964 | 0.374 | 0.441 | 0.408 |
| | k = 20 | NLR | 0.058 | 0.060 | 0.059 | 0.064 | 0.064 | 0.064 | 0.064 | 0.064 | 0.064 | 0.426 | 0.437 | 0.432 |
| | | RF | 0.051 | 0.060 | 0.056 | 0.057 | 0.071 | 0.064 | 0.057 | 0.071 | 0.064 | 0.368 | 0.442 | 0.405 |
| PTF5 | k = 10 | NLR | 0.051 | 0.053 | 0.052 | 0.053 | 0.054 | 0.054 | 0.764 | 0.796 | 0.780 | 0.380 | 0.397 | 0.389 |
| | | RF | 0.043 | 0.056 | 0.050 | 0.046 | 0.056 | 0.051 | 0.675 | 0.869 | 0.772 | 0.327 | 0.411 | 0.369 |
| | k = 15 | NLR | 0.051 | 0.053 | 0.052 | 0.053 | 0.055 | 0.054 | 0.764 | 0.790 | 0.777 | 0.381 | 0.399 | 0.390 |
| | | RF | 0.044 | 0.054 | 0.049 | 0.046 | 0.055 | 0.050 | 0.679 | 0.848 | 0.763 | 0.329 | 0.421 | 0.375 |
| | k = 20 | NLR | 0.051 | 0.053 | 0.052 | 0.053 | 0.055 | 0.054 | 0.765 | 0.789 | 0.777 | 0.381 | 0.399 | 0.390 |
| | | RF | 0.042 | 0.054 | 0.048 | 0.044 | 0.054 | 0.049 | 0.654 | 0.842 | 0.748 | 0.316 | 0.412 | 0.364 |
| PTF11 | k = 10 | NLR | 0.058 | 0.061 | 0.060 | 0.065 | 0.067 | 0.066 | 1.018 | 1.052 | 1.035 | 0.431 | 0.448 | 0.440 |
| | | RF | 0.050 | 0.061 | 0.056 | 0.047 | 0.057 | 0.052 | 0.770 | 0.978 | 0.874 | 0.370 | 0.443 | 0.406 |
| | k = 15 | NLR | 0.058 | 0.061 | 0.060 | 0.065 | 0.067 | 0.066 | 1.019 | 1.037 | 1.028 | 0.432 | 0.447 | 0.439 |
| | | RF | 0.050 | 0.060 | 0.055 | 0.047 | 0.057 | 0.052 | 0.770 | 1.009 | 0.889 | 0.369 | 0.450 | 0.410 |
| | k = 20 | NLR | 0.058 | 0.060 | 0.059 | 0.065 | 0.065 | 0.065 | 1.020 | 1.024 | 1.022 | 0.432 | 0.439 | 0.435 |
| | | RF | 0.049 | 0.061 | 0.055 | 0.046 | 0.056 | 0.051 | 0.745 | 0.964 | 0.855 | 0.361 | 0.443 | 0.402 |

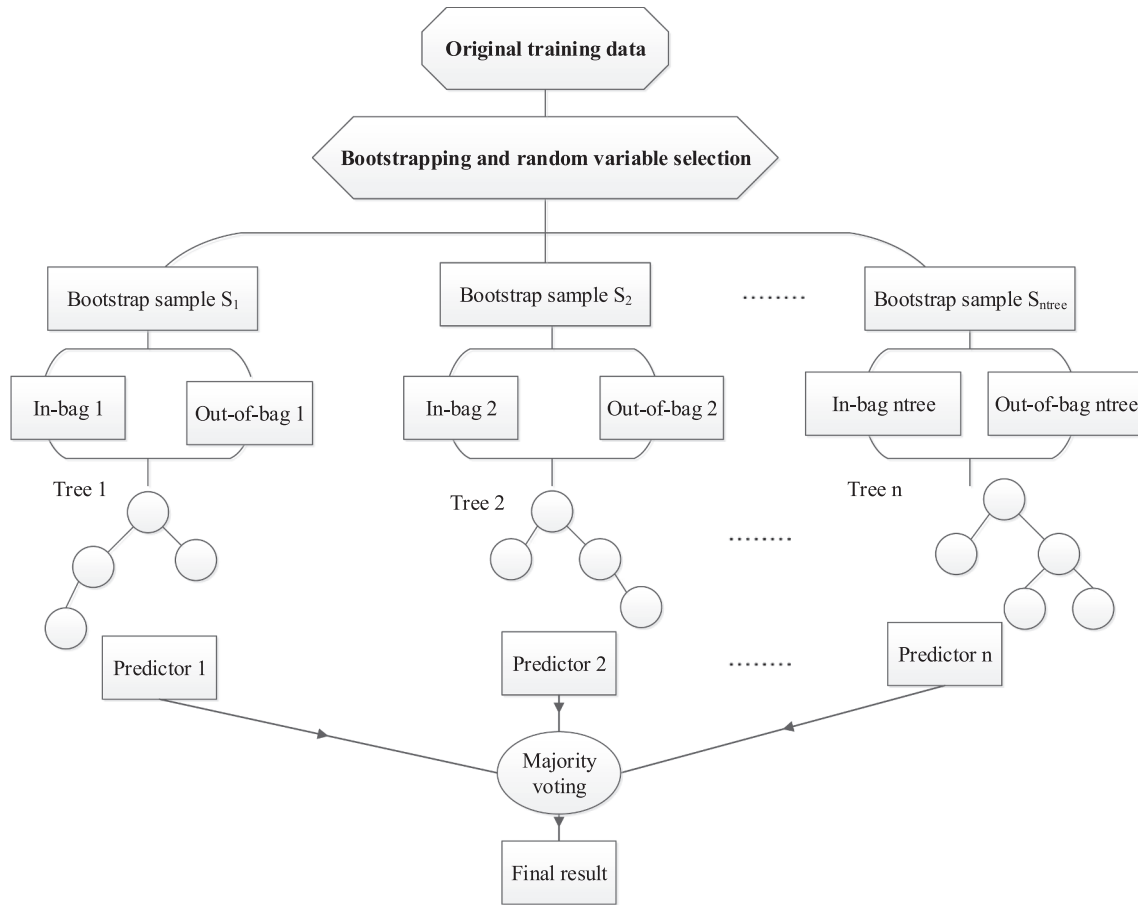


Fig. 2. An architecture of a random forest.

determination (R^2) (Wösten et al., 2001), were utilized to assess the predictive ability of the RF and NLR algorithms, which are defined as:

$$IRMSE (cm^3 cm^{-3}) = \left[\frac{1}{b-a} \int_a^b (\hat{y}_i - y_i)^2 d \log |h| \right]^{\frac{1}{2}} \quad (4)$$

$$IME (cm^3 cm^{-3}) = \frac{1}{b-a} \int_a^b (\hat{y}_i - y_i) d \log |h| \quad (5)$$

$$AIC = N \times \ln \left[\sum_{i=1}^N \frac{(y_i - \hat{y}_i)^2}{N} \right] + 2P \quad (6)$$

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y}_i)^2} \quad (7)$$

where h is the matric suction (kPa), y_i , \hat{y}_i and \bar{y}_i are the measured, predicted and average of the measured values of the water content, respectively, a and b values define the matric suction range over which the experimental curve is measured, i.e., 0 and 1500 kPa, respectively, and P and N are the number of parameters and the number of points that were considered in the SWRC, respectively. In calculating the AIC, N is the total number of points that were considered in the SWRC of all soil samples (i. e., N = number of soil samples \times number of paired points of the suction-water content for each soil sample), and i is paired points of the suctions-water content of the SWRC of each soil sample.

To evaluate the performance of each method in different PTFs, the effect of method as the first factor at two levels in the training step (i.e., NLR and RF methods) and at three levels in the testing step (i.e., NLR,

RF and Rosetta methods), and the different PTFs as the second factor at 15 levels (PTF1 to PTF15), were investigated using a two-way analysis of variance (ANOVA) with a randomized complete block design, based on the *IRMSE* of prediction of the SWRC. The *IRMSE* criterion calculates the total error, including bias and random errors, and is a more appropriate criterion for evaluating the accuracy and reliability of the RF and NLR methods compared to other criteria (Chai and Draxler, 2014). Therefore, to compare the predicting accuracy and reliability of the RF and NLR methods, the average values of the *IRMSE* was compared with Duncan's test by MathWorks (2018) software.

3. Results and discussion

3.1. Descriptive statistics of the soil properties

Table 2 summarizes some basic descriptive statistics for soil variables of the entire dataset used for the development of the PTFs. It can be seen that the average and maximum of clay content were 21.4 and 48%, respectively. The OM ranged from 0.17 to 4.41% with a mean of 1.84%, which was low due to the arid and semi-arid climates of Iran. The variation in soil texture is shown graphically in the United States Department of Agriculture (USDA) textural triangle (Fig. 3). Considering the distribution and range of the variables (Fig. 3 and Table 2), the dataset can be considered as representative of soils in arid and semi-arid regions of Iran.

