

Responses to the Referee

Major comments

A review of “Systematic comparison of five machine-learning methods” by Mo Zhang and Wenjiao Shi. The manuscripts describe a comparison between five machine learning methods for soil classification and interpolation of soil particle size fractions. It explores different transformed data. There are a few major problems with the manuscript:

Comment 1: The five machine-learning methods come falling from the air. That is to be regretted, as there is much more information available in the literature. The manuscript now provides only some technical aspects, and for instance not prior assumptions, restrictions in their use or their general comparability. The manuscript would largely benefit from a short mathematical introduction to the five techniques, from where it would become clear whether in fact comparable methods are compared, or that there is a comparison between apples and oranges.

Response: Thanks for the referee’s suggestion for the machine-learning methods. We have added mathematical introduction of these five machine-learning methods to show these techniques can be compared.

P19L4, Section 2.4.1—2.4.5 in our revised version:

(1) K-nearest neighbor (KNN)

K-nearest neighbor (KNN) is a simple non-parametric classifier based on known instance to label unknown instance (Cover and Hart, 1967). For the test set, K-nearest training set vectors were found, and maximum summed kernel densities were computed for classification. Moreover, continuous variables can also be predicted for regression with the average values of K-nearest neighbors. For a training set of observed data $L = \{(y_i, x_i), i = 1, \dots, n_L\}$, class $y_i \in \{1, \dots, c\}$, and the predictor values $x'_i = (x_{i1}, \dots, x_{ip})$. For a new observation (y, x) , the nearest neighbor $(y_{(1)}, x_{(1)})$ is based on the distance function which is as follows:

$$d(x, x_{(1)}) = \min_i (d(x, x_i)), \quad (1)$$

and $\hat{y} = y_{(1)}$ refers to the nearest neighbor, which is the prediction for y . Value $x_{(j)}$ and $y_{(j)}$ is the j th nearest neighbor of x and class of training set, respectively. Weighted KNN is an extended version of KNN, which considers the maximum of summed kernel densities and the K nearest vectors of training set for each row of the test set (the distances of the nearest neighbors) based on the Minkowski distance, more details can be found in Hechenbichler and Schliep (2004), the equation for Minkowski distance is as follows:

$$d(x_i, x_j) = (\sum_{s=1}^p |x_{is} - x_{js}|^q)^{1/q}, \quad (2)$$

where $d(x_i, x_j)$ refers to the Euclidean distance when $q = 2$ and the absolute distance results for $q = 1$. Therefore, the parameters of KNN contain the maximum value of k (kmax), the distances of the nearest neighbors (distance) and the types of kernel function (kernel). The KNN model is available in the R package “knn” (Schliep and Hechenbichler, 2016).

5 (2) Multilayer perceptron neural network (MLP)

Multilayer perceptron neural network (MLP), which is currently one of the most popular multilayer feed forward backpropagation networks, was selected to train artificial neural network (ANN) models in our study due to its rapid operation, small set of training requirements and ease of implementation (Subasi, 2007). MLP neurons can perform classification or regression depending on whether the response variable is categorical or continuous. The MLP has three sequential layers: input layer, hidden layer and output layer. In the hidden layer of MLP, each neuron j sums input environmental covariate in our study x_i after multiplying them by the connection weights w_{ji} respectively, and calculates its output y_j (soil PSF component or texture class) as a function of the sum:

$$y_j = f(\sum w_{ji}x_i), \quad (3)$$

where f is the activation function, which can be a linear (selected in our study) or logistic function. The sum of squared differences between the predicted values and observed values of the output results of neurons E is defined as follows:

$$E = \frac{1}{2} \sum_j (y_{pj} - y_{oj})^2, \quad (4)$$

where y_{pj} and y_{oj} is the predicted and observed value of output neuron j , respectively. Each w_{ji} is adjusted to reduce E and the adjustment of w_{ji} depends on the training algorithm (Basheer and Hajmeer, 2000). The resilient backpropagation algorithm was chosen because the learning rate of this algorithm is adaptive, avoiding oscillations and accelerating the learning process (Behrens and Scholten, 2006). The range of the data set should be standardized because MLPs operate in terms of the scale 0 to 1. MLP can be run using the R package “RSNNS” (Bergmeir and Benitez, 2012).

(3) Random forest (RF)

Random forest (RF) was developed by Breiman (2001), combining the bagging method (Breiman, 1996) with the random variable selection, and the principle was to merge a group of “weak learners” together to form a “strong learner”. Bootstrap sampling is used for each tree of RF, and the rules to binary split data are different for regression and classification problems. For classification, the Gini index is used to split the data; for regression, minimizing the sum of the squares of the mean deviations (M) can be selected to train each tree model, the equations are as follows:

$$Gini = 1 - \sum_{k=1}^K p_k^2, \quad (5)$$

$$Gini(D, A) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2), \quad (6)$$

$$M = \min_A [\min_{c_1} \sum_{x_i \in D_1(A)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in D_2(A)} (y_i - c_2)^2], \quad (7)$$

where p_k refers to the proportion of k th class in the data set on the current node, for feature $A = a$, data set D is divided into two parts (D_1 and D_2), D_1 describes the data set which meets the condition $A = a$ and D_2 is the opposite of D_1 ; $Gini(D, A)$ represents the uncertainty of set D after binary split; y_i is the predicted value of input value x_i , c_1 and c_2 is the mean of data set D_1 and D_2 , respectively. Benefits of using RFs are that the ensembles of trees are used without pruning.

5 In addition, RF is relatively robust to overfitting, and standardization or normalization are not necessary because it is insensitive to the range of value. Two parameters should be adjusted for RF model: the number of trees (ntree) and the number of features randomly sampled at each split (mtry). The RF model is available in the R package “randomForest” (Liaw and Wiener, 2002).

(4) Support vector machines (SVM)

10 The support vector machine (SVM), proposed by Cortes and Vapnik (1995), is a type of generalized linear classifier that is widely applied for classification and regression problems in soil science (Burges, 1998). The main principle of SVM is to classify different classes by constructing an optimal separating hyperplane in the feature space (so called “structural risk minimization”). Regression problems also can be solved by minimization of the structural risk using loss functions (Vapnik, 1998) in SVM, named support vector regression. For a data set $\{x_i, y_i\}$, $i = 1, \dots, k$, $x \in R$ and x refers to an n-

15 dimensional vector, $y \in \{-1, +1\}$ is the class corresponding to x , the equation for calculating a hyperplane of SVM is defined as follows:

$$\min_{w, b, \xi} \frac{1}{2} w^T \times w + C \sum_{i=1}^k \xi_i, \quad \text{s.t. } y_i(w^T \times \phi(x_i) + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, \dots, k, \quad (8)$$

where $\phi(x_i)$ refers to the mapping from the input space to the feature space, $C > 0$ is penalty factor (cost), w , b , and ξ

20 are the parameters need to be optimized during the process of model training, which can be determined by the Lagrange multipliers:

$$f(x) = \text{sgn}(y_i a_i k(x_i, x) + b^*) \quad (9)$$

where a_i refers to the support vector, $k(x_i, x)$ refers to the kernel function, and b^* is the bias. The advantages of SVMs are that they are effective in high dimensional spaces. Radial basis function was selected for SVM as the kernel function in our

25 study, and two other parameters need to be tuned, i.e., cost and gamma, controlling the tradeoff between the classification accuracy and complexity, and the ranges of radial effect, respectively. The SVM model is available in the R package “e1071” (Meyer et al., 2017).

(5) Extreme gradient boosting (XGB)

30 Extreme Gradient Boosting, put forward by Chen and Guestrin (2016), is an efficient method of implementation for gradient boosting frames, tree learning algorithms and efficient linear model solvers to solve both classification and regression problems (Chen et al., 2018). Like the boosted regression trees (Elith et al., 2008), it follows the principle of gradient enhancement;

however, more regularized model formalization is applied to XGB to control over-fitting, making it more remarkable. In addition, parallel calculations can be automatically executed during the training phase of the XGB model, presenting a great advantage in large data sets, as the XGB can be more than ten times faster than the existing gradient boosting model (Chen and Guestrin, 2016). The residuals of the first tree can be fitted by the second tree to enhance the model accuracy and the sum of the prediction of each tree generates the ultimate prediction. The general prediction function at step t is defined as follows:

$$f_i^{(t)} = \sum_{k=1}^t f_k(x_i) = f_i^{(t-1)} + f_t(x_i), \quad (10)$$

where $f_t(x_i)$ refers to the tree (learner) at step t , $f_i^{(t)}$ and $f_i^{(t-1)}$ refer to the predicted values at steps t and $t-1$, and x_i is the input value.

$$Obj^{(t)} = \sum_{k=1}^n l(\overline{y}_i, y_i) + \sum_{k=1}^n \Omega(f_i), \quad (11)$$

where l refers to the loss function, n is the number of data set, and Ω refers to the regularization term, which equation is defined as follows:

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \|\omega\|^2, \quad (12)$$

where ω refers to the score vector, λ is the parameter of regularization term, and γ is the minimum loss. There are seven parameters should be tuned in XGB, containing the learning rate (eta), the maximum depth of a tree (max_depth), the maximum number of boosting iterations (nrounds), the subsample ratio of columns (colsample_bytree), the subsample ratio of the training instance (subsample), the minimum loss reduction (gamma) and the minimum sum of instance weight (min_child_weight). The XGB model is available in the R package “xgboost” (Chen et al., 2018).

New Reference

- Basheer, I. A., and Hajmeer, M.: Artificial neural networks: fundamentals, computing, design, and application, Journal of Microbiological Methods, 43, 3-31, [https://doi.org/10.1016/S0167-7012\(00\)00201-3](https://doi.org/10.1016/S0167-7012(00)00201-3), 2000.
- Hechenbichler, K., and Schliep, K.: Weighted k-Nearest-Neighbor Techniques and Ordinal Classification, 2004.

Comment 2: Section 2.5 discusses log-ratio transformation methods.

- a. For one, here the term ‘methods’ is used (which I like) and at other places the term ‘approaches’ (which should be avoided throughout).
- b. But more importantly is the lack of mathematical rigor throughout. Line 13 gives a condition $\forall j = 1, \dots, j-1, j+1, \dots, D$. I do not at all understand this restriction. I think it is wrong.
- c. In equation (4) it is unclear why on the left hand side there is a term y_j mentioned: is that exceptional? I do not think so: it has to be removed.
- d. In equation (9) there seems to be the D -ith root: is that correct? What is the rationale behind this?
- e. Again in equation (9), the z is only defined for all except for the last component. Why is that the case?

In all, this section needs a very careful revision by a professional mathematician.

Response: Thanks for the referee's suggestion for log ratio transformation methods. We have improved this section in our revised version and the section was checked by a professional mathematician.

P22L26, 2.5 Log-ratio transformation methods:

5 **For comment 2 a**, we have replaced “approaches” with “methods”.

For comment 2 b, we apologize for our mistake about conditions and equations such as $\forall j = 1, \dots, j-1, j+1, \dots, D$ were **wrong**, and the right forms were:

For the composition of D elements $\mathbf{x} = [x_1, \dots, x_D]$, $x_j > 0$, $\forall j = 1, 2, \dots, D$, and $\sum_{j=1}^D x_j = 1$, the transformation equation for ALR is defined as follows:

10
$$alr(\mathbf{x}) = (\ln \frac{x_1}{x_j}, \dots, \ln \frac{x_{j-1}}{x_j}, \ln \frac{x_{j+1}}{x_j}, \dots, \ln \frac{x_D}{x_j}), \quad (13)$$

For comment 2 c, we have removed “ y_j ” in Eq. (4) in our revised version.

$$clr(\mathbf{x}) = (\ln \frac{x_1}{\sqrt[D]{\prod_{j=1}^D x_j}}, \dots, \ln \frac{x_D}{\sqrt[D]{\prod_{j=1}^D x_j}}), \quad (14)$$

For comment 2 d and e, the questions of Eq. (9) $z_i = \sqrt{\frac{D-i}{D-i+1}} \ln \frac{x_i}{\sqrt[D-i]{\prod_{j=i+1}^D x_j}}$, for $i = 1, \dots, D-1$ can be explained by the

rationale and interpretation of ILR transformation method. The isometric log ratio transformation method, proposed by
 15 Egozcue et al. (2003), in ILR transformation method, an orthonormal basis was chosen to project the compositions from S^D (simplex with respect to the Aitchison geometry) to R^{D-1} (real space with respect to the Euclidean geometry) isometrically. The choice of a specific orthonormal basis on S^D is important for the interpretation of coordinates (Fiserova and Hron, 2011), which can be explained by sequential binary partition (SBP) because compositions can be interpreted in terms of their groups (Egozcue and Pawłowsky-Glahn, 2005). In the SBP, the choice of construction of coordinates (so-called balances) is:

20 **(1)** First, the parts of the composition are divided into two groups: group coded by +1 and group coded by -1, and the first coordinate is obtained to describe the balance between the parts of +1 and -1 groups.

(2) Second and following steps, previous +1 and -1 groups are divided into two new groups, respectively, coding by +1 and -1 similarly until the components not involved are coded with 0. The balance of each step remains the same as before and the total number of steps is $D-1$ (the dimension of S^D , see Fiserova and Hron, 2011), finally. Therefore, the equation for coordinates
 25 in the k th step is as follows:

$$z_k = \sqrt{\frac{r_k s_k}{r_k + s_k}} \ln \left(\frac{(x_{i_1} x_{i_2} \dots x_{i_{r_k}})^{1/r_k}}{(x_{j_1} x_{j_2} \dots x_{j_{s_k}})^{1/s_k}} \right), \quad k = 1, \dots, D-1, \quad (15)$$

where z_k refers to the balance between two groups, i_1, i_2, \dots, i_{r_k} is the r_k parts of the +1 group, and j_1, j_2, \dots, j_{s_k} is the s_k parts of the -1 group. The balances therefore contain stepwise all the relevant information of compositions in +1 and -1 groups. It also can be explained in table, for example, in soil PSF compositional data ($D = 3$), a choice of SBP is shown.

Table 1 One choice of SBP of soil PSF data (D = 3).

order	sand	silt	clay	r	s	balance
1	+	-	-	1	2	$z_1 = \sqrt{\frac{2}{3}} \ln \frac{\text{sand}}{\sqrt{\text{silt} \times \text{clay}}}$
2	0	+	-	1	1	$z_2 = \sqrt{\frac{1}{2}} \ln \frac{\text{silt}}{\text{clay}}$

Thus, the interpretation above can explain these questions: D-ith root—derive from the Eq. (15) and Table 1, and the number of equations for ILR method is D-1, containing all information of compositions.

5 Note that the SBP can be applied blindly or can be based on expert knowledge (Fiserova and Hron, 2011), the SBP chose in our manuscript is shown in Table 1; however, there are different results such as the accuracy evaluation, maps of spatial prediction when different SBPs are used, which has been contained in our current research.

Reference

10 Egozcue, J. J., Pawlowsky-Glahn, V., Mateu-Figueras, G., and Barcelo-Vidal, C.: Isometric logratio transformations for compositional data analysis, *Math. Geol.*, 35, 279-300, 10.1023/a:1023818214614, 2003.

Egozcue, J. J., and Pawlowsky-Glahn, V.: Groups of parts and their balances in compositional data analysis, *Math. Geol.*, 37, 795-828, 10.1007/s11004-005-7381-9, 2005.

Fiserova, E., and Hron, K.: On the Interpretation of Orthonormal Coordinates for Compositional Data, *Math Geosci.*, 43, 455-15 468, 10.1007/s11004-011-9333-x, 2011.

Comment 3: The results section is not at all convincing in its current status. I could understand that a selection is made for the best of machine-learning method x transformation combination for the particular study area. That requires some space, but it then should be followed up by only one outcome, namely the best. As a scientist I am not interested in maps of an inferior quality. Hence, in figure 6, four of the five maps are redundant. The authors could use small sections of the maps, though, to show where the other techniques are critically sub-optimal, but not more than that.

20

Response: Thanks for the referee’s suggestion for the predicted classification maps. Because the difference of these maps covers the full study area, a certain small section of the maps cannot be selected. Further, the abundance index was used to describe how many soil texture classes were predicted for each machine-learning method in the study area. Therefore, it is necessary to compare the predicted classification maps of the whole study area. A number of literatures of spatial prediction comparing prediction maps using all the study areas are listed here (Buchanan et al., 2012; Niang et al., 2014; Wang and Shi, 2018). Moreover, even if more detailed information is produced on a given map, it is hard to argue which map is the best

25

among these five maps using different machine-learning methods, especially for regions where there is no data (cannot be verified). In the revised manuscript, we focused on objectively revealing these results generated from different methods, including the indicators of abundance indices, the value ranges of prediction maps, distribution characteristics and textural features, rather than subjectively describing which method is better than the others.

5

Reference

Buchanan, S., Triantafyllis, J., Odeh, I. O. A., and Subansinghe, R.: Digital soil mapping of compositional particle-size fractions using proximal and remotely sensed ancillary data, *Geophysics*, 77, WB201-WB211, 10.1190/geo2012-0053.1, 2012.

10 Niang, M. A., Nolin, M. C., Jegou, G., and Perron, I.: Digital Mapping of Soil Texture Using RADARSAT-2 Polarimetric Synthetic Aperture Radar Data, *Soil Sci. Soc. Am. J.*, 78, 673-684, 10.2136/sssaj2013.07.0307, 2014.

Wang, Z., and Shi, W. J.: Robust variogram estimation combined with isometric log-ratio transformation for improved accuracy of soil particle-size fraction mapping, *Geoderma*, 324, 56-66, 10.1016/j.geoderma.2018.03.007, 2018.

15 **Comment 4: Still in the results section, I am not convinced about the repetition on page 20 of what is already in the table. I want to know why one method and one transformation is the best, and in particular whether that is due to the specific case study or to an inherently superior combination of the transformation with the methods. What stands out is that RF_ORI is the best. Fine, point taken. But then in figure 7, I am only interested in the RF_ORI map, and all the other maps should be avoided. Similar remarks apply to table 3 and figures 8 and 9. The authors should make more**
20 **out of the data that they had to their disposal, than just creating sub-optimal maps! Figure 10, by the way, is interesting and may lead to a sub-optimal map (in terms of the RMSE). XGB is not so bad in terms of RMSE and much faster, hence if speed is an issue (when would that be?) then a researcher may opt for XGB. It should then be clear what s/he essentially loses in terms of representations on maps.**

25 **Response:** Thanks for the referee’s suggestions and questions. Table 4 (the comparisons of accuracy of different machine-learning models combined with original and transformed data) demonstrated the RF performed better than other four machine-learning methods. RF_ORI had the best performance of the MAE, RMSE and R^2 of sand, silt and clay components, and RF_ILR performed best in the STRESS. In fact, log ratio transformation methods can improve the overall evaluation indicators of compositions such as the STRESS. It is not because of the specific case study, the same conclusion was reached in our
30 previous study (Wang and Shi, 2018). Moreover, because of the large soil sampling data sets and regional scale study area, there are outliers (skewed distribution) and uncertainty in soil samples and spatial prediction. RF method has advantages than other machine-learning methods as we mentioned in our manuscript: “RF is relatively robust to overfitting, and standardization

or normalization are not necessary because it is insensitive to the range of value”. Therefore, this can be taken as a general conclusion.