3.2. Correlation of input and output variables

The simple correlation coefficients between all variables are depicted by matrix plot in Fig. 4. Correlation analysis was done between

Table 2

Some descriptive statistics of the measured soil variables and parameters of the van Genuchten model of the soil water retention curve for the entire dataset (223 soil samples).

| Variables ^a | Mean | CV (%) | Minimum | Maximum | P-value |
|--|--------|-----------|---------|---------|---------|
| Clay content (%) | 21.39 | 54.05 | 3.47 | 48.00 | 0.00 |
| Log (clay content) | 1.27 | 19.08 | 0.54 | 1.68 | 0.66 |
| Sand content (%) | 35.45 | 48.40 | 5.90 | 89.80 | 0.00 |
| Sand content* | -0.01 | -14350.94 | -3.40 | 3.14 | 0.90 |
| Bulk density (g cm ⁻³) | 1.43 | 10.97 | 1.03 | 1.84 | 0.83 |
| θ_{FC} (cm ³ cm ⁻³) [§] | 0.33 | 20.44 | 0.15 | 0.55 | 0.45 |
| θ_{PWP} (cm ³ cm ⁻³) | 0.18 | 26.21 | 0.04 | 0.31 | 0.90 |
| d_g (mm) | 0.07 | 86.62 | 0.00 | 0.21 | 0.00 |
| Log (d_g) | -1.33 | -27.91 | -2.34 | -0.67 | 0.77 |
| δ_g (-) | 11.57 | 29.39 | 4.54 | 19.97 | 0.00 |
| δ_g^* | -0.01 | -9872.87 | -2.53 | 1.80 | 0.96 |
| Total porosity (cm ³ cm ⁻³) | 0.46 | 13.26 | 0.31 | 0.61 | 0.67 |
| Organic matter content (%) | 1.84 | 53.68 | 0.17 | 4.41 | 0.00 |
| (Organic matter content) ^(1/4) | 1.13 | 14.83 | 0.64 | 1.45 | 0.86 |
| K_s (cm day ⁻¹) | 169.10 | 96.58 | 0.06 | 530 | 0.00 |
| (K_s) ^(1/4) | 3.23 | 30.37 | 0.50 | 4.80 | 0.59 |
| θ_r (cm ³ cm ⁻³) | 0.04 | 158.05 | 0.00 | 0.17 | 0.00 |
| θ_s (cm ³ cm ⁻³) | 0.52 | 16.26 | 0.35 | 0.70 | 0.56 |
| α (kPa ⁻¹) | 0.06 | 115.62 | 0.00 | 0.29 | 0.00 |
| α^* | 0.01 | 8889.14 | -2.93 | 2.19 | 0.93 |
| n | 1.24 | 9.80 | 1.08 | 1.48 | 0.00 |
| Ln ($n-1$) | -1.55 | -30.92 | -2.52 | -0.74 | 0.05 |

*Normalized form of sand content: $0.91 + 1.06 \times \ln((\text{sand content} - 4.3)/(100.2 - \text{sand content}))$; normalized form of δ_g : $-1.04657 + 1.39359 \times \text{Asinh}((\delta_g - 8.4)/3.04)$; and normalized form of α : $3.6 + 0.92 \times \ln((\alpha - 8.2 \times 10^{-6})/(1.6 - \alpha))$. P-value is a significance value for normality test.

§. A list of abbreviations is available in the notation box.

^a CV, coefficient of variation.

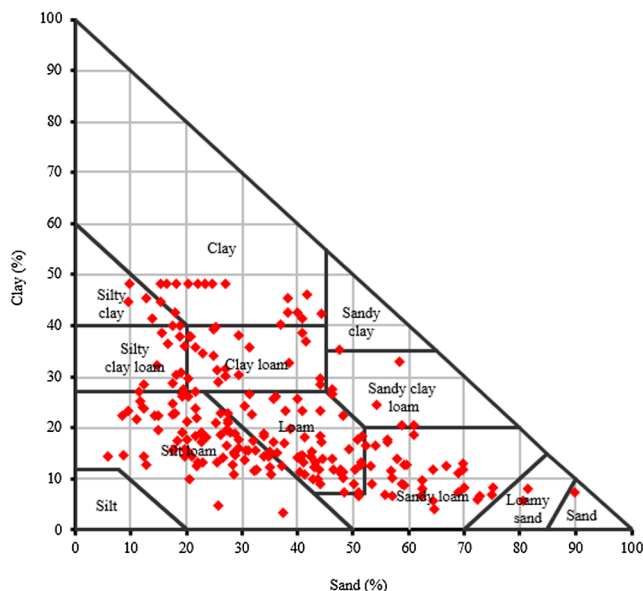


Fig. 3. Variation of soil texture classes for the dataset ($n = 223$) on the United States Department of Agriculture (USDA) textural triangle.

normalized input and output variables. The correlation test was not performed for the θ_r variable, because its value was zero in 138 out of 223 soil samples, as has been reported in other studies (Campbell and Horton Jr, 2002; Rawls et al., 1991; Tomasella et al., 2000) for θ_r variable. Clay and sand contents, θ_{FC} , θ_{PWP} , d_g and OM had the greatest significant correlations with the parameters of the van Genuchten model (Fig. 4), which was consistent with other studies (Dexter et al.,

2008; Nemes et al., 2006). For example, the correlation coefficient between clay content and θ_s ($r = 0.323$) is close to that between OM and θ_s ($r = 0.268$). Also, the results showed that there were significant correlations between θ_{PWP} and input variables of clay content (+), sand content (-), BD (-), OM (+) and K_s (-), and also between θ_{PWP} and θ_s (+) and n (-) parameters of the van Genuchten model (Fig. 4). Botula et al. (2012) also found the same observation for the correlation of θ_{PWP} with sand and clay contents and BD of tropical Lower Congo soils. Nevertheless, with regard to these correlation coefficients, clay and sand contents, θ_{FC} , d_g and OM can be used for developing PTFs to estimate the SWRC. On the contrary, there was no correlation between K_s and the van Genuchten model parameters. There are many cases, where two variables might not show a strong simple correlation, but may show a strong association in the regression, along with other predictors. In other words, the simple correlation coefficient is a way to show the relationship between independent and dependent variables, but it cannot show a model for the relationship between these two variables, when other independent variables have been used in a multiple regression (Simmons et al., 2011). The result of multiple regression analysis with backward selection method showed that the K_s variable remained in the PTF14 and PTF15 for all the van Genuchten model parameters. Some of the regression equations with backward selection method are shown in the following as examples:

$$\theta_r = -0.69 + 0.22 \times \text{Clay} + 0.278 \times \text{Sand} + 0.20 \times K_s, R = 0.31 ** \quad (8)$$

$$\alpha = -3.72 + 0.23 \times \text{Clay} + 0.17 \times \text{BD} + 0.282 \times K_s, R = 0.33 ** \quad (9)$$

$$n = -1.76 + 0.24 \times \text{Sand} + 0.164 \times K_s, R = 0.30 ** \quad (10)$$

On the other hand, the non-linear correlations between variables are very important in this study. Both the multiple NLR approach and RF data mining technique are non-linear prediction methods. Fig. 4 only shows simple linear correlation between variables, but there may be non-linear correlations between variables, which may affect the estimation of the dependent variables. For example, the results of non-linear correlations showed that K_s had strong correlations with θ_s and α of the van Genuchten model parameters by logarithmic ($\theta_s = 0.652 - 0.027 \times \ln K_s$, $R = 0.62 **$) and power ($\alpha = 0.007 \times K_s^{0.283}$, $R = 0.57 **$) equations, respectively, which were greater than their simple correlations

3.3. Development of the PTFs using the RF and NLR methods

Results of the multicollinearity analysis (VIF) are shown in Table 3. The VIF values showed low levels of multicollinearity among the independent variables ($VIF < 10$) (Khodaverdiloo et al., 2011).

3.3.1. Comparing the accuracy and reliability of the RF and NLR methods

Table 4 shows the results of the ANOVA of the IRMSE of prediction of the SWRC by different methods and PTFs. The effect of methods and PTFs, and their interaction, on the IRMSE was significant at $P < 0.01$, 0.01 and 0.05, respectively, in the training step, and at $P < 0.01$, 0.01 and 0.01, respectively, in the testing step. Therefore, we focus on the results and discussion of the comparison of the method \times PTF interaction effects.

Results of the prediction of the SWRC through the van Genuchten model using the NLR and RF-based PTFs are depicted in Figs. 5 and 6 for the training and testing steps, respectively. The accuracy and reliability are used to express the performance of the PTFs in the training and testing steps, respectively.

The results of the first to fourth steps of the training dataset (Fig. 5) showed that the RF method had better performance compared to the NLR method for the prediction of the SWRC in all PTFs in terms of the IRMSE and R^2 criteria and the differences were significant ($P < 0.05$) for PTFs 2, 3, 6, 7, 10, 13, 14 and 15 in terms of the IRMSE criterion.

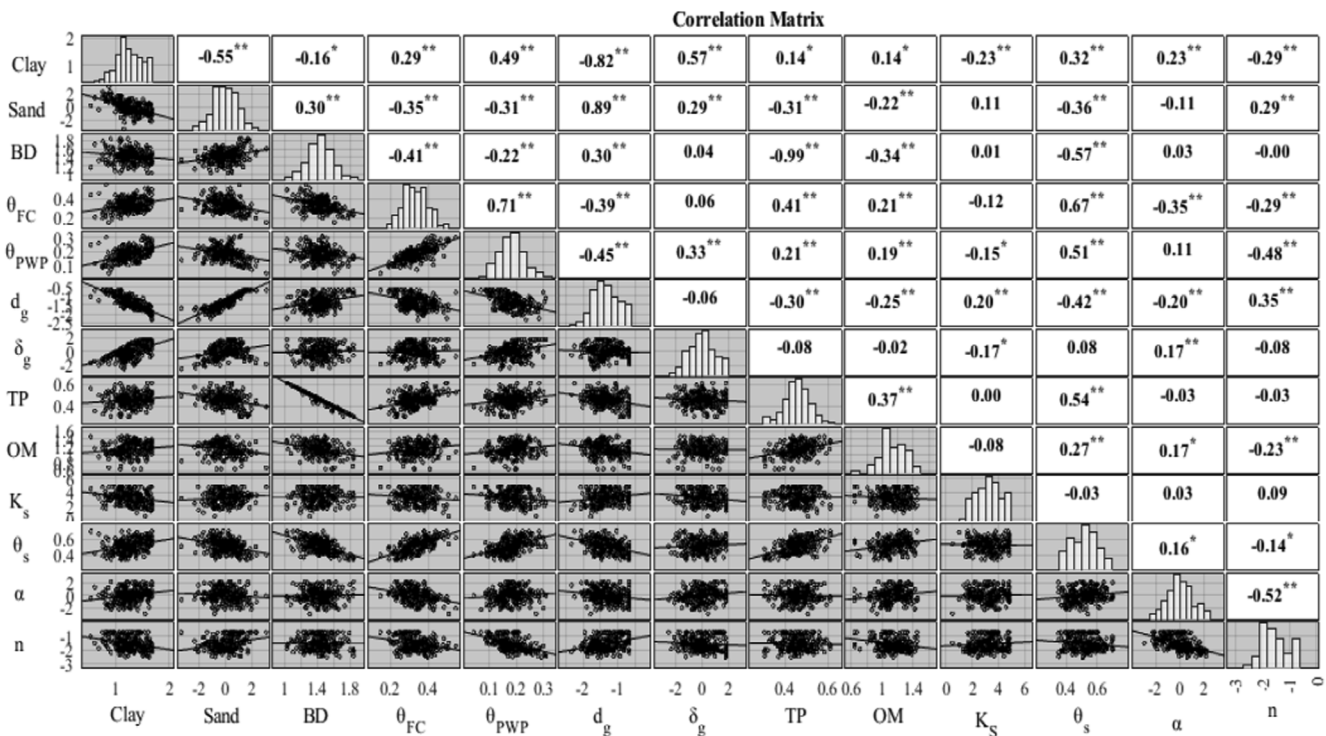


Fig. 4. Correlation matrix plot between input and output variables. ** Correlation is significant at the $P < 0.01$ level. * Correlation is significant at the $P < 0.05$ level. A list of abbreviations is available in the notation box.

Also, the accuracy of the RF method was better than that of the NLR method in 80% of the PTFs (with the exception of the PTFs 5, 9 and 12) in terms of the *AIC* criterion. In the training step, the values of the *IRMSE* of the first to fourth steps for the NLR model varied from 0.030 to 0.063 $\text{cm}^3 \text{cm}^{-3}$ and these were larger than those in the RF model, which ranged from 0.028 to 0.061 $\text{cm}^3 \text{cm}^{-3}$, respectively. Also, the values of the R^2 of the first to fourth steps for the RF model varied from 0.981 to 0.992, and this was larger than those in the NLR model, which ranged from 0.979 to 0.991 (Fig. 5).

The results of the first to fourth steps of the testing dataset (Fig. 6) showed that the NLR method had a better performance compared to the RF method on the prediction of the SWRC for PTFs 5, 8, 9 and 15 only in terms of the *IRMSE* criterion (significant at $P < 0.05$). In the other PTFs there were no significant differences between the *IRMSE* of the two methods and the R^2 and *AIC* criteria were comparable. In the testing step, the values of the *IRMSE* and *AIC* of the first to fourth steps

Table 4

Analysis of variance of the integral root mean square error (*IRMSE*) of the prediction of the soil water retention curve by different methods (nonlinear regression and random forest) and pedotransfer functions (PTFs 1–15) for both the train and test datasets.

| | Source | Degree freedom | Mean square | F-value | P-value |
|-------|-----------------------|----------------|-------------|---------|----------|
| Train | Repeat (Block) | 222 | 0.007 | 19.09 | < 0.0001 |
| | PTFs | 14 | 0.062 | 180.68 | < 0.0001 |
| | Methods | 1 | 0.038 | 109.69 | < 0.0001 |
| | PTFs \times Methods | 14 | 0.001 | 1.78 | 0.0356 |
| | Error | 6288 | 0.0003 | | |
| Test | Repeat (Block) | 222 | 0.010 | 16.04 | < 0.0001 |
| | PTFs | 14 | 0.073 | 117.22 | < 0.0001 |
| | Methods | 2 | 0.656 | 1056.43 | < 0.0001 |
| | PTFs \times Methods | 18 | 0.002 | 3.68 | < 0.0001 |
| | Error | 7398 | 0.0006 | | |

Table 3

The variance inflation factor (*VIF*) values for normalized form of the input variables.

| PTFs | Clay* (%) | Sand (%) | BD [§] (g cm^{-3}) | θ_{FC} ($\text{cm}^3 \text{cm}^{-3}$) | θ_{PWP} ($\text{cm}^3 \text{cm}^{-3}$) | d_g (mm) | δ_g (–) | TP ($\text{cm}^3 \text{cm}^{-3}$) | OM (%) | K_s (cm day^{-1}) |
|-------|-----------|----------|--|--|---|------------|----------------|-------------------------------------|--------|--------------------------------|
| PTF2 | 1.42 | 1.42 | | | | | | | | |
| PTF3 | 1.43 | 1.52 | 1.10 | | | | | | | |
| PTF4 | 1.45 | 1.56 | 1.25 | 1.31 | | | | | | |
| PTF5 | 1.79 | 1.58 | 1.27 | 2.48 | 2.56 | | | | | |
| PTF6 | | | | | | 1.00 | 1.00 | | | |
| PTF7 | | | 1.11 | | | 1.11 | 1.01 | | | |
| PTF8 | | | 1.25 | 1.33 | | 1.01 | 1.22 | | | |
| PTF9 | | | 1.28 | 2.50 | 2.73 | 1.34 | 1.22 | | | |
| PTF10 | 1.55 | 1.43 | | | | | | 1.11 | | |
| PTF11 | 1.58 | 1.46 | | 1.32 | | | | 1.26 | | |
| PTF12 | 1.60 | 1.79 | | 2.49 | 2.56 | | | 1.28 | | |
| PTF13 | 1.48 | 1.65 | 1.25 | | | | | | 1.14 | |
| PTF14 | 1.55 | 1.64 | 1.14 | | | | | | | 1.06 |
| PTF15 | 1.55 | 1.65 | 1.25 | | | | | | 1.15 | 1.06 |

*Normalized form of the input variables is available in Table 2.

§ A list of abbreviations is available in the notation box.