In addition, two sections (the performance of these maps and the table of evaluation indicators) are independent, which can draw different conclusions respectively. For the section of spatial prediction maps of soil PSF, we focused on the systematic comparison of different performance generated from different methods rather than only creating a certain map which performed best in the accuracy evaluation, similar structure layout of spatial prediction research also can be shown in Zhang et al. (2014) and the reference cited in the response of comment 3. This is because more information should be taken into account with respect to the systematic comparison. We have added more detailed description between different prediction maps to evaluate which methods are better. All the ranges of prediction maps of sand (approximately 9.0—90.0 %) were within the range of the original data (0.98—99.66 %), like the distribution of USDA triangles. Therefore, all models overestimated the low values and underestimated the high values. For sand content, maps of RF_ILR (7.9—94.7 %) and XGB_ORI (1.8—92.4 %) generated wider output distributions and were relatively closer to the range of the distribution of original data than other prediction maps such as KNN_ILR (7.3—88.6 %), KNN_ORI (7.8—80.8 %), MLP_ILR (8.8—90.8 %), MLP_ORI (9.0—90.3 %), RF_ORI (9.0—81.0 %), SVM_ILR (6.5—85.6 %), SVM_ORI (7.3—90.0 %) and XGB_ILR (5.0—88.5 %). From this point of view, RF and XGB were more accurate than others. Furthermore, for the soil texture classification (Fig. 8 and 9), LoSa and SaLo are obviously most confused classes. But they are fairly similar to each other so not a big problem probably. We have added the description of the similarity of LoSa and SaLo in discussion part to make readers have more profound impression.

With respect to the total computing time of different machine-learning methods, XGB did not lose much on accuracy (in terms of the RMSE etc.) but the computing time drops significantly. We have added this to our discussion.

P38L1: “All the ranges of prediction maps of sand (approximately 9.0—90.0 %) were within the range of the original data (0.98—99.66 %). RF_ILR (7.9—94.7 %) and XGB_ORI (1.8—92.4 %) generated wider output distributions and were relatively closer to the range of the distribution of original data than other prediction maps such as KNN_ILR (7.3—88.6 %), KNN_ORI (7.8—80.8 %), MLP_ILR (8.8—90.8 %), MLP_ORI (9.0—90.3 %), RF_ORI (9.0—81.0 %), SVM_ILR (6.5—85.6 %), SVM_ORI (7.3—90.0 %) and XGB_ILR (5.0—88.5 %).”

P46L6: “In fact, LoSa and SaLo were obviously most confused classes; however, they are fairly similar to each other (see Fig. 8).”

P44L25: “For the total computing time, RF revealed the longest time with respect to both classification (453.73 s) and regression (188.87 s); however, it is the most accurate among five machine-learning methods in our study. In addition, for trade-offs of the total computing time of model and sub-optimal accuracy, XGB was superior to any other model, reducing the computing time significantly, while maintaining the accuracy not drop too much.”

Reference

Wang, Z., and Shi, W. J.: Robust variogram estimation combined with isometric log-ratio transformation for improved accuracy of soil particle-size fraction mapping, *Geoderma*, 324, 56-66, 10.1016/j.geoderma.2018.03.007, 2018.

Zhang, S. W., Kong, W. F., Huang, Y. F., Shen, C. Y., and Ye, H. C.: Spatial Prediction of Topsoil Texture in a Mountain-plain Transition Zone Using Univariate and Multivariate Methods Based on Symmetry Logratio Transformation, *Intell. Autom. Soft Comput.*, 20, 115-129, 10.1080/10798587.2013.861966, 2014.

Comment 5: The discussion section has some interesting elements, but I could easily imagine an improvement when better focusing upon what has been achieved and how it should be interpreted, also in terms of the soils and the particle sizes. In particular, a transition of the methods to other areas should be discussed: how should we take it? Maybe do it hierarchically, i.e. first a quick and imprecise method, followed by a precise method?

Response: Thanks for the referee's suggestion for the discussion of a transition of the methods to other areas. We have added it to the discussion section of our revised revision.

P44L29: *"With respect to the generality results of a transition of these machine-learning methods to other areas, it can be considered hierarchically. First, for the quick and imprecise machine-learning methods, XGB was recommended (sub-optimal accuracy), which was fast at the expense of a loss in precision. Second, considering the precise methods, RF can deliver the most accurate results, but it takes the longest computing time. Therefore, XGB should be selected when researchers deal with larger data sets and regional scale study area; if they have enough time while want to produce more accurate results, RF is recommended."*

Comment 6: Finally, the abstract should be improved: the opening statements are too wordy, as a single sentence justification for the study is enough. The term 'notable consequences' is too vague for an abstract. The final main concluding sentence '... helps to elucidate the processing and selection of compositional data in spatial simulation' is not justified from the manuscript. We only see that one method x transformation combination is doing best, another combination is fast at the expense of a loss in precision. That seems to be a good conclusion, and in that sense the study is valuable as an interesting case study on soil analysis in a rather large area.

Response: Thanks for the referee's suggestion for the abstract which should be improved, and we have improved this in our revised version, such as wordy opening statements, vague description (*"notable consequences"*), subjective sentence (*"Our systematic comparison helps to elucidate the processing and selection of compositional data in spatial simulation"*).

P13L10: Abstract. *"Soil texture and soil particle size fractions (psf) play an increasing role in physical, chemical and hydrological processes. Many previous studies have used machine-learning and log ratio transformation methods for soil psf interpolation and soil texture classification to improve the prediction accuracy. However, few reports systematically analyzed*

and compared the classification and regression, the accuracies of original (untransformed) and log ratio methods, and the performance of direct and indirect soil texture classification using machine-learning methods. A total of 45 evaluation models generated from five different machine-learning models combined with original and three log ratio methods—additive log ratio, centered log ratio and isometric log ratio (ALR, CLR and ILR, respectively), to evaluate and compare the performance of soil texture classification and soil psf interpolation. The results demonstrated that log ratio methods modified the soil sampling data more symmetrically, and with respect to soil texture classification, random forest (RF) and extreme gradient boosting (XGB) showed better consequences with the overall accuracy (RF: 0.629, XGB: 0.611), kappa coefficients (RF: 0.238, XGB: 0.240) and precision-recall curve (PRC) analysis (RF: 0.646, XGB: 0.616). For soil psf interpolation, RF delivered the best performance among five machine-learning models with lowest root mean squared error (RMSE, sand: 15.09 %, silt: 13.86 %, clay: 6.31 %), mean absolute error (MAE, sand: 10.65 %, silt: 9.99 %, clay: 5.00 %), Aitchison distance (AD, 0.84) and standardized residual sum of squares (STRESS, 0.61), and highest coefficient of determination (R^2 , sand: 53.28 %, silt: 45.77 %, clay: 53.75 %). STRESS was improved using log ratio methods, especially CLR and ILR. There is a pronounced improvement (21.3 %) in the kappa coefficient using indirect soil texture classification compared to the direct approach. With respect to the evaluation of accuracy, RF was recommended as the best strategy among these five machine-learning models according to soil PSF interpolation and soil texture classification. In addition, from the point of view of total computing time of model and sub-optimal accuracy (trade-offs of accuracy and time), XGB was preferred than any other models. Log ratio transformation methods were needed in the evaluation of the indirect soil texture classification and maps of PSFs and texture classes. Our findings can provide a reference for other research of spatial prediction of soil PSF and texture combined with environmental covariates using machine-learning methods with skewed distribution soil PSF data in a large area.”

Minor comments

Comment 1: There is little need to describe the vegetation in the study area, please remove. Also, rainfall patterns are not so interesting when it comes to soil particle size fractions, but I may be mistaken here.

Response: Thanks for the referee’s suggestion. The environmental covariates such as vegetation types and the mean annual precipitation you mentioned were used as independent variables to train five machine-learning models in our study. Therefore, they were useful and we did not delete them.

Comment 2: The English needs a careful check: in general, the manuscript is well readable, but some fine-tuning may further improve the accessibility.

Response: Thanks for the referee’s suggestion for the English in our manuscript. We have improved the overall language of this article and we have checked and improved the writing in the revised version.

Comment 3: What are ‘close correlations’ (p. 2, l. 10)?

“Previous reports revealed that there are **close correlations** between the spatial variations of soil texture and landscape and topography” means that there are strong linear or nonlinear relationship between soil properties (soil texture class was included) and landscape or topography. It can be explained in Jenny’s famous equation—a mechanistic model for soil development:

5
$$S = f(c, o, r, p, t, \dots), \quad (16)$$

where S refers to soil, c (sometimes cl) represents climate, o organisms including humans, r relief, p parent material and t time. For the interpretation of r , the relationship between soil and topographic factors has been evident in previous reports. It also can be explained by the SCORPAN model, which is a better explanation:

$$S_c = f(s, c, o, r, p, a, n), \quad (17)$$

10 where S_c stands for soil classes or other properties of the soil at a point; c is climate, climatic properties of the environment; o is organisms, vegetation or fauna or human activity; r is topography, landscape attributes; p is parent material, lithology; a is age, the time factor; n is space, spatial position. More details can be found in McBratney et al. (2003). We have added the explanation of close correlations.

15 **P14L13:** “Previous reports revealed that there are close correlations of linear or nonlinear relationship between the spatial variations of soil texture and landscape and topography.”

Reference

McBratney, A. B., Santos, M. L. M., and Minasny, B.: On digital soil mapping, Geoderma, 117, 3-52, 10.1016/s0016-
20 7061(03)00223-4, 2003.

Comment 4: Whence the sentence ‘kriging methods (so-called geostatistics)’ (p.3, l. 18)? Kriging is a geostatistical interpolation method; in fact, it is a whole collection of those.

Response: Thanks for the referee’s suggestion. We have improved this sentence in our revised version.

25 **P15L26:** “However, most studies using log ratio approaches to simulate the spatial variation of soil psf were kriging method—a kind of geostatistical interpolation method.”

30 **Comment 5: I object to the use of the term ‘attribute’ for a variable (p. 6). It is a GIS term and not a scientific term, a much better term is ‘variable’. Further down in the paragraph we suddenly read about evapotranspiration data. The manuscript would benefit from a more careful separation between variable and data.**

Response: Thanks for the referee’s suggestion for more careful separation between variable and data in our manuscript. First, we have replaced ‘attribute’ with ‘variable’; second, ‘evapotranspiration data’ has been changed to ‘evapotranspiration variable’.

Comment 6: Notation is inconsistent for k-nearest neighbor (p. 6 ff.). It should either be capitalized throughout, or not at all.

Response: Thanks for the referee's suggestion and we have improved this in our revised version.

5

Comment 7: On page 9 there is the requirement stated that the sum of the components is equal to 1. Fine, but how is that guaranteed? Is a correction being made if it is not the case at a prediction location?

Response: Thanks for the referee's question. *'For soil psf compositional data (i.e. sand, silt and clay), the sum of the components is 1 (or 100 %), which should be guaranteed.'* It can be guaranteed by using log ratio transformation methods (ALR, CLR and ILR in our manuscript), the original soil PSF data was transformed from Aitchison space (simplex space) to Euclidean space (real space), log ratio transformed data then was modeled independently; finally, the results of prediction of soil PSF (in log ratio pattern, i.e. ILR_1 and ILR_2 for ILR transformation method) were back-transformed to the original space (in simplex pattern, i.e. sand, silt and clay), this process can guarantee the sum of soil PSF components was 1 (100 %). The back-transformed equations for these three log ratio methods were as follows:

10

$$\overline{alr}(x_j) = \frac{\exp(alr(x_j))}{\sum_{j=1}^D \exp(alr(x_j))} , \quad (18)$$

15

$$\overline{clr}(x_j) = \frac{\exp(clr(x_j))}{\sum_{j=1}^D \exp(clr(x_j))} , \quad (19)$$

$$Y(x_j) = \sum_{j=1}^D \frac{ilr(x_j)}{\sqrt{j \times (j+1)}} - \sqrt{\frac{j-1}{j}} \times ilr(x_j) , \quad (20)$$

$$ilr(x_0) = ilr(x_D) = 0 , \quad (21)$$

$$\overline{ilr}(x_j) = \frac{\exp(Y(x_j))}{\sum_{j=1}^D \exp(Y(x_j))} , \quad (22)$$

20

Furthermore, with respect to the original method (using original soil PSF data without transformation as dependent variable), the standardization function (Eq. 14 in our manuscript) was used to ensure the sum of predictions of soil PSF was 100%:

$$sand_s = \frac{sand}{(sand+silt+clay)} \times 100 \quad (23)$$

where, $sand_s$ is the content of sand after standardization, the same as silt and clay component. It cannot deal when the negative values are produced in wider regional scale area. Therefore, we recommended using log ratio transformation methods for soil PSF (compositional) data interpolation (spatial perdition). We have added the back-transformed equations for ALR, CLR and ILR transformation methods in **P24L10**.

25

Systematic comparison of five machine-learning methods in classification and interpolation of soil particle size fractions using different transformed data

Mo Zhang^{1,2}, Wenjiao Shi^{1,3}

5 ¹Key Laboratory of Land Surface Pattern and Simulation, Institute of Geographic Sciences and Natural Resources Research, Chinese Academy of Sciences, Beijing 100101, China

²School of Earth Sciences and Resources, China University of Geosciences, Beijing 100083, China

³College of Resources and Environment, University of Chinese Academy of Sciences, Beijing 100049, China

Correspondence to: Wenjiao Shi (shiwj@lreis.ac.cn)

10 **Abstract.** Soil texture and soil particle size fractions (psf) play an increasing role in physical, chemical and hydrological processes. Many previous studies have used machine-learning and log ratio transformation methods for soil psf interpolation and soil texture classification to improve the prediction accuracy. Digital soil mapping using machine-learning methods was widely applied to generate more detailed prediction of qualitative or quantitative outputs than traditional soil mapping methods in soil science. As compositional data, interpolation of soil psf combined with log ratio approaches methods was developed to
15 improve the prediction accuracy, which also can be used to indirectly derive soil texture. However, few reports systematically analyzed and compared the classification and regression, the ~~accuracies~~accuracy of original (untransformed) and log ratio ~~approaches~~methods, and the performance of direct and indirect soil texture classification using machine-learning methods. In this total, a total of 45 evaluation models generated from five different machine-learning models combined with original and three log ratio ~~approaches~~methods—additive log ratio, centered log ratio and isometric log ratio (ALR, CLR and ILR,
20 respectively), to evaluate and compare the performance of soil texture classification and soil psf interpolation. The results demonstrated that log ratio ~~approaches~~methods modified the soil sampling data more symmetrically, and with respect to soil texture classification, random forest (RF) and extreme gradient boosting (XGB) showed notable consequences. For soil psf interpolation, RF delivered the best performance among five machine-learning models with lowest root mean squared error (RMSE, sand: 15.09 %, silt: 13.86 %, clay: 6.31 %), mean absolute error (MAE, sand: 10.65 %, silt: 9.99 %, clay: 5.00 %),
25 Aitchison distance (AD, 0.84) and standardized residual sum of squares (STRESS, 0.61), and highest coefficient of determination (R^2 , sand: 53.28 %, silt: 45.77 %, clay: 53.75 %). STRESS was improved using log ratio ~~approaches~~methods, especially CLR and ILR. There is a pronounced improvement (21.3 %) in the kappa coefficient using indirect soil texture classification compared to the direct method. With respect to the evaluation of accuracy, RF was recommended as the best strategy among these five machine-learning models according to soil PSF interpolation and soil texture classification. In addition, from the point of view of total computing time of model and sub-optimal accuracy (trade-offs of accuracy and time), XGB was preferred than any other models. Log ratio transformation methods were needed in the evaluation of the indirect soil texture classification and maps of PSFs and texture classes. Our findings can provide a reference for other research of spatial
30

prediction of soil PSF and texture combined with environmental covariates using machine-learning methods with skewed distribution soil PSF data in a large area.~~Our systematic comparison helps to elucidate the processing and selection of compositional data in spatial simulation.~~

1 Introduction

5 Soil texture, classified by ranges of soil particle-size fractions (psf), is one of the most important attributes affecting the soil properties and the physical, chemical and hydrological processes covering soil porosity, soil fertility, water retention, infiltration, drainage and aeration. Measuring soil texture can be used for soil fertility management (Pahlavan-Rad and Akbarimoghaddam, 2018), water management (Thompson et al., 2012), maintenance of organic carbon (Bationo et al., 2007) and provision of ecosystem services (Adhikari and Hartemink, 2016). The soil psf, i.e., sand, silt and clay, are vital in most
10 hydrological, ecological, and environmental risk assessment models (Liess et al., 2012). The spatial distributions of soil texture and soil psf affect and control runoff generation, slope stability, depth of accumulation, and soluble salt content (McNamara et al., 2005; Follain et al., 2006; Yoo et al., 2006; Gochis et al., 2010; Crouvi et al., 2013).

Previous reports revealed that there are close correlations of linear or nonlinear relationship between the spatial variations of soil texture and landscape and topography (Gobin et al., 2001; Brown et al., 2004; Zhao et al., 2009; Liess et al., 2012).
15 Compared with traditional soil mapping methods, digital soil mapping (DSM) has an obvious advantage in that it is considerably more economical and efficient; additionally, soil maps using DSM yielded more details because of the development of data-mining algorithms and GIS tools and more extensive application of spatial remote sensing data, particularly in the regional and continental scale. DSM methods were applied by an increasing number of soil scientists to map soil properties using ancillary data (McBratney et al., 2003; Zeraatpisheh et al., 2017), the so-called environmental covariates,
20 which can be obtained from digital elevation models (DEM), remote sensing data, and categorical or geomorphology maps

1 Abbreviations: psf, soil particle-size fractions; HRB, Heihe River Basin; DSM, digital soil mapping; KNN, ~~k~~K-nearest neighbor; MLP, multilayer perceptron neural network; RF, random forest; SVM, support vector machines; XGB, extreme gradient boosting; ALR, additive log-ratio; CLR, centered log-ratio; ILR, isometric log-ratio; ORI, original; ROC, receiver operating characteristics; PRC, precision-recall curve; AUC, area under the ROC curve; AUPRC, area under the PRC; RMSE, root mean squared error; MAE, mean absolute error; R^2 , coefficient of determination; MAD, median absolute deviation; AD, Aitchison distance; STRESS, standardized residual sum of squares; KNN_ALR, KNN_CLR, KNN_ILR, KNN_ORI, MLP_ALR, MLP_CLR, MLP_ILR, MLP_ORI, RF_ALR, RF_CLR, RF_ILR, RF_ORI, SVM_ALR, SVM_CLR, SVM_ILR, SVM_ORI, XGB_ALR, XGB_CLR, XGB_ILR, XGB_ORI, KNN, MLP, RF, SVM, XGB combined with ALR, CLR, ILR, ORI respectively; CILo, clay loam; Lo, loam; LoSa, loamy sand; Sa, sand; SaCILo, sandy clay loam; SaLo, sandy loam; Si, silt; SiCILo, silty clay loam; SiLo, silt loam.