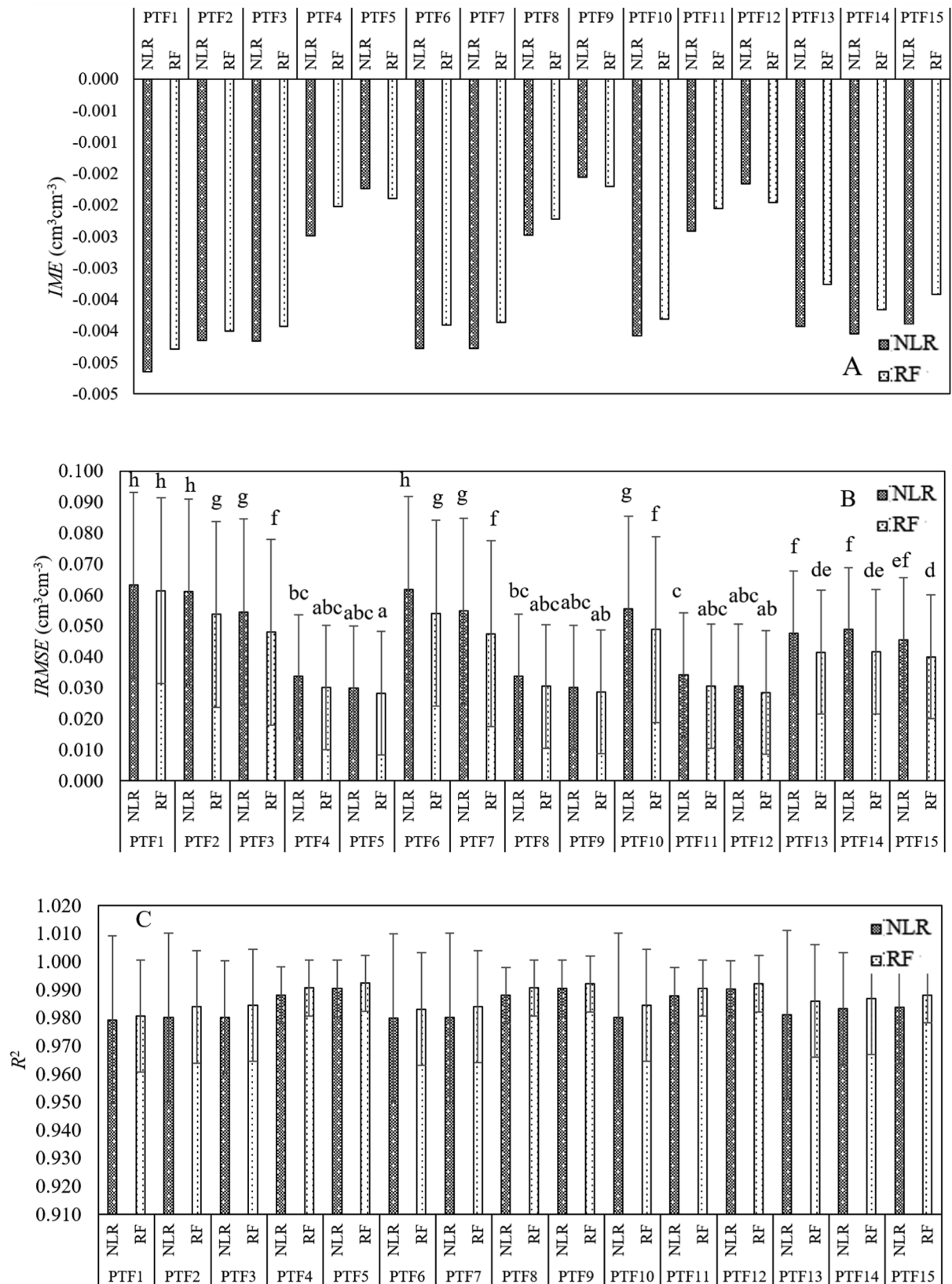


Fig. 5. Results of the prediction of the soil water retention curve (SWRC) through the van Genuchten model by the nonlinear regression (NLR) and random forests (RF) techniques for the training step as reflected in the integral mean error (IME), integral root mean square error (IRMSE), coefficient of determination (R^2), and Akaike's information criterion (AIC). Vertical lines indicate the standard deviations. Means with the same letter are not significantly different at the significance level of $P < 0.05$ (IRMSE only).

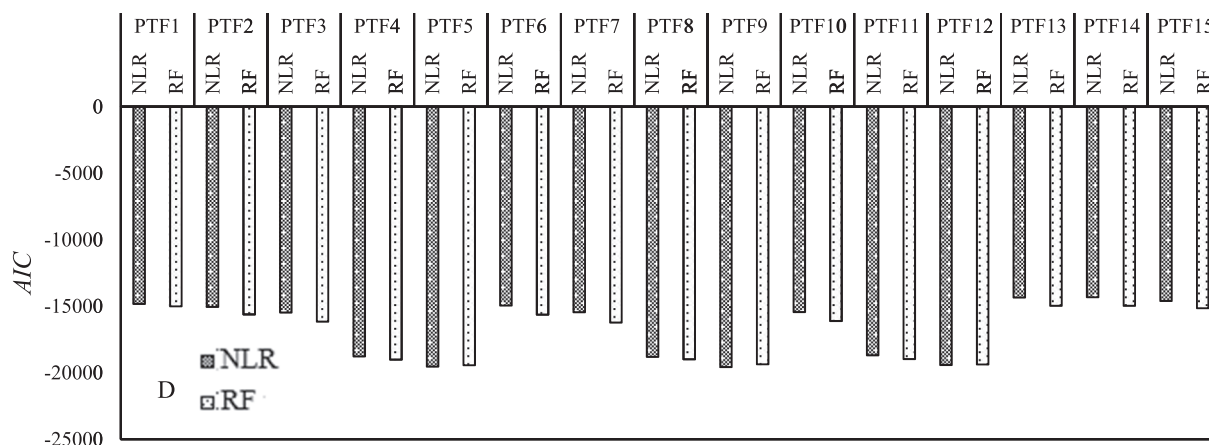


Fig. 5. (continued)

for the RF models varied from 0.038 to 0.065 $\text{cm}^3 \text{cm}^{-3}$ and from -13476.2 to -17646.8 , respectively, and these were comparable to those of the NLR models (with the exception of PTF1), which ranged from 0.032 to 0.064 $\text{cm}^3 \text{cm}^{-3}$ and from -14096.1 to -19234.1 , respectively (Fig. 6). Also, the values of the R^2 of the first to fourth steps for the NLR models varied from 0.979 to 0.989, and this was comparable to those of the RF models for all PTFs, which ranged from 0.977 to 0.987 (Fig. 6).

In each of the PTFs 1 to 5, the NLR and RF methods performed better ($P < 0.05$) than the Rosetta PTFs. Fig. 6(A) shows that the Rosetta-based PTFs had greater values of the *IME* criterion compared to the NLR and RF-based PTFs. The reason can be attributed to the various methods of optimizing parameters. The Rosetta method has only one ANN type with particular structure. In other words, the number of hidden layers (one) and neurons (six) and also the activation function (tangent hyperbolic) are constant for prediction of the SWRC in the Rosetta software. Therefore, the Rosetta method is not a dynamic approach for optimization, whereas the parameters of the RF method, such as number of splits and trees, and learning rate continuously and dynamically, change to achieve the best result of the objective function. The Rosetta method was developed from a large dataset, while the soils used in the present study were collected from a completely different climate area that was not represented in the Rosetta's database. Also, presented RF and NLR models were trained using this particular dataset while Rosetta had been trained using a different dataset. In other words, the results of the PTFs in the testing step were based on a soil dataset used for training. This could be a reason for Rosetta's poor performance compared with the RF and NLR methods. As a result, it seems that the universal portability of the Rosetta method can be limited. The testing results are in agreement with Touil et al. (2016) who found that the parametric-based PTFs of nonlinear models gave a better prediction than the Rosetta PTFs. The Fig. 5(A) and 6(A) showed that all of the *IME* values were negative for all PTFs at the training and testing steps. There are regular errors (bias) in the prediction of the SWRC that can be corrected by finding a correction coefficient, which would improve the accuracy and reliability of the estimations (Bayat et al., 2015).

The RF method in the training section gave better predictions of the SWRC compared to the NLR method (Fig. 5). The RF method produces low bias and variation in the data by majority voting compared to a single regression tree (Cheng et al., 2019; Matin and Chelgani, 2016). In this connection, the results of the standard deviations (SD) of evaluation criteria in each PTF for the training step (Fig. 5) showed that the RF method had a lower variation than the NLR method. Accordingly, the values of SD for the *IRMSE* and R^2 criteria were 0.024 and 0.022, respectively, for the NLR model and these were larger than those in the RF model, which were 0.020 and 0.017, respectively, for the training step. On the other hand, the RF method can be applied to high

dimensional datasets in regressions (Janitzka et al., 2016; Zhao et al., 2016).

As depicted in Fig. 6, unlike in the training section, the NLR method gave better predictions in the testing section compared to the RF method for the prediction of the SWRC. In other words, the reliability of the NLR method was better than that of the RF method in all the PTFs. The NLR equations can be more useful than the MLR method for the prediction of the SWRC due to their high flexibility (Williams et al., 1992). In other words, the NLR models have capacity to capture non-linear relationships in the dataset. Tomasella et al. (2000) successfully developed parametric PTFs for soils of the humid tropics using polynomials of n^{th} order. Medrado and Lima (2014) successfully developed NLR-based PTFs to predict the four parameters of the van Genuchten model for Brazilian soils. Also, Touil et al. (2016) developed parametric-PTFs to predict the SWRC using the NLR method from more readily-available properties such as soil texture, OM content, and BD for 242 soil samples of Algeria. They reported that the parametric-PTFs had better performance than Rosetta-based PTFs.

In the present study, in contrast to the NLR method which had less differences between the error values of the training and testing steps, the error values of the RF method in the testing dataset were much greater than those in the training dataset. These results can be due to overprediction phenomenon in the RF method. Gupta et al. (2017) expressed that one of the disadvantages of the RF method is the overprediction. In other words, the RF method is a 'greedy' method that easily leads to overprediction and instability in the testing step and solving this problem can be of great significance for improving the reliability of the RF method (Liu, 2014). Also, Ma et al. (2005) reported instability in results of the RF method. The forest size developed by the RF has not been clearly defined (Liu, 2014). Therefore, oversized scale can decrease the reliability and efficiency of the SWRC prediction. Hong et al. (2016) evaluated landslide susceptibility maps produced using the RF method and compared these maps with those from statistical-based methods, such as logistic regression, and their study revealed that the performance of the statistical-based methods was better than that of the RF method. A similar result was reported by Esposito et al. (2014). Generally, RFs are best suited for problems with many input variables and a reasonable sample size. According to our results (Figs. 5 and 6), performance of the PTFs was improved by increasing the number of input variables.