(Krasilnikov et al., 2011). Furthermore, some soil physicochemical attributes, such as soil organic carbon (SOC) and pH, were also permissible to obtain as environmental covariates (Camera et al., 2017). Wang and Shi (2017) also recommended that the soil psf prediction should consider the ancillary data, which can enhance the performance of interpolation.

Different machine-learning methods, such as boosting regression trees (Jafari et al., 2014; Yang et al., 2016), random forests (Hengl et al., 2015; Zeraatpisheh et al., 2017) and artificial neural networks (Bagheri Bodaghabadi et al., 2015; Taalab et al., 2015), have been most commonly employed in DSM models for both regression and classification combined with environmental covariates in soil science. Hengl et al. (2015) contrasted the performance of spatial predictions of soil properties, such as soil psf, using random forests and linear regression, and the results demonstrated that the random forests were superior to the linear regression with remarkable advantages of not only robust to noise but also low bias and variance. Hengl et al. (2017) improved the prediction of organic carbon, bulk density, pH and soil texture fractions on a global scale using machine-learning models – random forest, gradient boosting and multinomial logistic regression – indicating that random forest and gradient boosting outperformed linear models in large data sets. Taghizadeh-Mehrjardi et al. (2015) investigated the predictive power of soil classes using six machine learning-based classifiers and found that artificial neural network and decision trees performed better than any other models they mentioned with relatively high overall ~~aeuracies~~accuracy and kappa coefficients. Heung et al. (2016) evaluated a suite of 10 machine-learning models for predicting soil taxonomic units, and the consequences suggested that although the ~~k~~K-nearest neighbor and support vector machine had the highest accuracy, “tree learners” were preferred because of the interpretability of the results and the speed of parameterization. Most previous studies selected one or more machine-learning algorithms to simulate soil category or continuous variables for classification or regression problems. From this perspective, however, few studies systematically analyzed both soil texture classification and soil psf interpolation using multiple machine-learning methods.

The soil psf, which can be classified as soil texture, are not only continuous variables but also compositional data. We need to pay more attention to the latter case. Numerous different interpretations of the interpolation of compositional data in soil science have been suggested (Gobin et al., 2001; Lark and Bishop, 2007; Salazar et al., 2015), and the most extensively used were a combination of log ratio ~~approaches~~methods involving the additive log ratio (ALR) and the centered log ratio (CLR) put forward by Aitchison (1982), as well as the isometric log ratio (ILR) from Egozcue et al. (2003). However, most studies using log ratio ~~approaches~~methods to simulate the spatial variation of soil psf were kriging ~~methods—method—a kind of geostatistical interpolation method, so-called geostatistics~~, rather than machine-learning methods. Huang et al. (2014) combined multiple linear regression with ALR to improve the prediction precision of soil psf using electromagnetic data on a 1-m transect. Odeh et al. (2003) proposed that modified ALR ordinary kriging transcended compositional kriging and cokriging. Sun et al. (2014) contradistinguished compositional kriging, log ratio cokriging, cokriging, and ALR-cokriging, and produced proximate results. In contrast, Walvoort and de Gruijter (2001) thought compositional kriging had better performance than ALR ordinary kriging. Zhang et al. (2013) suggested compositional kriging was more appropriate for soil texture prediction than symmetry log ratio ordinary (or regression) kriging. Wang and Shi (2018) developed log ratio kriging combined with

robust variogram estimation, which was preferable to compositional kriging methods. However, few studies combined log ratio with machine-learning models for soil psf interpolation in soil science. Aside from those mentioned above, the lack of systematic comparison of accuracy, strengths and weaknesses between original (untransformed) and log ratio ~~approaches~~methods should be considered, especially in terms of combining with machine-learning methods.

Soil texture classification using machine-learning methods can be classified as a dependent variable; furthermore, it also can be derived indirectly from soil psf. Camera et al. (2017) reported that random forests were more remarkable than multinomial logistic regression in the direct soil texture classification. Wu et al. (2018) compared the support vector machines (SVM), artificial neural network (ANN), and classification tree (CT) models, demonstrating better prediction performance generated from SVM than from CT and ANN. For the indirect classification of soil texture, Poggio and Gimona (2017) combined hybrid geostatistical generalized additive models with ALR and modeled soil particle classes at medium resolution (250 m) in Scotland, expecting that vegetation index, morphological features and information about the phenological season were of vital significance as environmental covariates. Considering the particularity of compositional data, the consequences of soil psf classification and regression (indirect soil texture classification and soil psf interpolation, respectively) could be compared from the direct and indirect soil texture classification as a result of the relationship between soil texture and soil psf. Nevertheless, few studies systematically compared these using different machine-learning methods combined with original (untransformed) and log ratio transformed data for both direct and indirect soil texture classification.

In our study, five machine-learning models – ~~k~~K-nearest neighbor (KNN), multilayer perceptron neural network (MLP), random forest (RF), support vector machines (SVM), and extreme gradient boosting (XGB) – were included and applied for DSM of soil texture classification and soil psf interpolation. Furthermore, the original (untransformed) and log ratio transformed data were also combined with the machine-learning algorithms mentioned above for soil psf interpolation. Hence, the objectives of this study are (i) to compare different performance of five machine-learning models in direct soil texture classification, (ii) to evaluate the ~~ae~~accuracy of different log ratio ~~approaches~~methods and original (untransformed) method applied for soil psf from the perspective of compositional data using machine-learning models, and (iii) to estimate whether the ~~ae~~accuracy of indirect soil texture classification using original (untransformed) data and log ratio transformed data were improved compared with the direct soil texture classification.

2 Data and methods

2.1 Study area

The Heihe River Basin (HRB, 97 °6 ' -102 °3 ' E, 37 °43 ' ~ 42 °40 ' N) is situated in the Hexi Corridor, northwest of China, covering the Inner Mongolia Autonomous Region, Gansu and Qinghai provinces (Fig. 1a), which is the second largest inland river basin in China with an area of 146,700 km². The elevation and three reaches (i.e., upper, middle and lower) of the study

area are shown in Fig. 1b. For the upper reaches of HRB, the climate changes significantly with altitude; the mean annual precipitation is 350 mm, the mean annual temperature is from -5-4 °C and the annual average evaporation is 1000 mm. For the middle reaches of HRB, the mean annual precipitation declines between 250 and 50 mm, the annual average evaporation increases from 2000 (east) to 4000 mm (west), and the mean annual temperature is from 2.8 to 7.6 °C. The lower reaches of HRB are situated in Ejina Banner on the Alxa Plateau, which is an arid desert climate with annual precipitation under 50 mm and annual average evaporation above 3500 mm; the mean annual temperature is from 8 to 10 °C.

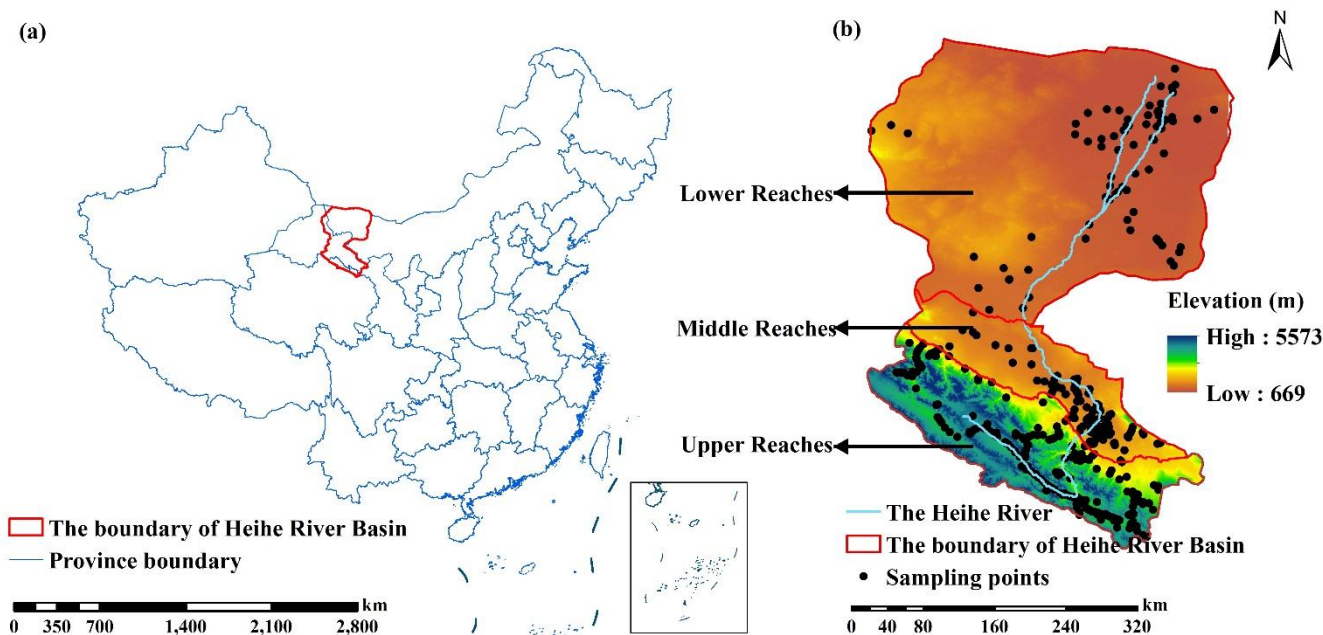


Figure 1. The (a) geographical location, (b) Heihe River, elevation and soil sampling points of Heihe River Basin, China.

The vegetation of the upper reaches of HRB is influenced from the southeast to northwest by hydrothermal conditions. The main vegetation types are alpine vegetation (4000-5000 m), alpine meadow vegetation belt (3000-4000 m), alpine shrub meadow (3200-3800 m), mountain forest meadow belt (2400-3200 m), mountain grassland belt (1800-2400 m), and desert base belt (less than 1800 m). The main vegetation types of the middle and lower reaches of the HRB are relatively fewer, including cultivated vegetation and desert, and the areas near the Heihe River on the lower reaches are shrub and steppe.

The main soil types are frigid desert soils (higher~~less~~ than 4000 m), alpine meadow soil and alpine steppe soil (3600-4000 m), gray cinnamon soil and chernozem (3200-3600 m), sierozem and chestnut soil (2600-3200 m), chestnut soil (2300-2600 m) and sierozem (1900-2300 m) on the upper reaches of the HRB. The main soil types on the middle reaches of HRB are aeolian sandy soil, frigid frozen soil and gray brown desert soil. The main soil types in the lower reaches of HRB are aeolian sandy soil, gray brown desert soil (northwest) and lithosol (northeast).

The main types of geomorphology on the upper reaches of HRB are modern glaciers, alpine and hilly, and ~~intermountain basin~~~~climatic basins~~. Narrow plains are distributed on the middle reaches of HRB. For the lower reaches, the main types of geomorphology are hilly (northwest), plain, sandy land and platform (east), and the area near Heihe River is a flood plain.

2.2 Soil sampling

A total of 640 soil sampling points was collected in the HRB from the Science Data Center of Cold and Arid Regions (WestDC) in China (<http://westdc westgis ac cn/>), involving 392 soil sampling points on the upper reaches and 248 soil sampling points on the middle and lower reaches of the HRB. The soil types, vegetation types, distribution of DEM and geomorphology types of the HRB were considered in soil sample collection according to the location and proportion of these types for the purpose of more representative spatial characteristics of soil psf using limited soil samples. There were more soil sampling points on the middle and upper reaches of HRB due to the more complicated soil types and vegetation types in these areas. In contrast, the types on the lower reaches are relatively similar with more desert in the northwest. Hence, the east of the lower reaches of the HRB contained more soil sampling points. All soil samples had information about soil psf (i.e., sand, silt and clay) and related environmental covariates using a laser diffraction approach and the extraction tool in ArcGIS, respectively, and the global position system (GPS) recorded the position information.

2.3 Environmental covariates and pre-processing

The environmental covariates, such as topographic attributes, remote sensing attributes, climate and position attributes, soil physicochemical attributes and categorical maps, are logically related to the distributions of soil psf. System for Automated Geoscientific Analysis (SAGA) GIS (Conrad et al., 2015) was used to compute ~~their~~~~the~~ topographic attributes from DEM, including slope, aspect, convergence index, general curvature, plane curvature, profile curvature and valley depth. Remote sensing attributes, including the normalized difference vegetation index (NDVI, Huete et al., 2002), the Brightness index (BI, Metternicht and Zinck, 2003), and the soil adjusted vegetation index (SAVI, Huete, 1988) were derived from the Landsat 7 based on band operation. We also collected climate attributes from the National Meteorological Information Center (NMIC, <http://data.cma.cn/>), such as the mean annual precipitation and the mean annual temperature. Latitude and longitude were also considered because of the large scale of the HRB. Mean annual surface evapotranspiration data (Wu et al., 2012) were gathered from WestDC (<http://westdc westgis ac cn/>), as ~~were well as~~ soil physicochemical attributes, such as soil organic carbon, saturated water content, field water holding capacity, wilt water content, saturated hydraulic conductivity, and soil thickness (Yi et al., 2015; Song et al., 2016; Yang et al., 2016), which can address the distributions of soil psf, as well. Additionally, the categorical maps were of significance, such as geomorphology types, soil types, land cover and vegetation types. For slope, the method of dividing the hierarchy rotates clockwise from the north (0°), and each 45° was an interval, including north

(337.5-22.5°), northeast (22.5-67.5°), east (67.5-112.5°), southeast (112.5-167.5°), south (167.5-202.5°), southwest (202.5-247.5°), west (247.5-292.5°), and northwest (292.5-337.5°).

2.4 Machine learning methods and parameters optimization

2.4.1 K-nearest neighbor (KNN)

5 K-nearest neighbor (KNN) is a simple non-parametric classifier based on known instance to label unknown instance (Cover and Hart, 1967). For the test set, ~~k~~K-nearest training set vectors were found, and maximum summed kernel densities were computed for classification. Moreover, continuous variables can also be predicted for regression with the average values of ~~k~~K-nearest neighbors. For a training set of observed data $L = \{(y_i, x_i), i = 1, \dots, n_L\}$, class $y_i \in \{1, \dots, c\}$, and the predictor values $x'_i = (x_{i1}, \dots, x_{ip})$. For a new observation (y, x) , the nearest neighbor $(y_{(1)}, x_{(1)})$ is based on the distance function which is as follows:

$$d(x, x_{(1)}) = \min_i (d(x, x_i)), \quad (1)$$

and $\hat{y} = y_{(1)}$ refers to the nearest neighbor, which is the prediction for y . Value $x_{(j)}$ and $y_{(j)}$ is the j th nearest neighbor of x and class of training set, respectively. Weighted KNN is an extended version of KNN, which considers the maximum of summed kernel densities and the K nearest vectors of training set for each row of the test set (the distances of the nearest neighbors) based on the Minkowski distance, more details can be found in Hechenbichler and Schliep (2004), the equation for Minkowski distance is as follows:

$$d(x_i, x_j) = (\sum_{s=1}^p |x_{is} - x_{js}|^q)^{1/q}, \quad (2)$$

where $d(x_i, x_j)$ refers to the Euclidean distance when $q = 2$ and the absolute distance results for $q = 1$. ~~Weighted KNN is an extended version of KNN that considers the distances of the nearest neighbors; therefore~~Therefore, the parameters of KNN contain the maximum value of k (kmax), the distances of the nearest neighbors (distance) and the types of kernel function (kernel). The KNN model is available in the R package “knn” (Schliep and Hechenbichler, 2016).

2.4.2 Multilayer perceptron neural network (MLP)

Multilayer perceptron neural network (MLP), which is currently one of the most popular multilayer feed forward backpropagation networks, was selected to train artificial neural network (ANN) models in our study due to its rapid operation, small set of training requirements and ease of implementation (Subasi, 2007). MLP neurons can perform classification or regression depending on whether the response variable is categorical or continuous. The MLP has three sequential layers: input layer, hidden layer and output layer. In the hidden layer of MLP, each neuron j sums input environmental covariate in our study x_i after multiplying them by the connection weights w_{ji} respectively, and calculates its output y_j (soil PSF component or texture class) as a function of the sum:

$$y_j = f(\sum w_{ji}x_i), \quad (3)$$

where f is the activation function, which can be a linear (selected in our study) or logistic function. The sum of squared differences between the predicted values and observed values of the output results of neurons E is defined as follows:

$$E = \frac{1}{2} \sum_j (y_{pj} - y_{oj})^2, \quad (4)$$

5 where y_{pj} and y_{oj} is the predicted and observed value of output neuron j , respectively. Each w_{ji} is adjusted to reduce E and the adjustment of w_{ji} depends on the training algorithm (Basheer and Hajmeer, 2000). The resilient backpropagation algorithm was chosen because the learning rate of this algorithm is adaptive, avoiding oscillations and accelerating the learning process (Behrens and Scholten, 2006). The range of the data set should be standardized because MLPs operate in terms of the scale 0 to 1. MLP can be run using the R package “RSNNS” (Bergmeir and Benitez, 2012).

10 2.4.3 Random forest (RF)

Random forest (RF) was developed by Breiman (2001), combining the bagging method (Breiman, 1996) with the random variable selection, and the principle was to merge a group of “weak learners” together to form a “strong learner”. Bootstrap sampling is used for each tree of RF, and the rules to binary split data are different for regression and classification problems. For classification, the Gini index is used to split the data; for regression, minimizing the sum of the squares of the mean deviations can be selected to train each tree model, the equations are as follows:

$$Gini = 1 - \sum_{k=1}^K p_k^2, \quad (5)$$

$$Gini(D, A) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2), \quad (6)$$

$$M = \min_A [\min_{c_1} \sum_{x_i \in D_1(A)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in D_2(A)} (y_i - c_2)^2], \quad (7)$$

15 where p_k refers to the proportion of k th class in the data set on the current node, for feature $A = a$, data set D is divided into two parts (D_1 and D_2), D_1 describes the data set which meets the condition $A = a$ and D_2 is the opposite of D_1 ; $Gini(D, A)$ represents the uncertainty of set D after binary split; y_i is the predicted value of input value x_i , c_1 and c_2 is the mean of data set D_1 and D_2 , respectively. - Benefits of using RFs are that the ensembles of trees are used without pruning. In addition, RF is relatively robust to overfitting, and standardization or normalization ~~are-is~~ not necessary because it is insensitive to the range of value. Two parameters should be adjusted for the RF model: the number of trees (ntree) and the number of features randomly sampled at each split (mtry). The RF model is available in the R package “randomForest” (Liaw and Wiener, 2002).