3.3.2. Evaluation of the effect of the basic soil properties on prediction performance of the SWRC

A significant improvement was achieved in the accuracy of PTF5 (with the inputs of Sand content + Clay content + BD + θ_{FC} + θ_{PWP}) compared to other PTFs (with the exception of PTFs 4, 8, 9, 11 and 12) by both NLR and RF methods in terms of the *IRMSE* criterion (Fig. 5).

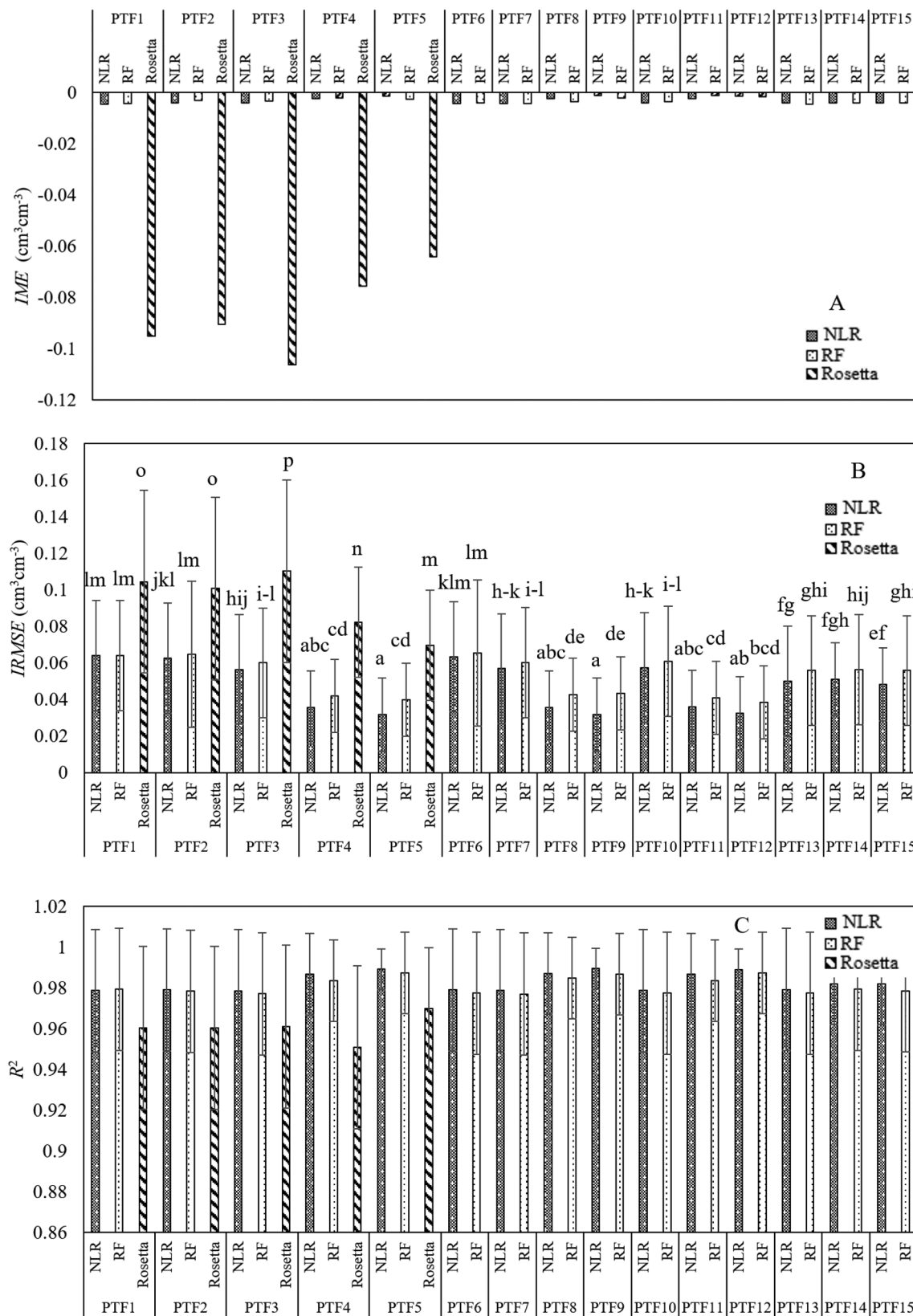


Fig. 6. Results of the prediction of the soil water retention curve (SWRC) through the van Genuchten model by the Rosetta software, nonlinear regression (NLR) and random forests (RF) techniques for the testing step as reflected in the integral mean error (IME), integral root mean square error (IRMSE), coefficient of determination (R^2), and Akaike's information criterion (AIC). Vertical lines indicate the standard deviations. Means with the same letter are not significantly different at the significance level of $P < 0.05$ (IRMSE only).

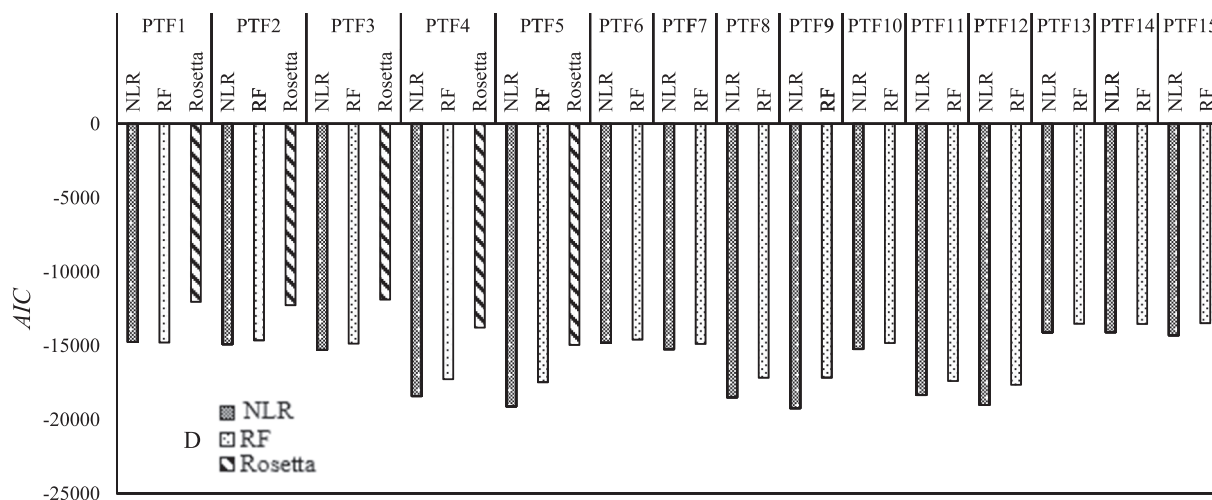


Fig. 6. (continued)

Among the PTFs of each method (RF or NLR), PTF5 had the greatest R^2 (0.992 and 0.991, respectively) and the smallest $IRMSE$ (0.028 and 0.03, respectively) and AIC (-19432 and -19571.1, respectively) in the training step of the prediction of the SWRC. In connection with the importance of input variables, an improvement was achieved in the reliability of the prediction of the SWRC by PTFs 9 (with the inputs of $d_g + \delta_g + BD + \theta_{FC} + \theta_{PWP}$) and 12 (with the inputs of Sand content + Clay content + TP + $\theta_{FC} + \theta_{PWP}$) from the second and third steps, using the NLR ($IRMSE = 0.032 \text{ cm}^3 \text{ cm}^{-3}$, $AIC = -19234.1$ and $R^2 = 0.989$) and RF ($IRMSE = 0.038 \text{ cm}^3 \text{ cm}^{-3}$, $AIC = -17646.8$ and $R^2 = 0.987$) methods, respectively, in comparison with the other PTFs of each method (Fig. 6). However, the differences of PTFs 9 and 12 were not significant ($P < 0.05$) with PTFs 4, 5, 8, 11 and 12 in the NLR method and with PTFs 4, 5, 8, 9 and 11 in the RF method, respectively, in terms of the $IRMSE$ criterion.

3.3.2.1. Effect of using different input variables of PSD and soil structure as predictors on the SWRC prediction. To evaluate the effect of using different descriptors of the PSD on the prediction of the SWRC, PTFs 2, 3, 4 and 5 (clay and sand contents) from the first step were compared with PTFs 6, 7, 8 and 9 (d_g and δ_g) from the second step, respectively. In the same way, to evaluate the effect of using different descriptors of soil structure on the prediction of the SWRC, PTFs 3, 4 and 5 (BD) were compared with PTFs 10, 11 and 12 (TP) from the third step, respectively. The accuracy and reliability of the prediction of the SWRC by both NLR and RF methods were not significantly different ($P < 0.05$) (Fig. 5B and 6B). For descriptors of soil structure, the accuracy and reliability of the prediction of the SWRC by both NLR and RF methods decreased in terms of the $IRMSE$ criterion for PTFs 10 to 12 from the third step compared to PTFs 3 to 5 (with the exception of PTFs 11 and 12 in the testing step for the RF method), respectively, when TP was used instead of BD in the list of input variables (Fig. 5B and 6B). However, the differences were not significant ($P < 0.05$).

The lack of significant differences between textural contents (clay and sand contents) and statistics (d_g and δ_g), and also between TP and BD on the SWRC prediction can be due to correlation of these parameters with the parameters of the van Genuchten model (Fig. 4). The SWRC is strongly influenced by the soil structure or pore-size distribution and soil texture at small and great matric suctions, respectively (Pachepsky et al., 2006). Therefore, input variables of the textural contents or statistics can influence the residual saturation region of the SWRC. However, soil water content at the dry end (high matric suctions) of the SWRC is primarily determined by textural contents (Hillel, 1998). Also, TP and BD are indicators of soil structure and had significant correlations with θ_s (Fig. 4). Indeed, TP was calculated by BD and particle density (Rab et al., 2011). The d_g and δ_g predictors were

derived from soil textural contents (Shirazi and Boersma, 1984). Therefore, these could be reasons for similar effects of textural contents and statistics and also TP and BD predictors on the prediction of the SWRC.

Many researchers used textural contents (Adhikary et al., 2008; Chakraborty et al., 2011; Minasny et al., 1999; Tomasella and Hodnett, 1998), d_g and δ_g (Rab et al., 2011; Scheinost et al., 1997; Ungaro et al., 2005), BD (Bayat et al., 2011; Pachepsky et al., 1998) and TP (Bayat et al., 2011; Pachepsky et al., 1998; Schaap et al., 1998) as effective predictors to derive point- and parametric-PTFs. Nemes et al. (2003), Schaap et al. (2001) and Schaap et al. (1998) reported that the variables of PTF5 have better capability on predicting the parameters of the van Genuchten (1980) model with an average $RMSE$ of 0.026, 0.044 and $0.058 \text{ cm}^3 \text{ cm}^{-3}$, respectively.

According to the results of the accuracy (Fig. 5) and reliability (Fig. 6) of PTFs 5, 9 and 12, it seems that certain points of the SWRC (e.g., θ_{FC}) can help to improve the prediction of the SWRC and this is in agreement with Schaap et al. (2001). These results indicate that the presence of at least one moisture point (e.g., θ_{FC}) can improve the prediction of the SWRC. In the first step, PTF5 with two moisture points ($\theta_{FC} + \theta_{PWP}$) and PTF4 with one moisture point (θ_{FC}) improved the prediction of the SWRC by 55, 48, 42% and 51, 44, 38% in terms of the $IRMSE$ criterion compared to the PTFs 1, 2 and 3, respectively, in the RF method in the training step. In the testing section of the second step, PTF9 with two moisture points ($\theta_{FC} + \theta_{PWP}$) and PTF8 with one moisture point (θ_{FC}) decreased the $IRMSE$ by 49, 44% and 44, 39% compared to PTFs 6 and 7, respectively, in the NLR method. The points above are also true for the RF-based PTF12 in the third step of the testing section. Many researchers successfully applied θ_{FC} and θ_{PWP} as effective predictors to derive point- and parametric-PTFs (Børgesen and Schaap, 2005; Nemes et al., 2003; Schaap et al., 2001; Touil et al., 2016; Twarakavi et al., 2009).

3.3.2.2. Effect of using OM and K_s as predictors on the SWRC prediction. To evaluate the effect of using OM and/or K_s and points of the SWRC on the prediction of the SWRC, the performances of PTFs 13, 14 and 15 were compared with those of PTFs 4 and 5. The accuracy and reliability of the prediction of the SWRC by both NLR and RF methods, significantly ($P < 0.05$) decreased in terms of the $IRMSE$, for the PTFs 13, 14 and 15 from the fourth step, when OM and/or K_s were used with textural contents and BD as inputs instead of θ_{FC} or both θ_{FC} and θ_{PWP} in the list of input variables, compared to PTFs 4 and 5 at the first step (Fig. 5B and 6B). Therefore OM and K_s were not as effective predictors as θ_{FC} and θ_{PWP} in the prediction of the SWRC, because θ_{FC} and θ_{PWP} are two points of the SWRC and enter direct information of the SWRC into the PTFs, whereas OM and K_s enter indirect information, and therefore

had less effect in the improvement of the estimation of the SWRC. These results agreed well with results obtained by Børgeesen and Schaap (2005). They reported that PTFs with the inputs of θ_{FC} and θ_{PWP} had smaller RMSE values than a PTF with the input of OM (0.038 versus 0.042) in the prediction of the SWRC. On the other hand, the results showed that by adding OM and/or K_s as predictors in the PTFs 13, 14 and 15, the accuracy (Fig. 5B) and reliability (Fig. 6B) of the prediction of the SWRC improved by 16, 13, 17 and 7.1, 6.3, 6.9%, respectively, compared to the PTF3 in terms of the *IRMSE* criterion in the RF method.

The SWRC depends mainly on the soil texture and structure (Hillel, 1998), with OM affecting the SWRC through development of soil structure (Nemes et al., 2005), important at low suctions. However, the OM retains water itself. Similarly, K_s can be a descriptive index of soil texture and porosity (Hillel, 1998). The correlation results showed that K_s can be strongly influenced by clay content and textural statistics (d_g and δ_g) (Fig. 4). Bayat et al. (2013b) applied OM and K_s to estimate water content at the measured matric suctions. They found that the OM and K_s can be most appropriately used in point-based PTFs to estimate water content at the matric suctions of 25 and 50 kPa. Also, the result of the present study agreed well with results obtained by Hollis et al. (1977) and Rawls et al. (1983). In this study, the OM and K_s in the PTFs 13, 14 and 15 were not effective predictors compared to θ_{FC} and θ_{PWP} in the PTFs 4 and 5, otherwise they had better results than PTF3.

4. Conclusion

Machine-learning tools have been widely applied for the prediction of the SWRC. The present study evaluated the capability and performance of the RF method as a novel machine learning tool and compared its performance with that of the NLR method on the prediction of the SWRC, using different combinations of easily-available soil properties. It was found that the RF method had a better performance ($P < 0.05$) than the NLR method in the training step of the prediction of the SWRC in term of the *IRMSE*, *AIC* and R^2 criteria. However, in the testing step, NLR had a better performance than RF. The poor performance of the RF compared to the NLR method could be due to overprediction in the former, resulting in instability in the testing step. The RF method can be sensitive to sparse areas on the prediction space. In other words, the performance and sensitivity of predictions, and the computational intensity of the RF method depends on the distribution and number of observations and input variables. Therefore, the method should be tested further with different datasets to evaluate its performance through soil and water investigations. An improvement was achieved in the accuracy of the prediction of the SWRC in the training step of the PTF5 (with the inputs of Sand content + Clay content + BD + θ_{FC} + θ_{PWP}) by both NLR and RF methods and also an improvement was achieved in the reliability of the PTF9 (with the inputs of d_g + δ_g + BD + θ_{FC} + θ_{PWP}) and PTF12 (with the inputs of Sand content + Clay content + TP + θ_{FC} + θ_{PWP}) by the NLR and RF methods compared to other PTFs, respectively. Considering that the PTFs 5, 9, and 12 had no significant difference from PTF4 (with the inputs of Sand content + Clay content + BD + θ_{FC}) and PTF8 (with the inputs of d_g + δ_g + BD + θ_{FC} + θ_{PWP}), these latter PTFs, with less and more-easily measured input variables, are suggested to be the best PTFs for the prediction of the SWRC. Also, PTFs without predictors of θ_{FC} and θ_{PWP} , such as the PTF3 (with the inputs of Sand content + Clay content + BD) and PTF7 (with the inputs of d_g + δ_g + BD), can be effective models for the prediction of the SWRC.

CRedit authorship contribution statement

M. Rastgou: Data curation, Writing - original draft, Visualization, Investigation, Formal analysis. **H. Bayat:** Conceptualization, Methodology, Supervision, Project administration, Funding acquisition. **M. Mansoorizadeh:** Software, Validation. **Andrew S. Gregory:** Writing - review & editing.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

This work was funded by Bu-Ali Sina University, Hamedan, Iran. The authors are deeply grateful to anonymous reviewers and the editor for their helpful comments on the manuscript.