2.4.4 Support vector machines (SVM)

The support vector machine (SVM), proposed by Cortes and Vapnik (1995), is a type of generalized linear classifier that is widely applied for classification and regression problems in soil science (Burges, 1998). The main principle of SVM is to

classify different classes by constructing an optimal separating hyperplane in the feature space (so called “structural risk minimization”). Regression problems also can be solved by minimization of the structural risk using loss functions (Vapnik, 1998) in SVM, named support vector regression. For a data set $\{x_i, y_i\}, i = 1, \dots, k, x \in R$ and x refers to an n -dimensional vector, $y \in \{-1, +1\}$ is the class corresponding to x , the equation for calculating a hyperplane of SVM is defined as follows:

$$\min_{w, b, \xi} \frac{1}{2} w^T \times w + C \sum_{i=1}^k \xi_i \quad \text{s.t. } y_i(w^T \times \phi(x_i) + b) \geq 1 - \xi_i, \xi_i \geq 0, i = 1, \dots, k \quad (8)$$

where $\phi(x_i)$ refers to the mapping from the input space to the feature space, $C > 0$ is penalty factor (cost), w, b , and ξ are the parameters need to be optimized during the process of model training, which can be determined by the Lagrange multipliers:

$$f(x) = \text{sgn}(y_i a_i k(x_i, x) + b^*) \quad (9)$$

where a_i refers to the support vector, $k(x_i, x)$ refers to the kernel function, and b^* is the bias. The advantages of SVMs are that they are effective in high dimensional spaces. Radial basis function was selected for SVM as the kernel function in our study, and two other parameters need to be tuned, i.e., cost and gamma, controlling the tradeoff between the classification accuracy and complexity, and the ranges of radial effect, respectively. The SVM model is available in the R package “e1071” (Meyer et al., 2017).

2.4.5 Extreme gradient boosting (XGB)

Extreme Gradient Boosting, put forward by Chen and Guestrin (2016), is an efficient method of implementation for gradient boosting frames, tree learning algorithms and efficient linear model solvers to solve both classification and regression problems (Chen et al., 2018). Like the boosted regression trees (Elith et al., 2008), it follows the principle of gradient enhancement; however, more regularized model formalization is applied to XGB to control over-fitting, making it more remarkable. In addition, parallel calculations can be automatically executed during the training phase of the XGB model, presenting a great advantage in large data sets, as the XGB can be more than ten times faster than the existing gradient boosting model (Chen and Guestrin, 2016). The residuals of the first tree can be fitted by the second tree to enhance the model accuracy and the sum of the prediction of each tree generates the ultimate prediction. The general prediction function at step t is defined as follows:

$$f_i^{(t)} = \sum_{k=1}^t f_k(x_i) = f_i^{(t-1)} + f_t(x_i), \quad (10)$$

where $f_t(x_i)$ refers to the tree (learner) at step t , $f_i^{(t)}$ and $f_i^{(t-1)}$ refer to the predicted values at steps t and $t - 1$, and x_i is the input value.

$$Obj^{(t)} = \sum_{k=1}^n l(\bar{y}_i, y_i) + \sum_{k=1}^n \Omega(f_i), \quad (11)$$

where l refers to the loss function, n is the number of data set, and Ω refers to the regularization term, which equation is defined as follows:

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \|\omega\|^2, \quad (12)$$

where ω refers to the score vector, λ is the parameter of regularization term, and γ is the minimum loss. There are seven parameters should be tuned in XGB, containing the learning rate (eta), the maximum depth of a tree (max_depth), the max number of boosting iterations (nrounds), the subsample ratio of columns (colsample_bytree), the subsample ratio of the training instance (subsample), the minimum loss reduction (gamma) and the minimum sum of instance weight (min_child_weight). The XGB model is available in the R package “xgboost” (Chen et al., 2018).

2.4.6 Parameters optimization

The parameters of machine-learning models we mentioned above need to be adjusted, and the numbers of these parameters of models are different. For instance, XGB has seven parameters and is one of the most complicated models; on the other hand, for the MLP, in the case where we have chosen the algorithm, the only parameter that should be tuned is the size of the MLP model.

R package “caret” (Kuhn, 2018) provides an effective grid-search method that can automatically adjust the parameters by setting the adjustment grid, avoiding the uncertainty of artificial adjustment for some models (e.g., XGB) with more parameters. A set of parameters with the lowest RMSE or the highest R^2 for regression and the highest overall accuracy or kappa coefficient for classification by cross-validation can be selected to be the best parameters. However, in the presence of many adjustment parameters, it may be inefficient due to the long training time. Thus, we used the other package of “randomForest” for RF and “kknn” for KNN, which can also restructure the parameters for these two models.

In our study, eleven dependent variables (i.e., ten for regression and one for classification) were trained with environmental covariates (independent variables) for the sake of parameter adjustment for each model, including “sand, silt, clay, ilr1, ilr2, alr1, alr2, clr1, clr2, clr3” and “class”. Subsequently, the parameters were definitely computed; here, we just give the relative ranges of the parameters after adjustment for most dependent variables; for example, in KNN the kmax was 15, the distance was 1, and the kernel was rectangular; in MLP, the size fluctuated between 5 and 10; in RF, the ntree was 1000 and mtry fluctuated from 9 to 11; in SVM, gamma was 0.01 and cost was 1; and in XGB, the range of parameters of max_depth (3-4), eta (0.05-0.1), colsample_bytree (0.6-0.8), nrounds (30), subsample (0.8-1), gamma (0-0.4), and min_child_weight (0.6-0.8) were obtained after conditioning.

2.5 Log-ratio transformation methods

For soil psf compositional data (i.e. sand, silt and clay), the sum of the components is 1 (or 100 %), which should be guaranteed. Soil ~~particle-size~~psf data, including three dimensions, are typical compositional data. The closed number system can be explained as follows: the individual variables in the data set are not independent of each other; moreover, they are related by being expressed as a percentage (Filzmoser et al., 2009). In the Euclidean space, one dimension (variable) would be omitted

for the original method to guarantee no information loss because of the constant-sum constraint. Therefore, the Euclidean space is not appropriate for the analysis of soil psf data. The most widely used [approaches methods](#) are log ratio [approaches methods](#) (Aitchison, 1982), consisting of the additive log ratio, centered log ratio and isometric log ratio (ALR, CLR and ILR for short, respectively) from Aitchison (1982) and Egozcue et al. (2003).

5 For the composition of D elements $\mathbf{x} = [x_1, \dots, x_D]$, $x_j > 0$, $\forall j = 1, \dots, j-1, j+1, 2, \dots, D$, and $\sum_{j=1}^D x_j = 1$, the transformation equation for ALR is defined as follows:

$$alr(\mathbf{x}) = (\ln \frac{x_1}{x_j}, \dots, \ln \frac{x_{j-1}}{x_j}, \ln \frac{x_{j+1}}{x_j}, \dots, \ln \frac{x_D}{x_j}), \quad (413)$$

For soil psf ($D = 3$) in our study, the transformation equations for ALR are:

$$alr(1) = \ln \frac{sand}{clay}, \quad (214)$$

$$10 \quad alr(2) = \ln \frac{silt}{clay}, \quad (315)$$

All of the information regarding the soil psf was contained in $alr(1)$ and $alr(2)$; however, the ALR has been criticized because the choice of denominator is subjective, which can influence the results (Bacon-Shone, 2011). The CLR transformation method can remove this arbitrariness, and the equation is defined as follows

$$clr(\mathbf{x}) = (\cancel{y_1}, \dots, \cancel{y_j}, \dots, \cancel{y_D}) = (\ln \frac{x_1}{\sqrt[D]{\prod_{j=1}^D x_j}}, \dots, \ln \frac{x_D}{\sqrt[D]{\prod_{j=1}^D x_j}}), \quad (416)$$

15 ~~where y_j is the j th component.~~ Similarly, for the soil psf, the transformation equations for CLR are:

$$clr(1) = \ln \frac{sand}{\sqrt[3]{sand \times silt \times clay}}, \quad (517)$$

$$clr(2) = \ln \frac{silt}{\sqrt[3]{sand \times silt \times clay}}, \quad (618)$$

$$clr(3) = \ln \frac{clay}{\sqrt[3]{sand \times silt \times clay}}, \quad (719)$$

20 In the CLR transformation method, the geometric mean composed of all compositions of soil psf is the denominator, and one-to-one mapping of equations and soil psf could be implemented. Nevertheless, the CLR is inapplicable for multivariate analysis because the sum of the dimensions of CLR is 0, and thus the results are collinear. These problems can be overcome by using ILR, which transforms all the information into $D-1$ orthogonal log contrasts (Abdi et al., 2015). The transformation equations for ILR are defined as follows:

$$z = (z_1, \dots, z_{D-1}) = ilr(x), \quad (820)$$

$$25 \quad z_i = \sqrt{\frac{D-i}{D-i+1}} \ln \frac{x_i}{\sqrt[D-i]{\prod_{j=i+1}^D x_j}}, \text{ for } i = 1, \dots, D-1. \quad (921)$$

where z_i is the i th component. The ILR transformation equations for soil psf in our study can also be defined as follows:

$$ilr(1) = \sqrt{\frac{2}{3}} \ln \frac{sand}{\sqrt{silt \times clay}}, \quad (4022)$$

$$ilr(2) = \sqrt{\frac{1}{2}} \ln \frac{silt}{clay}, \quad (423)$$

For a more uniform comparison of the descriptive statistics, the ordering of three components of soil psf followed sand-silt-clay, and we added the third equation for the ALR and ILR. Although all the information could be included in the first two equations, note that in the process of model training interpolation, only the first two equations were used for ALR and ILR:

$$5 \quad alr(3) = \ln \frac{clay}{sand}, \quad (424)$$

$$ilr(3) = \sqrt{\frac{2}{3}} \ln \frac{clay}{\sqrt{sand \times silt}}, \quad (425)$$

The equations for $alr(1), alr(2), alr(3)$ were equivalent to $alr(sand), alr(silt), alr(clay)$ in ALR, the same as in ILR. The back-transformed equations for ALR, CLR and ILR were recommended in our previous research (Wang and Shi, 2017), and were computed in the “compositions” R package (van den Boogaart and Tolosana-Delgado, 2008), which were defined as follows:

$$10 \quad \overline{alr}(x_j) = \frac{\exp(alr(x_j))}{\sum_{j=1}^D \exp(alr(x_j))}, \quad (26)$$

$$\overline{clr}(x_j) = \frac{\exp(clr(x_j))}{\sum_{j=1}^D \exp(clr(x_j))}, \quad (27)$$

$$Y(x_j) = \sum_{j=1}^D \frac{ilr(x_j)}{\sqrt{j \times (j+1)}} - \sqrt{\frac{j-1}{j}} \times ilr(x_j), \quad (28)$$

$$ilr(x_0) = ilr(x_D) = 0, \quad (29)$$

$$15 \quad \overline{ilr}(x_j) = \frac{\exp(Y(x_j))}{\sum_{j=1}^D \exp(Y(x_j))}, \quad (30)$$

For the original (untransformed) method, the standardization function was used to ensure predictions of soil psf were between 0 and 100 and that their sum was 100%:

$$sand_s = \frac{sand}{(sand + silt + clay)} \times 100, \quad (431)$$

where, $sand_s$ is the content of sand after standardization, the same as silt and clay component.

20 2.6 Validation

2.6.1 Validation method

A total of 45 methods that we simulated are presented in Table 1; five machine-learning models were combined with one original (ORI) and three log ratio approaches methods (ALR, CLR, ILR). Five machine-learning methods were applied for direct soil texture classification; additionally, these methods were combined with original (untransformed) and log ratio transformed data for a total of 40 methods for indirect soil texture classification (20) and soil psf interpolation (20). The data were randomly divided into two sets to guarantee prediction accuracies accuracy; for instance, one (70 % = 448 soil samples)

was employed for training models and the other (30 % = 192 soil samples) was set aside for validation. This process was repeated 30 times for soil texture classification and soil psf interpolation, and different indicators were chosen to evaluate different performances of models (or methods).

Table 1. The method system of soil texture classification and soil psf interpolation.

Methods	Soil texture classification		Soil psf interpolation
	Direct classification	Indirect classification	—
Original data (ORI)	KNN, MLP, RF, SVM, XGB	KNN_ORI, MLP_ORI, RF_ORI, SVM_ORI, XGB_ORI	
Log-ratio transformed data (ALR, CLR, ILR)	—	KNN_ALR, KNN_CLR, KNN_ILR, MLP_ALR, MLP_CLR, MLP_ILR, RF_ALR, RF_CLR, RF_ILR, SVM_ALR, SVM_CLR, SVM_ILR, XGB_ALR, XGB_CLR, XGB_ILR,	

5 **2.6.2 Validation indicators for soil texture classification**

The overall accuracy (Brus et al., 2011) and kappa coefficient were selected to evaluate the overall effects of soil texture types predicted by different models. Moreover, the receiver operating characteristic (ROC) curve, precision-recall curve (PRC), area under the ROC curve (AUC), area under the precision-recall curve (AUPRC) and abundance index were applied to evaluate the performance of different soil texture types.

10 The overall accuracy represents all samples of soil texture types correctly classified by machine-learning models, divided by the total number of samples of soil texture types used in the validation. The higher overall accuracy, the more accurate soil map (Brus et al., 2011):

$$Overall\ Accuracy = \frac{TP+TN}{TP+TN+FP+FN},$$
(4532)

15 where T, F, P and N denote True, False, Positive, and Negative and TP, TN, FP, FN were true positive, true negative, false positive, and false negative, respectively. When the numbers of samples in different classes are imbalanced in the data set, the kappa coefficient can explain the agreement of classes (Marchetti et al., 2011), which is calculated based on the confusion matrix, the equation is defined as:

$$kappa = \frac{p_o - p_e}{1 - p_e},$$
(4633)

20 where, p_o is the probability of observed agreement (overall accuracy) and p_e is the probability of agreement when two classes are unconditionally independent. The strength of the kappa coefficients is interpreted in the following manner: 0.01-0.20: slight, 0.21-0.40: fair, 0.41-0.60: moderate, 0.61-0.80: substantial, 0.81-1.00: almost perfect (Landis and Koch, 1977).

The probabilities of different soil texture types (sum to 1) obtained during the training and predicting processes of machine-learning models were selected to calculate the sensitivity, specificity, precision and recall:

$$Sensitivity = recall = \frac{TP}{TP+FN}, \quad (1734)$$

$$Specificity = \frac{TN}{TN+FP}, \quad (1835)$$

$$5 \quad Precision = \frac{TP}{TP+FP}, \quad (1936)$$

In general, sensitivity, precision and recall indicate the extent of identifying positive cases, and specificity demonstrates the extent of identifying the negative cases of models. ROC analysis is commonly used in two-class problems. However, soil texture types are more than two classes. In our point of view, a one-vs-rest strategy was employed to produce different ROC graphs for each soil texture type.

$$10 \quad P_i = c_i, \quad (2037)$$

$$N_i = \cup j \neq i c_j \in C, \quad (2138)$$

where C is the set including all classes, P_i is the positive class, N_i is the negative class, including all classes except c_i in ROC graph i (Fawcett, 2006).

15 In practice, the weakness of the ROC curve is that it cannot indicate the differences among the models in the cases of imbalanced samples between positive and negative. Soil texture data are a class-imbalanced data set of positive and negative, and the negative classifier would be overvalued under these circumstances because of the overabundance of majority (negative) examples, additionally revealing overly optimistic findings (Davis and Goadrich, 2006). However, precision and recall curves (PRC) are more informative than ROC curves in dealing with class-imbalanced data (Fu et al., 2017). The R package “precrec” (Saito and Rehmsmeier, 2017) generated ROC and PRC curves and computed AUC and AUPRC for each soil texture type.

20 This process was repeated 30 times and eventually, the average ROC and PRC curves with their average areas under these curves were obtained.

Abundance index was applied to describe the proportion of all soil texture types and well-classified soil texture types in the prediction map, which was defined as follows:

$$Abundance\ index = p/t, \quad (2239)$$

25 where p is all soil texture types in the prediction map and t is well-classified soil texture type(s) in test sets. For the sake of ensuring the balance of the soil texture types, all nine soil texture types were involved in test sets, covering clay loam (CILo: 12), loam (Lo: 57), loamy sand (LoSa: 18), sand (Sa: 23), sandy clay loam (SaCILo: 4), sandy loam (SaLo: 58), silt (Si: 31), silty clay loam (SiCILo: 37), and silt loam (SiLo: 400); most were SiLo (62.5%) and the fewest were SaCILo (0.63%).

2.6.3 Validation indicators for soil psf interpolation

30 The accuracy and performance of machine-learning models mentioned above for the original (untransformed) and different log ratio transformation [approaches](#) were evaluated using five statistical indicators, containing coefficient of

determination (R²), root mean square error (RMSE), mean absolute error (MAE), Aitchison distance (AD, Aitchison, 1992), and standardized residual sum of squares (STRESS, Martin-Fernandez et al., 2001). The equations for the validation indicators R², RMSE, MAE, AD and STRESS are as follows:

$$R^2 = \frac{\sum_{i=1}^n (Y_{i,m} - Y_{i,e})^2}{\sum_{i=1}^n (Y_{i,m} - \bar{Y}_{i,m})^2} \quad (2340)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_{i,m} - Y_{i,e})^2}, \quad (2441)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |Y_{i,m} - Y_{i,e}|, \quad (2542)$$

where $Y_{i,m}$, $Y_{i,e}$, $\bar{Y}_{i,m}$ and n are the measured, predicted and the mean of measured soil psf and the number of observations (soil sampling points for validation). Closer to 1 and higher values of R² and the lower values of RMSE and MAE show better performance of models and methods.

$$AD = \left[\sum_{i=1}^D \left[\log \frac{x_i}{g(x)} - \log \frac{X_i}{g(X)} \right]^2 \right]^{1/2}, \quad (2643)$$

$$STRESS = \left[\frac{\sum_{i < j} (AD_{x,ij} - AD_{X,ij})^2}{\sum_{i < j} (AD_{x,ij})^2} \right]^{1/2}, \quad (2744)$$

where x is the observed value; X is the predicted value; D is the number of dimensions (for soil psf is 3); $g(x)$ denotes the geometric mean $(x_1 \dots x_D)^{1/D}$; $AD_{x,ij}$ and $AD_{X,ij}$ are the ADs between the observed soil psf and the predicted soil psf at sites i and j . Both present that model performances are better when the values are lower.

2.6.4 Indirect soil texture classification by soil psf interpolation

Seventy percent of the 640 soil sampling points were used for training each machine-learning model, and the remaining 30 % were used for the soil psf interpolation; thereafter, we transformed the content of three components (sand, silt and clay) into the soil texture types in the USDA soil texture classification using the R package “soiltexture” (Moeys, 2018). Eventually, the overall accuracy and kappa coefficient were computed and evaluated. This process was repeated 30 times, and the averages of these consequences were employed to compare the classification performance ~~of~~ for each model. The direct and indirect soil texture classifications were also compared with the overall accuracy and kappa coefficient. The training and testing sets for each time were the same by setting seeds, and all calculations and analysis were performed with the freely available software R (R Core Team, 2018).

2.7 Statistical analysis for the original and log-ratio transformed data

The mean, median, minimum (Min), maximum (Max), median absolute deviation (MAD), skewness (Skew), kurtosis and Kolmogorov-Smirnov test ($p > 0.05$) were employed for descriptive statistical analysis of the original (untransformed) and log

ratio transformed soil psf data. The arithmetic mean of log-ratio transformation data should be back-transformed to the original space. For $X = [X_1, \dots, X_n]$, the MAD can be calculated according to the Eq. (28) as below:

$$MAD(X) = \text{median}(|X_i - \text{median}(X)|). \quad (2845)$$

3 Results

5 3.1 The descriptive statistics for the original and log-ratio transformed soil psf data

With respect to the original (untransformed) data of sand, the mean fraction (30.64 %) was much higher than that of median fraction (25.10 %); conversely, both silt and clay were the opposite, with lower mean fractions (silt: 55.79 %, clay: 13.57 %) than median fractions (silt: 59.47 %, clay: 13.78 %). For the log ratio transformed data, the means of sand (28.69 %) and silt (60.54 %) were closer to the median values of the original data, aside from clay, with mean of 10.78 %.

10 All MADs of log ratio transformed data were much smaller than those of the original data in all cases; for instance, ILR contained the best value of MAD for sand (0.66) and clay (0.44), and CLR generated the lowest MAD for silt (0.43) among different log ratio [approachesmethods](#) (Fig. 2). All log ratio [approachesmethods](#) had lower skews (ALR: 0.77, CLR: 0.88, ILR: -1.20) than those of the original data (1.24) ~~for-of~~ sand. Moreover, CLR (-0.4) declined the original skew (-0.93) ~~for-of~~ silt. However, it was negligible for log ratio transformation data compared with the original skew of clay (0.4). The kurtosis of all

15 log ratio [approachesmethods](#) ~~was-were~~ much higher compared with the consequences generated from original (untransformed) data. In terms of the k-s test ($p < 0.05$), although the p values of the original (untransformed) and different log ratio transformed data were not significant and all histograms were not subject to normal distribution, log ratios made the images of the data more symmetric (Fig. 2).

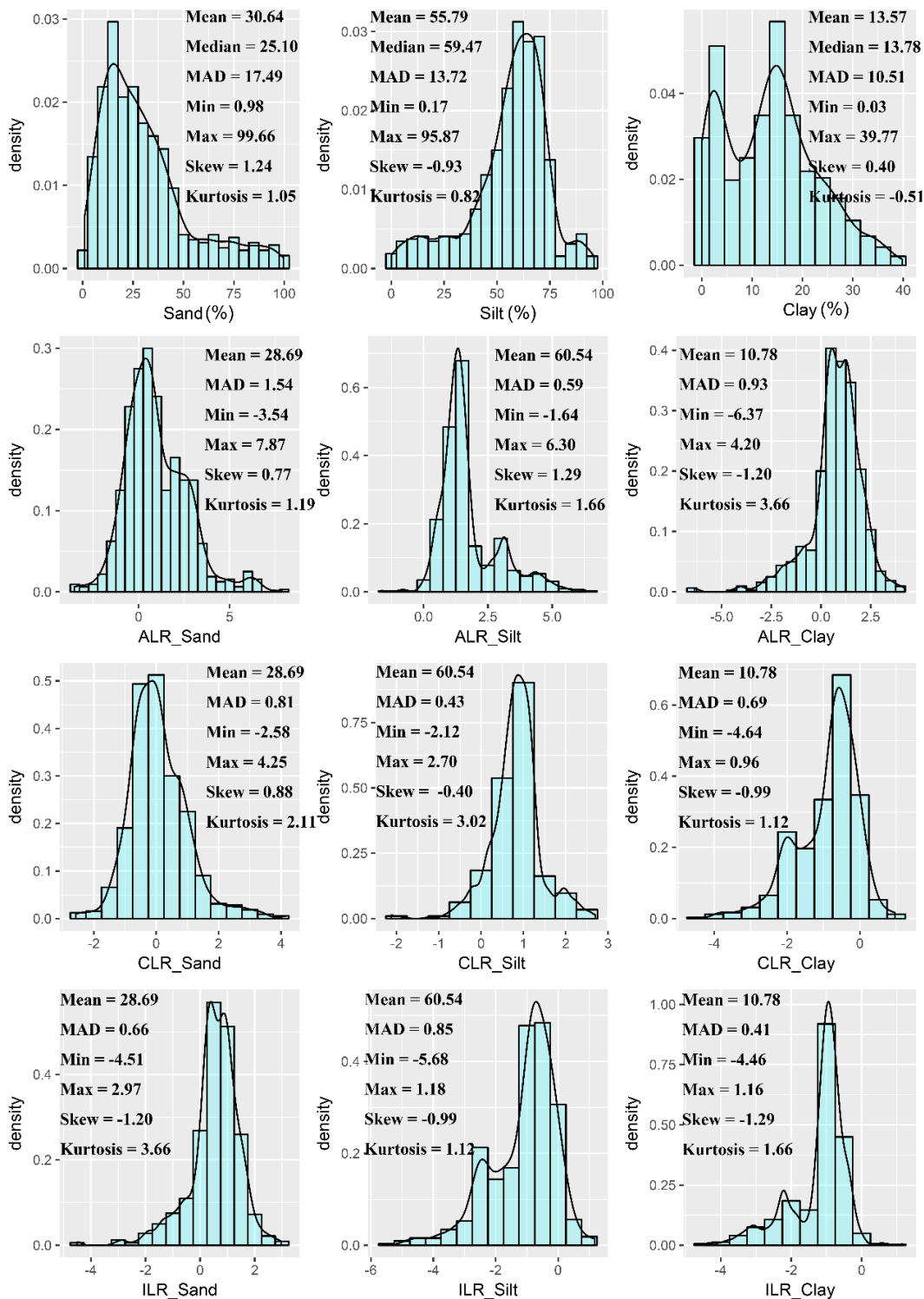


Figure 2. Descriptive statistical analysis for the original (untransformed) and logratio transformed soil sampling data. Note that the means generated from log ratio transformed data were back-transformed to the real space.

3.2 Comparison of the machine learning models in the classification of soil texture types

3.2.1 Comparison of the validation indicators for soil texture classification

The overall accuracy of each model ranged from 0.610 to 0.647 (Fig. 3a). SVM had the highest overall accuracy (0.647) among the five models, followed closely by the ~~accuracies~~accuracy of KNN (0.631) and RF (0.629). XGB (0.611) and MLP (0.610) were relatively lower among these models. The highest kappa coefficient was generated from XGB (0.240), followed by RF (0.238), KNN (0.234) and MLP (0.230), and the worst performer was SVM, with kappa coefficient dropping to 0.186 (Fig. 3b).

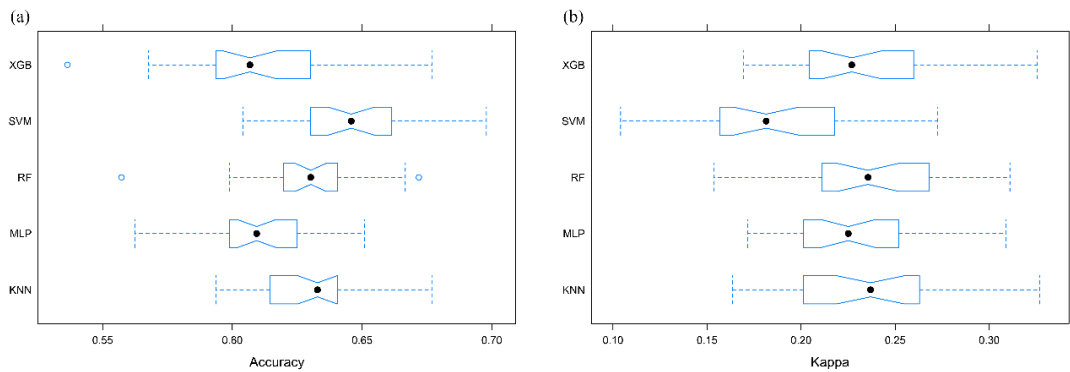


Figure 3. (a) The overall ~~accuracies~~accuracy and (b) kappa coefficients for different machine learning models of KNN, MLP, RF, SVM and XGB.

The AUC with regard to each soil texture type of 640 soil sampling points predicted from five different models demonstrated that the ranking of the AUC was RF>XGB>SVM>KNN>MLP in the case of fewer soil sampling points (CILo, LoSa, Sa, SaCILo and Si). However, in the case of the types with more soil sampling points (Lo, SaLo, SiLo, SiCILo), the ROC curves exhibited roughly the same shape for each model (Fig. 4); therefore, the order of performance was as follows: RF>SVM>XGB>MLP>KNN.

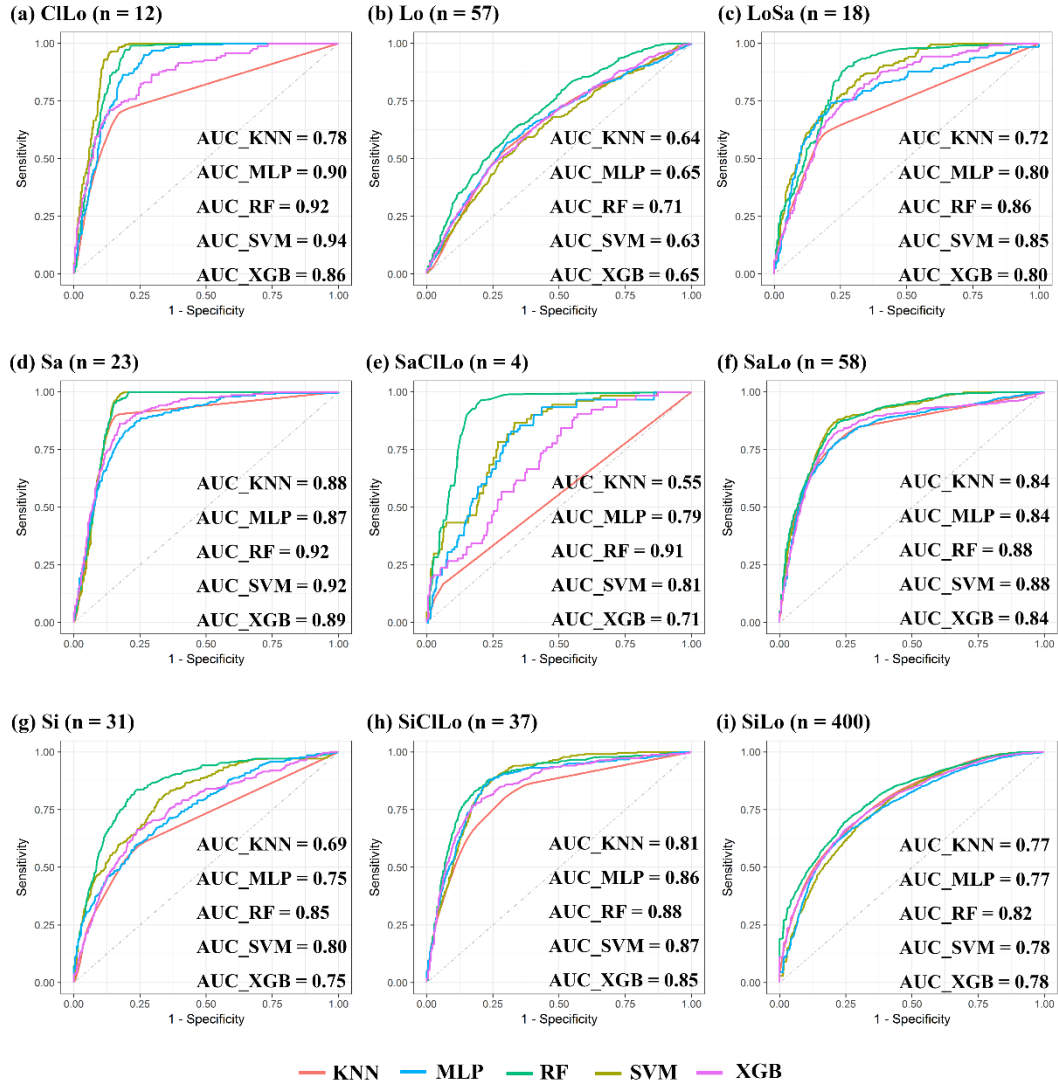
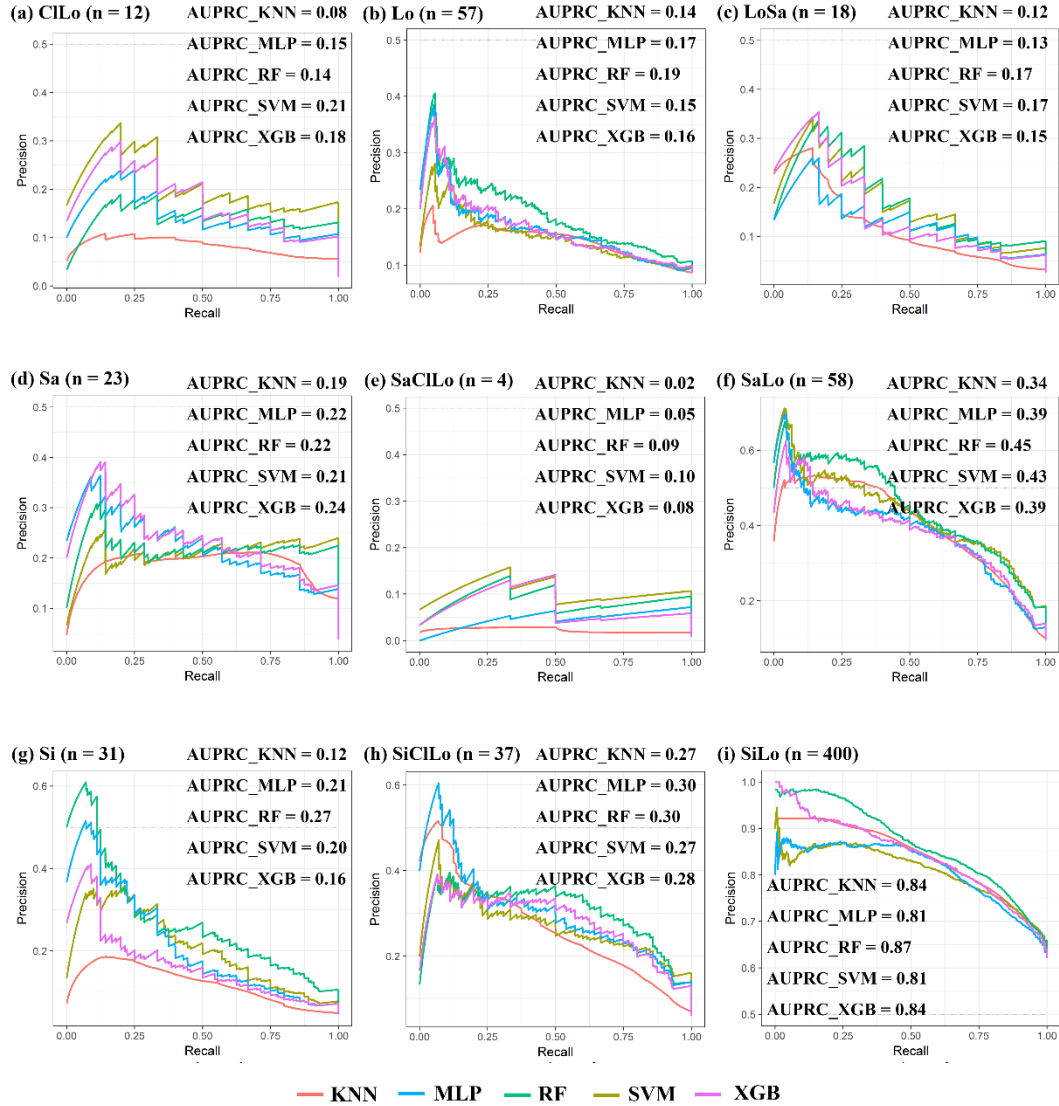


Figure 4. The AUC for different machine learning methods of each soil texture type (a) CILo (b) Lo (c) LoSa (d) Sa (e) SaCILo (f) SaLo (g) Si (h) SiCILo (i) SiLo; n was the sampling points for different soil texture types.

We combined the PRCs with five machine-learning methods to evaluate the performance of these models with respect to predicting each soil texture type using soil psf imbalanced data with different samples of soil texture types (Fig. 5). We found that the AUPRC of types with fewer positive examples were typically small, especially in the case of SaCILo (only four samples), which resulted in unsatisfying consequences because the lack of soil sampling points made models learn poorly during the training process. Hence, the soil texture types (Lo, SaLo, SiLo, SiCILo) with more positive examples delivered superior results to those with fewer positive examples. Moreover, these soil texture types had significant differences in

AUPRCs. For example, SiLo, which had the largest number of samples, was the most effective among these nine types. The total AUPRC calculated by the weights of samples for AUPRC of each type was applied to evaluate the effect of each model, and the order was as follows: RF (0.646)>XGB (0.616)>KNN (0.601)>MLP (0.600)>SVM (0.599).



5 **Figure 5.** The AUPRC for different machine learning methods of each soil texture type (a) ClLo (b) Lo (c) LoSa (d) Sa (e) SaClLo (f) SaLo (g) Si (h) SiClLo (i) SiLo; n was the sampling points for different soil texture types.

3.2.2 Comparison of the prediction maps for soil texture classification

Prediction maps of soil texture types in the HRB using machine-learning models delivered quite different spatial distributions in the overall performance of different models (Fig. 6). The abundance indices pointed out that all models could not predict the type of SaCILo; in other words, KNN and XGB predicted 8 of 9 types, followed closely by RF (7 of 9 types) and MLP (6 of 9 types). However, SVM predicted only two types, which was an unsatisfactory result associated with the lowest kappa coefficient (Fig. 3). Additionally, the prediction effects of different models were different in the distributions of soil texture types in the HRB. The consequences of RF and XGB illustrated that the main soil texture types in the northwest of the lower reaches of HRB were mostly LoSa, while other prediction models produced SaLo. On the upper reaches of the HRB, soil texture types generated from RF were more abundant and more in accordance with the real environment.