References

- Adhikary, P.P., Chakraborty, D., Kalra, N., Sachdev, C., Patra, A., Kumar, S., Tomar, R., Chandna, P., Raghav, D., Agrawal, K., 2008. Pedotransfer functions for predicting the hydraulic properties of Indian soils. *Soil Res.* 46, 476–484.
- Akaike, H., 1974. A new look at the statistical model identification. *IEEE Trans. Autom. Control* 19, 716–723.
- Araya, S.N., Ghezzehei, T.A., 2019. Using machine learning for prediction of saturated hydraulic conductivity and its sensitivity to soil structural perturbations. *Water Resour. Res.* 55, 5715–5737.
- Bayat, H., Ersahin, S., Hepper, E.N., 2013a. Improving estimation of specific surface area by artificial neural network ensembles using fractal and particle size distribution curve parameters as predictors. *Environ. Model. Assess.* 18, 605–614.
- Bayat, H., Neyshabouri, M., Mohammadi, K., Nariman-Zadeh, N., 2011. Estimating water retention with pedotransfer functions using multi-objective group method of data handling and ANNs. *Pedosphere* 21, 107–114.
- Bayat, H., Neyshabouri, M.R., Mohammadi, K., Nariman-Zadeh, N., Irannejad, M., 2013b. Improving water content estimations using penetration resistance and principal component analysis. *Soil Tillage Res.* 129, 83–92.
- Bayat, H., Sedaghat, A., Sinegani, A.A.S., Gregory, A.S., 2015. Investigating the relationship between unsaturated hydraulic conductivity curve and confined compression curve. *J. Hydrol.* 522, 353–368.
- Berry, W.D., 1993. *Understanding Regression Assumptions*. Sage Publications, London.
- Blake, G., Hartge, K., 1986. Bulk density, methods of soil analysis: Part 1. Physical and mineralogical methods, Madison, Wisconsin, USA. *Soil Sci. Soc. Am. J.*
- Børgeesen, C.D., Schaap, M.G., 2005. Point and parameter pedotransfer functions for water retention predictions for Danish soils. *Geoderma* 127, 154–167.
- Botula, Y.-D., Cornelis, W., Baert, G., Van Ranst, E., 2012. Evaluation of pedotransfer functions for predicting water retention of soils in Lower Congo (DR Congo). *Agric. Water Manag.* 111, 1–10.
- Botula, Y.-D., Cornelis, W.M., Baert, G., Mafuka, P., Van Ranst, E., 2013. Particle size distribution models for soils of the humid tropics. *J. Soils Sediments* 13, 686–698.
- Bouma, J., 1989. Using soil survey data for quantitative land evaluation. *Adv. Soil Sci.* Springer 177–213.
- Breiman, L., 1984. *Classification and Regression Trees*. Routledge, New York.
- Breiman, L., 2001. Random forests. *Machine Learn.* 45, 5–32.
- Bruce, R.R., Luxmoore, R.J., 1986. Water retention: field methods. In: Klute, A. (Ed.), *Methods of Soil Analysis: Part 1—Physical and Mineralogical Methods*. Soil Science Society of America, American Society of Agronomy, Madison, WI, pp. 663–686.
- Campbell, G.S., Horton Jr, R., 2002. *Methods of soil analysis: Part 4. Phys. Methods*. Soil Sci. Soc. Am.
- Chai, T., Draxler, R.R., 2014. Root mean square error (RMSE) or mean absolute error (MAE)?—Arguments against avoiding RMSE in the literature. *Geosci. Model. Dev.* 7, 1247–1250.
- Chakraborty, D., Mazumdar, S., Garg, R., Banerjee, S., Santra, P., Singh, R., Tomar, R., 2011. Pedotransfer functions for predicting points on the moisture retention curve of Indian soils. *Indian J. Agr. Sci.* 81, 1030.
- Cheng, L., Chen, X., De Vos, J., Lai, X., Witlox, F., 2019. Applying a random forest method approach to model travel mode choice behavior. *Travel Behav. Soc.* 14, 1–10.
- Dexter, A., Czyż, E., Richard, G., Reszkowska, A., 2008. A user-friendly water retention function that takes account of the textural and structural pore spaces in soil. *Geoderma* 143, 243–253.
- Dobaro, M.R., Cousin, I., Le Bas, C., Martin, M.P., 2019. Pedotransfer functions for predicting available water capacity in French soils, their applicability domain and associated uncertainty. *Geoderma* 336, 81–95.
- Efron, B., Tibshirani, R.J., 1994. *An Introduction To The Bootstrap*. CRC Press.
- Espósito, C., Barra, A., Evans, S.G., Scarascia Mugnozza, G., Delaney, K., 2014. Landslide susceptibility analysis by the comparison and integration of random forest and logistic regression methods; application to the disaster of Nova Friburgo-Rio de Janeiro, Brasil (January 2011). *EGU General Assembly Conference Abstracts*.
- Gee, G.W., Or, D., 2002. 2.4 Particle-Size Analysis. In: Dane, J.H., Topp, C.G. (Eds.), *Methods of Soil Analysis: Part 4 Physical Methods*. Soil Science Society of America, Madison, WI, pp. 255–293.
- Gunarathna, M., Sakai, K., Nakandakari, T., Momii, K., Kumari, M., 2019a. Machine learning approaches to develop pedotransfer functions for tropical Sri Lankan soils. *Water* 11, 1940.
- Gunarathna, M., Sakai, K., Nakandakari, T., Momii, K., Kumari, M., Amarasekara, M., 2019b. Pedotransfer functions to estimate hydraulic properties of tropical Sri Lankan soils. *Soil Till. Res.* 190, 109–119.