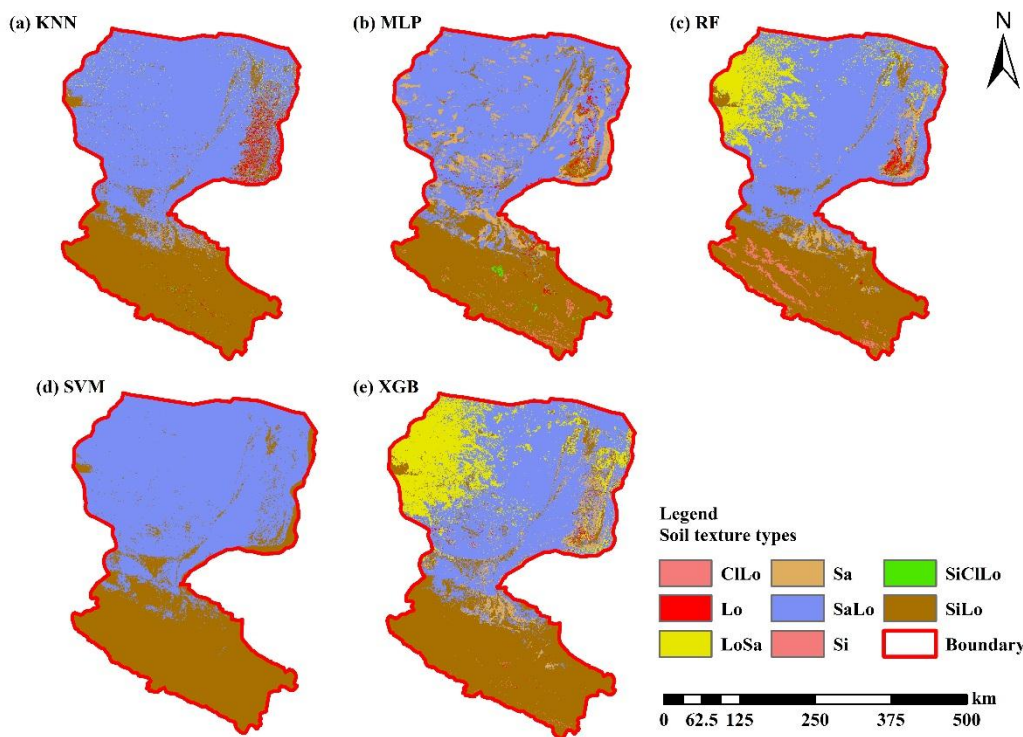


Figure 6. Soil texture classification prediction map of different soil texture types of (a) KNN, (b) MLP, (c) RF, (d) SVM and (e) XGB.

3.3 Comparison of the machine learning models combined with log-ratio transformed methods in the interpolation of soil psf

3.3.1 Comparison of the validation indicators for interpolation of soil psf

We compared the performance of each machine-learning model combined with the original (untransformed) and the log ratio transformed data of soil psf. The results indicated that the ~~accuracies~~accuracy of STRESS of the methods combined with log ratio transformed data were superior to other approaches~~methods~~ using original (untransformed) data (Table 2). With respect to KNN, MLP, RF and XGB, the RMSE, MAE, R^2 and AD generated from original (untransformed) data outperformed log ratio transformed data; for SVM, log ratio transformed data delivered superior improvement. For instance, SVM_CLR and SVM_ILR had higher R^2 and lower RMSE and MAE than SVM_ORI of sand, silt and clay.

By comparison among different log ratio transformed data of the same machine-learning model, ILR and CLR outperformed ALR in these models, other than MLP, showing a slight difference. As shown in Table 2, KNN_CLR demonstrated the most remarkable performance among the three KNN models using different log ratio transformed data with highest R^2 (sand: 48.48 %; silt: 38.37 %; clay: 41.43 %) and lowest RMSE (sand: 15.82 %; silt: 14.77 %; clay: 7.09 %) and MAE (sand: 11.21 %; silt: 10.74 %; clay: 5.58 %). Furthermore, CLR and ILR generated relatively similar consequences for each model of RF and SVM; with respect to XGB, XGB_ILR showed the best performance with all indicators we measured, aside from RMSE (6.75 %) and MAE (5.36 %) of clay, and STRESS (0.63).

We also compared five different machine-learning models using the same log ratio transformation approaches~~methods~~. In the case of ALR, ALR_RF had talent, with the lowest RMSE (sand: 15.50 %; silt: 14.43 %; clay: 6.62 %) and MAE (sand: 10.90 %; silt: 10.52 %; clay: 5.24 %), the highest R^2 (sand: 50.57 %; silt: 41.23 %; clay: 48.90 %), and the lowest AD (0.86) and STRESS (0.61), followed by SVM_ALR, XGB_ALR, KNN_ALR and MLP_ALR. Regarding CLR and ILR, RF also produced the most preferable performance followed by SVM, XGB, KNN and MLP. For original (untransformed) data, RF outperformed other models in accordance with log ratio approaches~~methods~~, and the next were XGB, SVM, KNN and MLP. Therefore, it is clear that RFs demonstrated the most extraordinary indicators of RMSE, MAE, R^2 and AD from the untransformed model and the best STRESS from the log ratio models (RF_ALR, RF_CLR and RF_ILR).

Table 2. The comparisons of ~~accuracy~~accuracy of different machine-learning models combined with original (untransformed) and transformed data.

	RMSE (%)			MAE (%)			R ² (%)			AD	STRESS
	Sand	Silt	Clay	Sand	Silt	Clay	Sand	Silt	Clay		
KNN_ALR	16.05	15.04	7.12	11.35	10.93	5.59	47.02	36.11	41.07	0.90	0.62
KNN_CLR	15.82	14.77	7.09	11.21	10.74	5.58	48.48	38.37	41.43	0.88	0.62
KNN_ILR	15.82	14.82	7.14	11.22	10.84	5.60	48.46	37.88	40.74	0.88	0.64
KNN_ORI	15.51	14.47	7.05	11.12	10.51	5.49	50.59	40.92	42.24	0.84	0.66
MLP_ALR	15.83	15.07	7.43	11.42	11.06	5.97	48.50	35.82	35.79	0.92	0.66
MLP_CLR	15.84	15.07	7.41	11.45	11.05	5.96	48.42	35.86	36.19	0.92	0.66
MLP_ILR	15.84	15.07	7.40	11.46	11.04	5.95	48.40	35.85	36.32	0.92	0.66
MLP_ORI	15.80	14.72	6.96	11.50	10.85	5.52	48.75	38.84	43.72	0.90	0.68
RF_ALR	15.50	14.43	6.62	10.90	10.52	5.24	50.57	41.23	48.90	0.86	0.61
RF_CLR	15.28	14.22	6.61	10.70	10.25	5.21	51.95	42.89	49.16	0.86	0.61
RF_ILR	15.27	14.25	6.66	10.66	10.26	5.26	51.99	42.60	48.28	0.86	0.61
RF_ORI	15.09	13.86	6.31	10.65	9.99	5.00	53.28	45.77	53.75	0.84	0.66
SVM_ALR	15.66	14.59	6.76	11.66	10.88	5.34	49.61	39.87	46.89	0.88	0.66
SVM_CLR	15.27	14.36	6.87	11.01	10.41	5.41	52.12	41.85	45.14	0.87	0.65
SVM_ILR	15.29	14.37	6.84	10.92	10.43	5.42	51.99	41.69	45.58	0.87	0.65
SVM_ORI	15.30	14.38	6.92	10.94	10.32	5.43	51.98	41.71	44.45	0.87	0.67
XGB_ALR	15.82	14.92	6.72	11.32	11.01	5.35	48.57	37.23	47.44	0.88	0.64
XGB_CLR	15.70	14.80	6.75	10.96	10.67	5.39	49.23	38.10	46.90	0.88	0.62
XGB_ILR	15.45	14.57	6.75	10.91	10.52	5.36	50.88	40.01	47.01	0.88	0.63
XGB_ORI	15.15	14.05	6.47	10.88	10.15	5.15	52.85	44.27	51.36	0.86	0.68

3.3.2 Comparison of the interpolation maps of soil psf

Interpolation maps of soil psf (sand, silt and clay) using log ratio transformed data (ILR) and original (untransformed) data were represented in Figs. 7, S1 and S2. At first glance, there was a negligible difference between ILR and ORI based on the same machine-learning model. However, the maps generated from models combined with ILR transformed data showed closer ranges to the original soil sampling data in the case of sand (0.98-99.66 %), silt (0.17-95.87 %) and clay (0.03-39.77 %), and the texture features were more suitable for the distributions of the real environment (Figs. 7, S1 and S2). With respect to different machine-learning models, RF and XGB delivered more detailed information about texture features in prediction maps than did KNN, SVM and MLP.

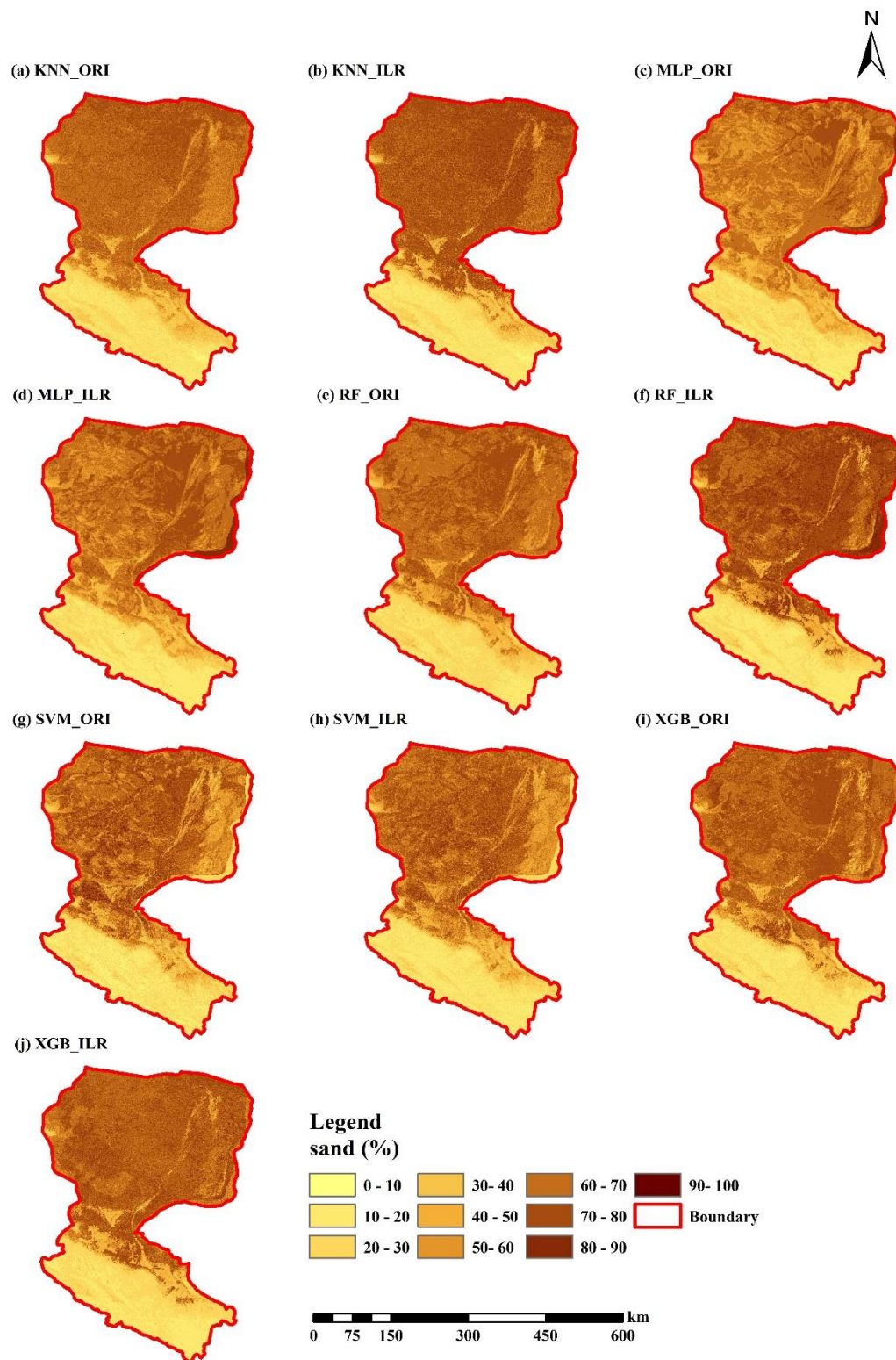


Figure 7. The interpolation maps of sand fraction. All the ranges of prediction maps of sand (approximately 9.0—90.0 %) were within the range of the original data (0.98—99.66 %). RF_ILR (7.9—94.7 %) and XGB_ORI (1.8—92.4 %) generated wider output distributions and were relatively closer to the range of the distribution of original data than other prediction maps such as KNN_ILR (7.3—88.6 %), KNN_ORI (7.8—80.8 %), MLP_ILR (8.8—90.8 %), MLP_ORI (9.0—90.3 %), RF_ORI (9.0—81.0 %), SVM_ILR (6.5—85.6 %), SVM_ORI (7.3—90.0 %) and XGB_ILR (5.0—88.5 %).

3.4 Comparison of direct and indirect soil texture classification

3.4.1 Comparison of the validation indicators for direct and indirect soil texture classification

Compared with the classification performance of the five machine-learning models using original (untransformed) data, the overall ~~accuracies~~accuracy and kappa coefficients of models ~~combined~~ using log ratio transformed data were improved, especially RF and XGB, which combined with all three log ratio ~~approaches~~methods were superior to the interpolation methods using original data. Table 3 ~~shows~~showed that the overall accuracy (0.631) and kappa coefficient (0.245) of the original method in KNN models were better than any other log ratio transformed methods. In summary, the ILR transformation method of five machine-learning models showed the highest overall accuracy among three log ratio transformation ~~approaches~~methods (KNN: 0.628; MLP: 0.614; RF: 0.631; SVM: 0.631; XGB: 0.632), which also demonstrated the best performance with regard to kappa coefficients (KNN: 0.244; RF: 0.291; SVM: 0.239; XGB: 0.252), except for MLP (ALR: 0.216; CLR: 0.216; ILR: 0.214). We also compared direct classification (Fig. 3) with indirect classification and found that the highest values of overall accuracy of indirect classification (KNN: 0.631; MLP: 0.614; RF: 0.628; SVM: 0.638; XGB: 0.632) were slightly decreased in comparison of direct classification (KNN: 0.631; MLP: 0.610; RF: 0.629; SVM: 0.647; XGB: 0.611) for RF and SVM, and improved or kept stable for MLP and XGB, and KNN, respectively. In turn, the kappa coefficients were greatly modified using indirect classification (KNN: 0.245; MLP: 0.216; RF: 0.291; SVM: 0.239; XGB: 0.252) compared with direct classification (KNN: 0.234; MLP: 0.230; RF: 0.238; SVM: 0.186; XGB: 0.240), other than MLP; peculiarly, RF_ILR increased the kappa coefficient to 0.291 (21.3 % improvement) while keeping accuracy stable, ~~which showed~~showing the highest kappa coefficient among these methods.

Table 3. Overall ~~accuracies~~accuracy and kappa coefficients calculated from soil texture classification by the interpolated maps from five models using original (untransformed) data and log ratio transformed data.

Methods	Overall accuracy	Kappa coefficient
KNN_ALR	0.623	0.236
KNN_CLR	0.627	0.241
KNN_ILR	0.628	0.244
KNN_ORI	0.631	0.245
MLP_ALR	0.614	0.216

MLP_CLR	0.614	0.216
MLP_ILR	0.614	0.214
MLP_ORI	0.611	0.216
RF_ALR	0.619	0.284
RF_CLR	0.625	0.276
RF_ILR	0.628	0.291
RF_ORI	0.619	0.279
SVM_ALR	0.591	0.205
SVM_CLR	0.630	0.227
SVM_ILR	0.631	0.239
SVM_ORI	0.638	0.232
XGB_ALR	0.610	0.226
XGB_CLR	0.612	0.240
XGB_ILR	0.632	0.252
XGB_ORI	0.619	0.239

3.4.2 The prediction performance of soil texture types from different methods

The distributions of soil texture classes using original (untransformed) data and ILR transformed data are illustrated in the USDA soil texture triangle (Fig. 8). The triangle of the original data (Fig. 8a) shows wider ranges of spatial dispersion than the interpolation data using machine-learning models, revealing the properties of aggregate from the sides to the center of triangles. With respect to these machine-learning models, RF showed the most dispersed feature in accordance with the original data. The distributions predicted from models combined with ILR transformed data were more discrete and more associated with the original soil psf data than those resulting from ORI [approachesmethods](#). The results of prediction represented striking differences in that the error ratio (red color) of soil sampling points on types of LoSa, SaLo and Lo (left side of triangles) were significantly more than those on types of SiLo and Si (the right side of triangles) for most models, especially KNN and MLP.

The log ratio [approachesmethods](#) overestimated the content of silt in the process of transformation (Fig. 2); in this way, these points were biased to the right of the USDA soil texture triangle based on overall contraction (regression smoothing effects), crossing the classification boundary and becoming other soil texture types. RF_ILR (Fig. 8f) delivered the highest right ratio (RR) among these models, and the classification accuracy was enhanced using the ILR method (83.9%) compared with the ORI method (81.7%). In the case of other models, the differences between original and log ratio [approachesmethods](#) were negligible. We also compared the RR of indirect classification models with those of direct classification, demonstrating all RRs of direct classification were higher (KNN: 67.97 %; MLP: 75.16 %; RF: 100 %; SVM: 66.09 %; XGB: 81.09 %),

especially RF and XGB; however, we removed this evaluation indicator because the same data sets were employed in the processes of training and predicting.

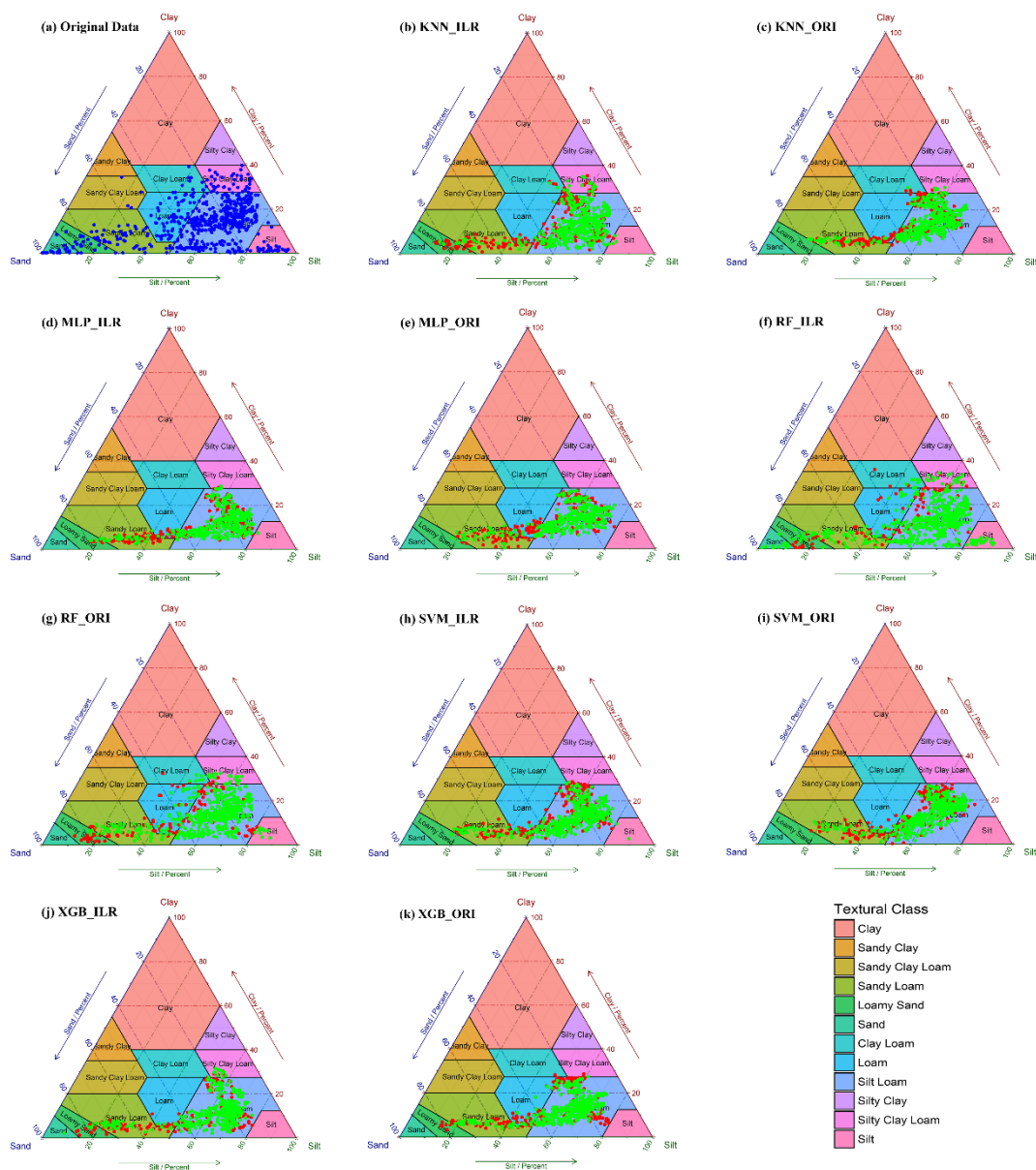


Figure 8. Soil texture types of 640 soil samples shown in USDA texture triangle. The results of soil psf were generated from
 5 (a) original (untransformed) data, (b) KNN_ILR (65.0 %), (c) KNN_ORI (65.9 %), (d) MLP_ILR (63.3 %), (e) MLP_ORI (63.6 %), (f) RF_ILR (83.9 %), (g) RF_ORI (81.7 %), (h) SVM_ILR (66.1 %), (i) SVM_ORI (66.4 %), (j) XGB_ILR (67.8 %), and (k) XGB_ORI (68.0 %). Note that the predicted right-ratios (RRs) of the soil texture types were in the bracket after interpolators

3.4.3 Comparison of prediction maps of direct and indirect soil texture classification

Fig. 9 ~~indicated~~~~shows~~ the ~~similarities~~~~similarity~~ of the three log-ratio transformation methods. The soil texture maps predicted using original data is different from those generated ~~from~~~~by~~ log-ratio transformed data, and the classification maps from the machine learning models combined with the log-ratio transformed data had more detailed information. Three log-ratio transformation methods of the same machine learning model ~~are~~~~were~~ similar in the number of each type predicted; however, there are some differences between methods using original data and those using log-ratio transformed data. All machine learning models combined with original data predicted more types of Lo and SaLo, and less types of LoSa and Si, which could also be presented in Fig. 9. The performance of different machine learning models, especially in the ~~lower~~~~fewer~~ reaches of the Heihe River Basin ~~was~~~~were~~ also compared, for log-ratio transformation methods, for KNN, KNN_ALR and KNN_CLR predicted more type of LoSa than KNN_ILR in the north of lower reaches; for each model of MLP and RF, the differences were slight; more types of Lo in the northwest of lower reaches and less LoSa near the Heihe River were generated by SVM_ALR, compared with SVM_CLR and SVM_ILR; for XGB, the performance of three maps were different due to the prediction of LoSa. We also compared the prediction of the soil texture types by direct classification (Fig. 6) with those generated ~~by~~~~from~~ indirect classification using the same machine learning model, ~~revealing~~~~and~~~~found~~ completely difference between them on the lower reaches of Heihe River Basin, such as the distribution of LoSa; on the middle and upper reaches of Heihe River Basin, all the prediction maps were similar, mainly distributed with SiLo.

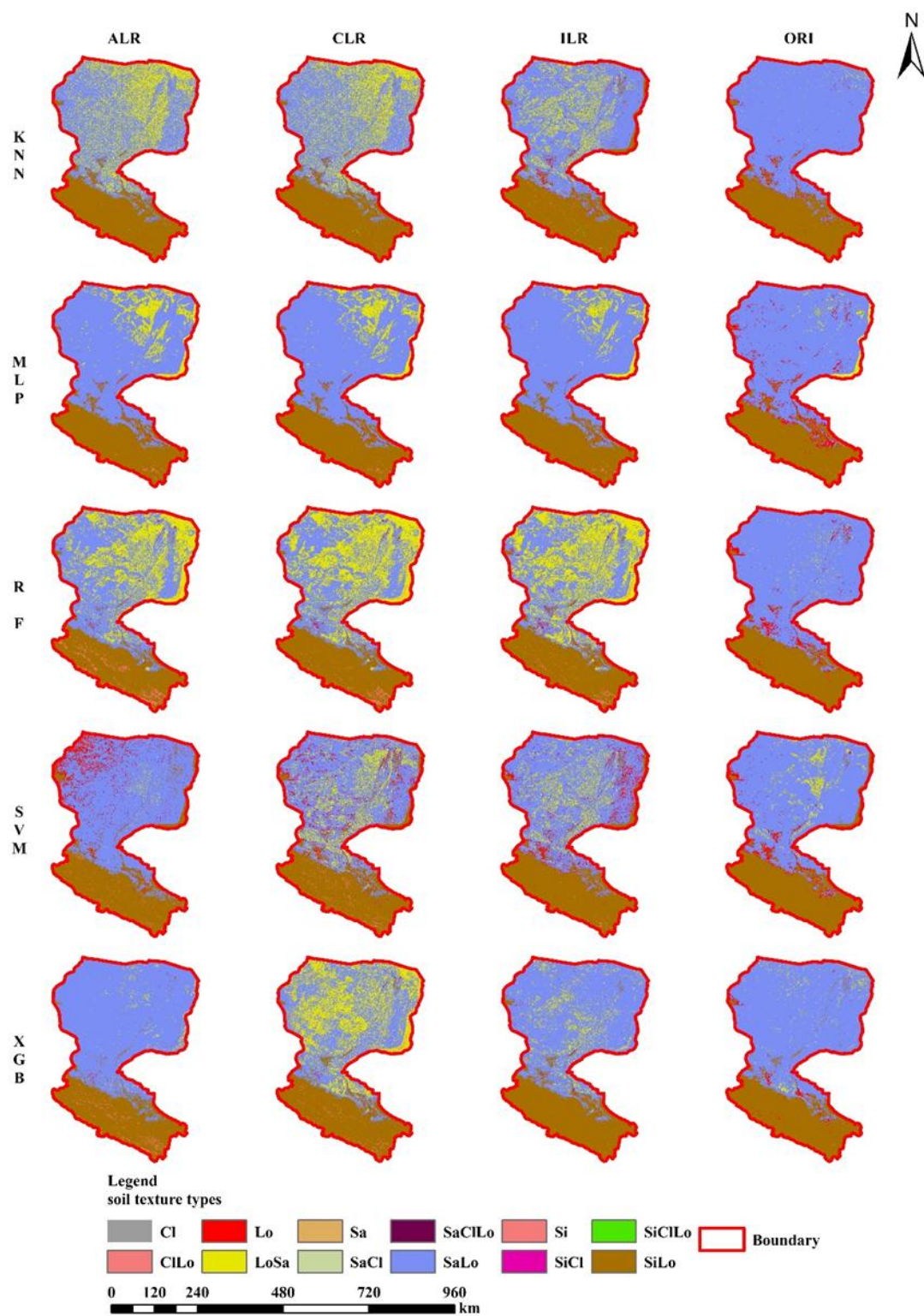


Figure 9. Soil texture classification prediction maps by soil psf interpolation (ALR, CLR, ILR log-ratio transformation methods and the original method) of KNN, MLP, RF, SVM and XGB.

3.4.4 Comparison of time-spending for each model in soil texture classification and soil psf interpolation

Time spending for models was computed to compare the efficiency of different machine-learning models in soil texture classification and soil psf interpolation (Fig. 10). Because the differences in time ~~spending-spent~~ among ORI and log ratio ~~approachesmethods~~ were similar, time spent of ILR was selected for soil psf interpolation. For the different models, RFs required the longest time for both classification (453.73 s) and regression (188.87 s), which may cause it to lose advantages when dealing with big data sets. KNN (classification: 4.2 s, regression: 23.6 s) and SVM (classification: 4.15 s, regression: 12.4 s) both showed shorter time in not only classification but also regression. Likewise, XGB (classification: 21.6 s, regression: 17.13 s) was much more stable and used less time, and the data processes were simpler compared with MLP (classification: 47.28 s, regression: 152.31 s). Moreover, XGB# delivered better performance than KNN and SVM in prediction maps of HRB, demonstrating an effective way of dealing with larger data.

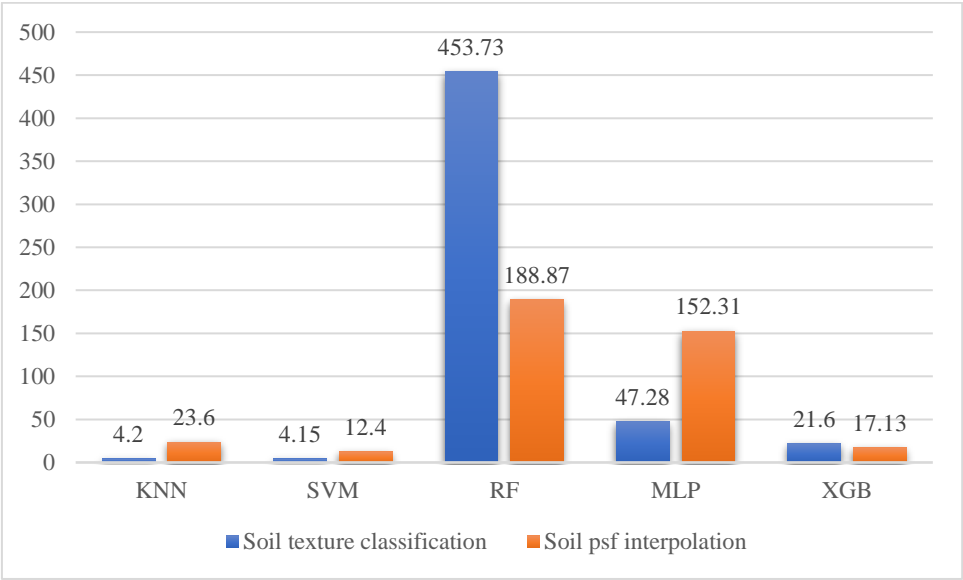


Figure 10. Average time spent running 30 times for KNN, MLP, RF, SVM and XGB of soil texture classification and soil psf interpolation.

4 Discussion

4.1 The systematic comparison of the five machine learning models

As mentioned previously, we compared the performance of different machine-learning methods containing KNN, MLP, RF, SVM and XGB. The results ~~demonstrate~~ demonstrated that SVM had the highest overall accuracy and XGB generated the highest kappa coefficient with respect to direct soil texture classification; considering the comprehensive evaluation of AUC and AUPRC, RF showed the best performance among these machine-learning models. In the case of soil psf interpolation, the indicators of RMSE, MAE, R^2 , AD and STRESS showed that RFs outperformed other machine-learning models, which also indicated additional information in prediction maps of sand, silt and clay as well as ~~as and~~ models of XGB. For the indirection classification of soil texture, the USDA soil texture triangles generated from RF were the closest to the distribution of the original data (Fig. 8a), with the highest classification right ratio. Prediction maps of indirect soil texture classification were also considered, ~~demonstrating~~ moreover, RF and MLP models were more suitable for the real environment, especially the models combined with log ratio transformation approaches ~~methods~~. Time ~~spending~~ spent of different machine-learning models showed that KNN, SVM and XGB required less time than RF and MLP to fit large data sets.

The comparisons of machine-learning models were also mentioned in previous reports. Heung et al. (2016) demonstrated that tree learners, such as RFs, delivered better performance than KNN and SVM due to the advantages of the interpretability of the results for classification problems in soil science; tree learners (decision trees) were also shown by Taghizadeh-Mehrjardi et al. (2015), indicating that the decision trees and ANN outperformed KNN, RF and SVM. ANNs, however, were typically complicated, which was true for our study due to the standardization and back transformation of MLP. In contrast, Wu et al. (2018) proposed that SVM revealed reliable consequences in direct soil texture classification, which was quite different from our results. In general, as binary classifiers, multi-class tasks can be handled as well using SVM; however, this is no longer the case in our study, as only 2 types of soil texture were generated from SVM, showing unsatisfactory results in both kappa coefficients and prediction maps. The consequences may be explained by the imbalanced data of soil texture types. For more information about tree learners in soil science for regression, Hengl et al. (2017) found lower R^2 using XGB than RF on a global-scale prediction. Zeraatpisheh et al. (2018) put forward the lowest RMSE and the highest R^2 using RF compared with multiple linear regression and regression trees for the prediction of clay, and this conclusion was similar to our study. For the total computing time, RF revealed the longest time with respect to both classification (453.73 s) and regression (188.87 s); however, it is the most accurate among five machine-learning methods in our study. In addition, for trade-offs of the total computing time of model and sub-optimal accuracy, XGB was superior to any other model, reducing the computing time significantly, while maintaining the accuracy not drop too much. With respect to the generality results of a transition of these machine-learning methods to other areas, it can be considered hierarchically. First, for the quick and imprecise machine-learning methods, XGB was recommended (sub-optimal accuracy), which was fast at the expense of a loss in precision. Second, considering the precise methods, RF can deliver the most accurate results, but it takes the longest computing time. Therefore,

XGB should be selected when researchers deal with larger data sets and regional scale study area; if they have enough time while want to produce more accurate results, RF is recommended.

4.2 The systematic comparison of the models combined with three log-ratio transformed data and original data

We compared the performance of models combined with three types of log ratio transformed data and original (untransformed) data for soil psf interpolation and indirect soil texture classification, and the results showed that the models using original data performed better in the case of indicators, such as RMSE, MAE, R^2 and AD, while the models using log ratio transformed data improved the STRESS. The interpolation maps of soil psf using the ILR method illustrated closer ranges of soil sampling data than those based on the ORI method. With respect to the indirect soil texture classification, models using log ratio transformed data improved the overall ~~aeuracies~~accuracy and kappa coefficients, such as RF and XGB. The USDA soil texture triangles showed more discrete distribution and more accordance with soil sampling data using the ILR transformation method. Better performance was shown in soil texture classification prediction maps generated from log ratio transformed data. Among the three log ratio ~~approaches~~methods, ILR and CLR were superior to ALR for the reason of more accurate indicators of soil psf interpolation and indirect soil texture classification, as well the performance of prediction maps. Additionally, log ratio ~~approaches~~methods modified soil sampling data to become more symmetric (Filzmoser et al., 2009); however, this improvement was not greatly effective. Fig.2 illustrated that soil sampling data for sand and clay were right-skewed, and silt was left-skewed because the silt component was predominant. The ALR transformed method enhanced soil sampling data of sand; nevertheless, the ALR_sand was still right-skewed, similar to the CLR_sand, presenting the lack of adjustment. In contrast, the ILR_sand changed from right-skewed to left-skewed; from this point of view, the over-adjustment was revealed. Similarly, the lack of adjustments was also shown in CLR_silt and ILR_silt; over-adjustments included ALR_silt, ALR_clay, CLR_clay and ILR_clay, making images that were different from normal distribution, and the p values of k-s tests were not significant. In our previous research (Wang and Shi, 2017), the ILR method had better performance than ALR and CLR, with the highest R^2 and lowest AD. The CLR method also performed well due to the lowest RMSE and mean error (ME) among the three log ratio ~~approaches~~methods. When comparing the original (untransformed) and log ratio ~~approaches~~methods, kriging approaches based on the log ratio delivered slightly decreased ~~aeuracies~~accuracy, which was similar to the conclusion in our study.

4.3 The systematic comparison of the direct and indirect classification for soil psf

Indirect classification showed not only better performance with respect to accuracy evaluation but also more accordance with the real environment than direct classification. The highest kappa coefficient generated from indirect classification (RF_ILR: 0.291) demonstrated obvious improvement (approximately 21.3 %) compared with that of direct classification (XGB: 0.240), keeping the highest overall accuracy stable (-1.4 %) at the same time (direct: 0.647; indirect: 0.638, respectively).

Compared with the real soil texture distribution and environment of the HRB, SiLo overlaid the upper reaches of HRB, and SaLo and Lo were in the south of the upper reaches of HRB showed strip distribution. Moreover, an uncovered area was detected in the northwest of the lower reaches of HRB, where it cannot be predicted due to a lack of information (soil samples) input in the process of model training. The main soil texture types of the lower reaches of the HRB were SiLo, LoSa and small amounts of SaLo and Lo distributed in uncovered area. The main soil texture types predicted by direct classification using machine-learning models were SaLo and SiLo; RF and XGB delivered much more LoSa than other direct classification models. However, all these models predicted that the main soil type of the lower reaches of HRB was SaLo, which was not fitted for the real environment (LoSa). In fact, LoSa and SaLo were obviously most confused classes; however, they are fairly similar to each other (see Fig. 8). In addition, because of the limitation of the train sets, direct classification can only predict types in the training data; in contrast, indirect classification broke such limitations, and new prediction types arose due to the transformation from soil psf to soil texture types. Moreover, more suitable matching performance with the real environment should be considered, such as the log ratio [approachesmethods](#) of MLP, RF, KNN_ ALR, KNN_ ILR and XGB_CLR. The direct soil texture classification generated relative unsatisfactory consequences. Although the indirect soil texture classification outperformed the direct one, kappa coefficients for indirect classification at fair-level (0.21-0.40) also need to be enhanced. Hence, soil sampling data appear to be comprehensively meaningful, considering accuracy improvement. In the case of soil sampling data, the laser diffraction approach we mentioned above was applied to obtain the discrete representation of particle size curves based on the given quantiles of these curves, i.e., soil particle size fractions (psf, sand, silt and clay). Subsequently, soil psf data were separately modeled for prediction and validation. Another perspective of soil psf should be considered, i.e., the probability density functions of particle size curves (so-called functional compositions), which are non-negative values that integrate to 1 (or 100 %) and can be considered as compositional data with infinitesimal parts (Menafoglio et al., 2014). Unlike conventional approaches, the viewpoints of functional compositions are beneficial to acquiring complete and continuous information rather than discrete information (sand, silt and clay) and soil texture and soil particle size fractions can be extracted using the stochastic simulation of soil particle-size curves (Menafoglio et al., 2016b). Previous studies applied such functional-compositional data for the simulation of particle size curves combined with geostatistical or machine-learning methods such as kriging and bayes [approachesmethods](#) (Menafoglio et al., 2016a) in hydrogeology, demonstrating more remarkable results compared with traditional methods. Therefore, which data should be used is the key points of accuracy improvement in future research.