- Gupta, B., Rawat, A., Jain, A., Arora, A., Dhama, N., 2017. Analysis of various decision tree algorithms for classification in data mining. *Int. J. Comput. Appl.* 163, 15–19.
- Haghverdi, A., Leib, B.G., Cornelis, W.M., 2015. A simple nearest-neighbor technique to predict the soil water retention curve. *Trans. ASABE* 58, 697–705.
- Hillel, D., 1998. *Environmental Soil Physics: Fundamentals, Applications, and Environmental Considerations*. Academic Press.
- Hocking, R.R., 2013. *Methods and Applications of Linear Models: Regression And The Analysis Of Variance*. John Wiley & Sons.
- Hollis, J., Jones, R., Palmer, R., 1977. The effects of organic matter and particle size on the water-retention properties of some soils in the West Midlands of England. *Geoderma* 17, 225–238.
- Hong, H., Pourghasemi, H.R., Pourtaghi, Z.S., 2016. Landslide susceptibility assessment in Lianhua County (China): a comparison between a random forest data mining technique and bivariate and multivariate statistical models. *Geomorphology* 259, 105–118.
- IBM, C., 2016. *IBM SPSS Statistics for Windows, Version 24.0*. Armonk, NY: IBM Corp.
- Ibrahim, I.A., Khatib, T., 2017. A novel hybrid model for hourly global solar radiation prediction using random forests technique and firefly algorithm. *Energy Convers. Manag.* 138, 413–425.
- Janitza, S., Tutz, G., Boulesteix, A.-L., 2016. Random forest for ordinal responses: prediction and variable selection. *Comput. Statist. Data Anal.* 96, 57–73.
- Kholosi, M., Alhamdoosh, M., Douaik, A., Gabriels, D., Cornelis, W., 2016. Enhanced pedotransfer functions with support vector machines to predict water retention of calcareous soil. *Eur. J. Soil Sci.* 67, 276–284.
- Khodaverdilloo, H., Homae, M., van Genuchten, M.T., Dashtaki, S.G., 2011. Deriving and validating pedotransfer functions for some calcareous soils. *J. Hydrol.* 399, 93–99.
- Klute, A., 1986. Water Retention: Laboratory Methods. In: Klute, A. (Ed.), *Methods of Soil Analysis: Part 1—Physical and Mineralogical Methods*. Soil Science Society of America, American Society of Agronomy, Madison, WI, pp. 635–662.
- Klute, A., Dirksen, C., 1986. Hydraulic Conductivity and Diffusivity: Laboratory Methods. In: Klute, A. (Ed.), *Methods of Soil Analysis: Part 1—Physical and Mineralogical Methods*. Soil Science Society of America, American Society of Agronomy, Madison, WI, pp. 687–734.
- Koekkoek, E., Bootink, H., 1999. Neural network models to predict soil water retention. *Eur. J. Soil Sci.* 50, 489–495.
- Lamorski, K., Pachepsky, Y., Sławiński, C., Walczak, R., 2008. Using support vector machines to develop pedotransfer functions for water retention of soils in Poland. *Soil Sci. Soc. Am. J.* 72, 1243–1247.
- Lamorski, K., Sławiński, C., Moreno, F., Barna, G., Skierucha, W., Arrue, J.L., 2014. Modelling soil water retention using support vector machines with genetic algorithm optimisation. *Sci. World J.* 2014 (740521), 1–10.
- Liaw, A., Wiener, M., 2002. Classification and regression by random forest. *R News* 2, 18–22.
- Ließ, M., Glaser, B., Huwe, B., 2012. Uncertainty in the spatial prediction of soil texture: comparison of regression tree and Random Forest models. *Geoderma* 170, 70–79.
- Liu, Y., 2014. Random forest algorithm in big data environment. *Comput. Model. New Tech.* 18, 147–151.
- Ma, Y., Cukic, B., Singh, H., 2005. A Classification Approach To Multi-Biometric Score Fusion, International Conference on Audio-and Video-Based Biometric Person Authentication. Springer, pp. 484–493.
- MathWorks, MATLAB: the language of technical computing Inc 2018 Natick Massachusetts, United States.
- Matin, S., Chelgani, S.C., 2016. Estimation of coal gross calorific value based on various analyses by random forest method. *Fuel* 177, 274–278.
- Medrado, E., Lima, J.E., 2014. Development of pedotransfer functions for estimating water retention curve for tropical soils of the Brazilian savanna. *Geoderma Regional* 1, 59–66.
- Merdun, H., Çınar, Ö., Meral, R., Apan, M., 2006. Comparison of artificial neural network and regression pedotransfer functions for prediction of soil water retention and saturated hydraulic conductivity. *Soil Tillage Res.* 90, 108–116.
- Minasny, B., McBratney, A.B., Bristow, K.L., 1999. Comparison of different approaches to the development of pedotransfer functions for water-retention curves. *Geoderma* 93, 225–253.
- Mualem, Y., 1976. A new model for predicting the hydraulic conductivity of unsaturated porous media. *Water Resour. Res.* 12, 513–522.
- Nemes, A., Rawls, W.J., Pachepsky, Y.A., 2005. Influence of organic matter on the estimation of saturated hydraulic conductivity. *Soil Sci. Soc. Am. J.* 69, 1330–1337.
- Nemes, A., Rawls, W.J., Pachepsky, Y.A., 2006. Use of the nonparametric nearest neighbor approach to estimate soil hydraulic properties. *Soil Sci. Soc. Am. J.* 70, 327–336.
- Nemes, A., Schaap, M., Wösten, J., 2003. Functional evaluation of pedotransfer functions derived from different scales of data collection. *Soil Sci. Soc. Am. J.* 67, 1093–1102.
- Neyshaburi, M.R., Bayat, H., Mohammadi, K., Nariman-Zadeh, N., Irannejad, M., 2015. Improvement in estimation of soil water retention using fractal parameters and multiobjective group method of data handling. *Arch. Agron. Soil Sci.* 61, 257–273.
- Nguyen, P.M., Haghverdi, A., De Pue, J., Botula, Y.-D., Le, K.V., Waegeman, W., Cornelis, W.M., 2017. Comparison of statistical regression and data-mining techniques in estimating soil water retention of tropical delta soils. *Biosyst. Eng.* 153, 12–27.
- Pachepsky, Y., Rawls, W., Gimenez, D., Watt, J., 1998. Use of soil penetration resistance and group method of data handling to improve soil water retention estimates. *Soil Tillage Res.* 49, 117–126.
- Pachepsky, Y.A., Rawls, W., 1999. Accuracy and reliability of pedotransfer functions as affected by grouping soils. *Soil Sci. Soc. Am. J.* 63, 1748–1757.
- Pachepsky, Y.A., Rawls, W., Lin, H., 2006. *Hydropedology and pedotransfer functions*. *Geoderma* 131, 308–316.
- Pachepsky, Y.A., Timlin, D., Vallyay, G., 1996. Artificial neural networks to estimate soil water retention from easily measurable data. *Soil Sci. Soc. Am. J.* 60, 727–733.
- Rab, M., Chandra, S., Fisher, P., Robinson, N., Kitching, M., Aumann, C., Imhof, M., 2011. Modelling and prediction of soil water contents at field capacity and permanent wilting point of dryland cropping soils. *Soil Res.* 49, 389–407.
- Rajkai, K., Kabos, S., Van Genuchten, M.T., 2004. Estimating the water retention curve from soil properties: comparison of linear, nonlinear and concomitant variable methods. *Soil Tillage Res.* 79, 145–152.
- Rawls, W., Brakensiek, D., Soni, B., 1983. Agricultural management effects on soil water processes part I: Soil water retention and Green and Ampt infiltration parameters. *Trans. ASAE* 26, 1747–1752.
- Rawls, W., Gish, T., Brakensiek, D., 1991. Estimating soil water retention from soil physical properties and characteristics. *Adv. Soil Sci.* Springer 213–234.
- Rawls, W.J., Brakensiek, D., 1985. Prediction of soil water properties for hydrologic modeling. *Watershed management in the eighties*. ASCE 293–299.
- Schaap, M.G., Leij, F.J., van Genuchten, M.T., 1998. Neural network analysis for hierarchical prediction of soil hydraulic properties. *Soil Sci. Soc. Am. J.* 62, 847–855.
- Schaap, M.G., Leij, F.J., van Genuchten, M.T., 2001. Rosetta: A computer program for estimating soil hydraulic parameters with hierarchical pedotransfer functions. *J. Hydrol.* 251, 163–176.
- Scheinost, A., Sinowski, W., Auerswald, K., 1997. Regionalization of soil water retention curves in a highly variable soilscape, I. Developing a new pedotransfer function. *Geoderma* 78, 129–143.
- Seo, S., 2006. A review and comparison of methods for detecting outliers in univariate data sets, Thesis for Master of Science in Field of Public Health University of Pittsburgh, pp. 1–59.
- Shirazi, M.A., Boersma, L., 1984. A unifying quantitative analysis of soil texture. *Soil Sci. Soc. Am. J.* 48, 142–147.
- Shwetha, P., Varjia, K., 2015. Soil water retention curve from saturated hydraulic conductivity for sandy loam and loamy sand textured soils. *Aquat. Procedia* 4, 1142–1149.
- Simmons, J.P., Nelson, L.D., Simonsohn, U., 2011. False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychol. Sci.* 22, 1359–1366.
- Szabó, B., Szatmári, G., Takács, K., Laborczi, A., Makó, A., Rajkai, K., Pásztor, L., 2019. Mapping soil hydraulic properties using random forest based pedotransfer functions and geostatistics. *Hydrol. Earth Syst. Sci.* 23, 2615–2635.
- Tietje, O., Tapkenhinrichs, M., 1993. Evaluation of pedo-transfer functions. *Soil Sci. Soc. Am. J.* 57, 1088–1095.
- Tomasella, J., Hodnett, M.G., 1998. Estimating soil water retention characteristics from limited data in Brazilian Amazonia. *Soil Sci.* 163, 190–202.
- Tomasella, J., Hodnett, M.G., Rossato, L., 2000. Pedotransfer functions for the estimation of soil water retention in Brazilian soils. *Soil Sci. Soc. Am. J.* 64, 327–338.
- Tóth, B., Makó, A., Toth, G., 2014. Role of soil properties in water retention characteristics of main Hungarian soil types. *J. Cent. Eur. Agric.* 15, 137–153.
- Touil, S., Degre, A., Chabaca, M.N., 2016. Sensitivity analysis of point and parametric pedotransfer functions for estimating water retention of soils in Algeria. *Soil* 2, 647.
- Twarakavi, N.K., Šimůnek, J., Schaap, M., 2009. Development of pedotransfer functions for estimation of soil hydraulic parameters using support vector machines. *Soil Sci. Soc. Am. J.* 73, 1443–1452.
- Ungaro, F., Calzolari, C., Busoni, E., 2005. Development of pedotransfer functions using a group method of data handling for the soil of the Pianura Padana-Veneta region of North Italy: water retention properties. *Geoderma* 124, 293–317.
- van Genuchten, M.T., 1980. A closed-form equation for predicting the hydraulic conductivity of unsaturated soils. *Soil Sci. Soc. Am. J.* 44, 892–898.
- Verhagen, J., 1997. Site specific fertilizer application for potato production and effects on N-leaching using dynamic simulation modelling. *Agric. Ecosyst. Environ.* 66, 165–175.
- Verikas, A., Gelzinis, A., Bacauskiene, M., 2011. Mining data with random forests: A survey and results of new tests. *Pattern Recogn.* 44, 330–349.
- Walkley, A., Black, I.A., 1934. An examination of the Degtjareff method for determining soil organic matter, and a proposed modification of the chromic acid titration method. *Soil Sci.* 37, 29–38.
- Wassar, F., Gandolfi, C., Rienzi, M., Chiaradia, E.A., Bernardoni, E., 2016. Predicted and measured soil retention curve parameters in Lombardy region north of Italy. *Int. Soil Water Conserv. Res.* 4, 207–214.
- Wiesmeier, M., Barthold, F., Blank, B., Kögel-Knabner, I., 2011. Digital mapping of soil organic matter stocks using Random Forest modeling in a semi-arid steppe ecosystem. *Plant Soil* 340, 7–24.
- Williams, J., Ross, P., Bristow, K.L., 1992. Prediction of the Campbell water retention function from texture, structure, and organic matter. *Indirect Methods For Estimating The Hydraulic Properties Of Unsaturated Soils*. University of California, Riverside.
- Wösten, J., Pachepsky, Y.A., Rawls, W., 2001. Pedotransfer functions: bridging the gap between available basic soil data and missing soil hydraulic characteristics. *J. Hydrol.* 251, 123–150.
- Zaklouta, F., Stanculescu, B., 2012. Real-time traffic-sign recognition using tree classifiers. *IEEE Trans. Intell. Transp. Syst.* 13, 1507–1514.
- Zhao, P., Su, X., Ge, T., Fan, J., 2016. Propensity score and proximity matching using random forest. *Contemp. Clin. Trials* 47, 85–92.