5 Conclusion

We systematically compared a total of 45 models for direct and indirect soil texture classification, and soil psf interpolation using five machine-learning [approachesmethods](#) combined with original (untransformed) and three different log ratio transformed data in the HRB. The results indicate-indicated that as flexible and stable models, tree learners such as RF delivered powerful performance in both classification and regression and were superior to other machine-learning models mentioned

above. As a new and sub-optimal machine-learning method in soil science, XGB appeared to be more meaningful and more computationally efficient when dealing with large data sets. In addition, the log ratio [approachesmethods](#) had advantages of modifying STRESS in soil psf interpolation. Moreover, the indirect soil texture classification outperformed the direct one, especially when combined with the log ratio [approachesmethods](#). The indirect soil texture classification generated preferable consequences in both cases of accuracy indicators and prediction maps. More appropriate environmental covariates and interpolation techniques, more symmetric distribution of soil sampling data (or multiple perspectives of compositional data selection), and systematic parameter adjustment algorithms of compositional data are key to improving accuracy in the future.

Data availability. The soil sampling data (<http://westdc.westgis.ac.cn/DOI:10.3972/heihe.009.2013.db>; <http://westdc.westgis.ac.cn/DOI:10.3972/heihe.009.2013.db>; <http://westdc.westgis.ac.cn/DOI:10.3972/heihe.00135.2016.db>; <http://westdc.westgis.ac.cn/DOI:10.3972/hiwater.147.2013.db>; <http://westdc.westgis.ac.cn/DOI:10.3972/heihe.037.2014.db>; <http://westdc.westgis.ac.cn/DOI:doi:10.3972/heihe.0034.2013.db>; <http://westdc.westgis.ac.cn/DOI:10.3972/heihe.093.2013.db>) and part of environmental covariates data can be accessed through <http://westdc.westgis.ac.cn/> (last access: 29 October 2018). The meteorological data can be accessed through <http://data.cma.cn/> (last access: 29 October 2018).

Author contributions. WS contributed to soil data sampling, oversaw the design of the entire project. MZ performed the analysis and wrote the manuscript. Both authors contributed to writing this paper and interpreting data.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. This study was supported by the National Natural Science Foundation of China (Grant No. 41771111 and 41771364), Fund for Excellent Young Talents in Institute of Geographic Sciences and Natural Resources Research, Chinese Academy of Sciences (2016RC201), and the Youth Innovation Promotion Association, CAS (No. 2018071).

References

Abdi, D., Cade-Menun, B. J., Ziadi, N., and Parent, L.-E.: Compositional statistical analysis of soil p-31-nmr forms, *Geoderma*, 257, 40-47, <https://doi.org/10.1016/j.geoderma.2015.03.019>, 2015.

- Adhikari, K., and Hartemink, A. E.: Linking soils to ecosystem services - a global review, *Geoderma*, 262, 101-111, <https://doi.org/10.1016/j.geoderma.2015.08.009>, 2016.
- Aitchison, J.: The statistical-analysis of compositional data, *Journal of the Royal Statistical Society Series B-Methodological*, 44, 139-177, 1982.
- 5 Aitchison, J.: On criteria for measures of compositional difference, *Mathematical Geology*, 24, 365-379, <https://doi.org/10.1007/bf00891269>, 1992.
- Bacon-Shone, J. H.: A short history of compositional data analysis, in: *Compositional data analysis: Theory and applications*, Wiley, Chichester, West Sussex, 3, 2011.
- Bagheri Bodaghabadi, M., Antonio Martinez-Casasnovas, J., Salehi, M. H., Mohammadi, J., Esfandiarpour Borujeni, I.,
 10 Toomanian, N., and Gandomkar, A.: Digital soil mapping using artificial neural networks and terrain-related attributes, *Pedosphere*, 25, 580-591, [https://doi.org/10.1016/s1002-0160\(15\)30038-2](https://doi.org/10.1016/s1002-0160(15)30038-2), 2015.
- Basheer, I. A., and Hajmeer, M.: Artificial neural networks: fundamentals, computing, design, and application, *Journal of Microbiological Methods*, 43, 3-31, [https://doi.org/10.1016/S0167-7012\(00\)00201-3](https://doi.org/10.1016/S0167-7012(00)00201-3), 2000.
- Bationo, A., Kihara, J., Vanlauwe, B., Waswa, B., and Kimetu, J.: Soil organic carbon dynamics, functions and management
 15 in west african agro-ecosystems, *Agricultural Systems*, 94, 13-25, <https://doi.org/10.1016/j.agsy.2005.08.011>, 2007.
- Behrens, T., and Scholten, T.: Chapter 25 a comparison of data-mining techniques in predictive soil mapping, in: *Developments in soil science*, edited by: Lagacherie, P., McBratney, A. B., and Voltz, M., Elsevier, 353-617, 2006.
- Bergmeir, C., and Benitez, J. M.: Neural networks in R using the stuttgart neural network simulator: RSNNS, *Journal of Statistical Software*, 46, 1-26, 2012.
- 20 Breiman, L.: Bagging predictors, *Machine Learning*, 24, 123-140, <https://doi.org/10.1023/a:1018054314350>, 1996.
- Breiman, L.: Random forests, *Machine Learning*, 45, 5-32, <https://doi.org/10.1023/a:1010933404324>, 2001.
- Brown, D. J., Clayton, M. K., and McSweeney, K.: Potential terrain controls on soil color, texture contrast and grain-size deposition for the original catena landscape in uganda, *Geoderma*, 122, 51-72, <https://doi.org/10.1016/j.geoderma.2003.12.004>, 2004.
- 25 Brus, D. J., Kempen, B., and Heuvelink, G. B. M.: Sampling for validation of digital soil maps, *European Journal of Soil Science*, 62, 394-407, <https://doi.org/10.1111/j.1365-2389.2011.01364.x>, 2011.
- Burges, C. J. C.: A tutorial on support vector machines for pattern recognition, *Data Mining and Knowledge Discovery*, 2, 121-167, <https://doi.org/10.1023/a:1009715923555>, 1998.
- Camera, C., Zomeni, Z., Noller, J. S., Zissimos, A. M., Christoforou, I. C., and Bruggeman, A.: A high resolution map of soil
 30 types and physical properties for cyprus: A digital soil mapping optimization, *Geoderma*, 285, 35-49, <https://doi.org/10.1016/j.geoderma.2016.09.019>, 2017.
- Chen, T., and Guestrin, C.: Xgboost: A scalable tree boosting system, *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Francisco, California, USA, 2016.

- Chen, T., He, T., Benesty, M., Khotilovich, V., Tang, Y., Cho, H., Chen, K., Mitchell, R., Cano, I., Zhou, T., Li, M., Xie, J., Lin, M., Geng, Y., and Li, Y.: Xgboost: Extreme gradient boosting, R package version 0.71.2, available at: <https://CRAN.R-project.org/package=xgboost> (last access: 18 November 2018), 2018.
- Conrad, O., Bechtel, B., Bock, M., Dietrich, H., Fischer, E., Gerlitz, L., Wehberg, J., Wichmann, V., and Böhner, J.: System for automated geoscientific analyses (saga) v. 2.1.4, *Geosci. Model Dev.*, 8, 1991-2007, <https://doi.org/10.5194/gmd-8-1991-2015>, 2015.
- Cortes, C., and Vapnik, V.: Support-vector networks, *Machine Learning*, 20, 273-297, <https://doi.org/10.1023/a:1022627411411>, 1995.
- Cover, T. M., and Hart, P. E.: Nearest neighbor pattern classification, *Ieee Transactions on Information Theory*, 13, 21, <https://doi.org/10.1109/tit.1967.1053964>, 1967.
- Crouvi, O., Pelletier, J. D., and Rasmussen, C.: Predicting the thickness and aeolian fraction of soils in upland watersheds of the mojave desert, *Geoderma*, 195, 94-110, <https://doi.org/10.1016/j.geoderma.2012.11.015>, 2013.
- Davis, J., and Goadrich, M.: The relationship between precision-recall and roc curves, *Proceedings of the 23rd international conference on Machine learning*, Pittsburgh, Pennsylvania, USA, 2006.
- Egozcue, J. J., Pawlowsky-Glahn, V., Mateu-Figueras, G., and Barcelo-Vidal, C.: Isometric logratio transformations for compositional data analysis, *Mathematical Geology*, 35, 279-300, <https://doi.org/10.1023/a:1023818214614>, 2003.
- Elith, J., Leathwick, J. R., and Hastie, T.: A working guide to boosted regression trees, *Journal of Animal Ecology*, 77, 802-813, <https://doi.org/10.1111/j.1365-2656.2008.01390.x>, 2008.
- Fan, J., Wang, X., Wu, L., Zhou, H., Zhang, F., Yu, X., Lu, X., and Xiang, Y.: Comparison of support vector machine and extreme gradient boosting for predicting daily global solar radiation using temperature and precipitation in humid subtropical climates: A case study in china, *Energy Conversion and Management*, 164, 102-111, <https://doi.org/10.1016/j.enconman.2018.02.087>, 2018.
- Fawcett, T.: An introduction to roc analysis, *Pattern Recognition Letters*, 27, 861-874, <https://doi.org/10.1016/j.patrec.2005.10.010>, 2006.
- Filzmoser, P., Hron, K., and Reimann, C.: Univariate statistical analysis of environmental (compositional) data: Problems and possibilities, *Science of the Total Environment*, 407, 6100-6108, <https://doi.org/10.1016/j.scitotenv.2009.08.008>, 2009.
- Follain, S., Minasny, B., McBratney, A. B., and Walter, C.: Simulation of soil thickness evolution in a complex agricultural landscape at fine spatial and temporal scales, *Geoderma*, 133, 71-86, <https://doi.org/10.1016/j.geoderma.2006.03.038>, 2006.
- Fu, G.-H., Xu, F., Zhang, B.-Y., and Yi, L.-Z.: Stable variable selection of class-imbalanced data with precision-recall criterion, *Chemometrics and Intelligent Laboratory Systems*, 171, 241-250, <https://doi.org/10.1016/j.chemolab.2017.10.015>, 2017.
- Gobin, A., Campling, P., and Feyen, J.: Soil-landscape modelling to quantify spatial variability of soil texture, *Physics and Chemistry of the Earth Part B-Hydrology Oceans and Atmosphere*, 26, 41-45, [https://doi.org/10.1016/s1464-1909\(01\)85012-7](https://doi.org/10.1016/s1464-1909(01)85012-7), 2001.

- Gochis, D. J., Vivoni, E. R., and Watts, C. J.: The impact of soil depth on land surface energy and water fluxes in the north american monsoon region, *Journal of Arid Environments*, 74, 564-571, <https://doi.org/10.1016/j.jaridenv.2009.11.002>, 2010.
- Hechenbichler, K., and Schliep, K.: Weighted k-Nearest-Neighbor Techniques and Ordinal Classification, 2004.
- 5 Hengl, T., Heuvelink, G. B. M., Kempen, B., Leenaars, J. G. B., Walsh, M. G., Shepherd, K. D., Sila, A., MacMillan, R. A., de Jesus, J. M., Tamene, L., and Tondoh, J. E.: Mapping soil properties of africa at 250 m resolution: Random forests significantly improve current predictions, *Plos One*, 10, <https://doi.org/10.1371/journal.pone.0125814>, 2015.
- Hengl, T., de Jesus, J. M., Heuvelink, G. B. M., Gonzalez, M. R., Kilibarda, M., Blagotic, A., Shangguan, W., Wright, M. N., Geng, X., Bauer-Marschallinger, B., Guevara, M. A., Vargas, R., MacMillan, R. A., Batjes, N. H., Leenaars, J. G. B.,
- 10 Ribeiro, E., Wheeler, I., Mantel, S., and Kempen, B.: Soilgrids250m: Global gridded soil information based on machine learning, *Plos One*, 12, <https://doi.org/10.1371/journal.pone.0169748>, 2017.
- Heung, B., Ho, H. C., Zhang, J., Knudby, A., Bulmer, C. E., and Schmidt, M. G.: An overview and comparison of machine-learning techniques for classification purposes in digital soil mapping, *Geoderma*, 265, 62-77, <https://doi.org/10.1016/j.geoderma.2015.11.014>, 2016.
- 15 Huang, J., Subasinghe, R., and Triantafyllis, J.: Mapping particle-size fractions as a composition using additive log-ratio transformation and ancillary data, *Soil Science Society of America Journal*, 78, 1967-1976, <https://doi.org/10.2136/sssaj2014.05.0215>, 2014.
- Huete, A., Didan, K., Miura, T., Rodriguez, E. P., Gao, X., and Ferreira, L. G.: Overview of the radiometric and biophysical performance of the modis vegetation indices, *Remote Sensing of Environment*, 83, 195-213,
- 20 [https://doi.org/10.1016/s0034-4257\(02\)00096-2](https://doi.org/10.1016/s0034-4257(02)00096-2), 2002.
- Huete, A. R.: A soil-adjusted vegetation index (savi), *Remote Sensing of Environment*, 25, 295-309, [https://doi.org/10.1016/0034-4257\(88\)90106-x](https://doi.org/10.1016/0034-4257(88)90106-x), 1988.
- Jafari, A., Khademi, H., Finke, P. A., Van de Wauw, J., and Ayoubi, S.: Spatial prediction of soil great groups by boosted regression trees using a limited point dataset in an arid region, southeastern iran, *Geoderma*, 232, 148-163,
- 25 <https://doi.org/10.1016/j.geoderma.2014.04.029>, 2014.
- Krasilnikov, P. V., Garcia-Calderon, N. E., Ibanez-Huerta, A., Bazan-Mateos, M., and Hernandez-Santana, J. R.: Soils in the dynamic tropical environments: The case of sierra madre del sur, *Geomorphology*, 135, 262-270, <https://doi.org/10.1016/j.geomorph.2011.02.013>, 2011.
- Kuhn, M.: Caret: Classification and regression training, R package version 6.0-80, available at: <https://CRAN.R-project.org/package=caret> (last access: 18 November 2018), 2018.
- 30 Landis, J. R., and Koch, G. G.: Measurement of observer agreement for categorical data, *Biometrics*, 33, 159-174, <https://doi.org/10.2307/2529310>, 1977.
- Lark, R. M., and Bishop, T. F. A.: Cokriging particle size fractions of the soil, *European Journal of Soil Science*, 58, 763-774, <https://doi.org/10.1111/j.1365-2389.2006.00866.x>, 2007.

- Liaw, A., and Wiener, M.: Classification and regression by randomforest, R News, 2, 18-22, 2002.
- Liess, M., Glaser, B., and Huwe, B.: Uncertainty in the spatial prediction of soil texture comparison of regression tree and random forest models, Geoderma, 170, 70-79, <https://doi.org/10.1016/j.geoderma.2011.10.010>, 2012.
- Marchetti, A., Piccini, C., Santucci, S., Chiuchiarelli, I., and Francaviglia, R.: Simulation of soil types in teramo province (central italy) with terrain parameters and remote sensing data, Catena, 85, 267-273, <https://doi.org/10.1016/j.catena.2011.01.012>, 2011.
- Martin-Fernandez, J. A., Olea-Meneses, R. A., and Pawlowsky-Glahn, V.: Criteria to compare estimation methods of regionalized compositions, Mathematical Geology, 33, 889-909, <https://doi.org/10.1023/a:1012293922142>, 2001.
- McBratney, A. B., Santos, M. L. M., and Minasny, B.: On digital soil mapping, Geoderma, 117, 3-52, [https://doi.org/10.1016/s0016-7061\(03\)00223-4](https://doi.org/10.1016/s0016-7061(03)00223-4), 2003.
- McNamara, J. P., Chandler, D., Seyfried, M., and Achet, S.: Soil moisture states, lateral flow, and streamflow generation in a semi-arid, snowmelt-driven catchment, Hydrological Processes, 19, 4023-4038, <https://doi.org/10.1002/hyp.5869>, 2005.
- Menafoglio, A., Guadagnini, A., and Secchi, P.: A kriging approach based on aitchison geometry for the characterization of particle-size curves in heterogeneous aquifers, Stoch. Environ. Res. Risk Assess., 28, 1835-1851, <https://doi.org/10.1007/s00477-014-0849-8>, 2014.
- Menafoglio, A., Guadagnini, A., and Secchi, P.: Stochastic simulation of soil particle-size curves in heterogeneous aquifer systems through a bayes space approach, Water Resources Research, 52, 5708-5726, 10.1002/2015wr018369, 2016a.
- Menafoglio, A., Secchi, P., and Guadagnini, A.: A class-kriging predictor for functional compositions with application to particle-size curves in heterogeneous aquifers, Math Geosci., 48, 463-485, <https://doi.org/10.1007/s11004-015-9625-7>, 2016b.
- Metternicht, G. I., and Zinck, J. A.: Remote sensing of soil salinity: Potentials and constraints, Remote Sensing of Environment, 85, 1-20, [https://doi.org/10.1016/s0034-4257\(02\)00188-8](https://doi.org/10.1016/s0034-4257(02)00188-8), 2003.
- Meyer, D., Dimitriadou, E., Hornik, K., Andreas, W., and Friedrich, L.: E1071: Misc functions of the department of statistics, probability theory group (formerly: E1071), tu wien, R package version 1.6-8, available at: <https://CRAN.R-project.org/package=e1071> (last access: 18 November 2018), 2017.
- Moeys, J.: Soiltexture: Functions for soil texture plot, classification and transformation, R package version 1.4.6, available at: <https://CRAN.R-project.org/package=soiltexture> (last access: 18 November 2018), 2018.
- Odeh, I. O. A., Todd, A. J., and Triantafyllis, J.: Spatial prediction of soil particle-size fractions as compositional data, Soil Science, 168, 501-515, <https://doi.org/10.1097/00010694-200307000-00005>, 2003.
- Pahlavan-Rad, M. R., and Akbarimoghaddam, A.: Spatial variability of soil texture fractions and ph in a flood plain (case study from eastern iran), Catena, 160, 275-281, <https://doi.org/10.1016/j.catena.2017.10.002>, 2018.
- Poggio, L., and Gimona, A.: 3d mapping of soil texture in scotland, Geoderma Regional, 9, 5-16, <https://doi.org/10.1016/j.geodrs.2016.11.003>, 2017.

- R Core Team: R: A language and environment for statistical computing, R Foundation for Statistical Computing, Vienna, Austria. available at: <https://www.R-project.org/> (last access: 18 November 2018), 2018.
- 2018.
- Saito, T., and Rehmsmeier, M.: Precrec: Fast and accurate precision-recall and roc curve calculations in r, *Bioinformatics*, 33, 145-147, <https://doi.org/10.1093/bioinformatics/btw570>, 2017.
- Salazar, E., Giraldo, R., and Porcu, E.: Spatial prediction for infinite-dimensional compositional data, *Stochastic Environmental Research and Risk Assessment*, 29, 1737-1749, <https://doi.org/10.1007/s00477-014-1010-4>, 2015.
- Schliep, K., and Hechenbichler, K.: Kknn: Weighted k-nearest neighbors, R package version 1.3.1, available at: <https://CRAN.R-project.org/package=kknn> (last access: 18 November 2018), 2016.
- 10 Song, X.-D., Brus, D. J., Liu, F., Li, D.-C., Zhao, Y.-G., Yang, J.-L., and Zhang, G.-L.: Mapping soil organic carbon content by geographically weighted regression: A case study in the heihe river basin, china, *Geoderma*, 261, 11-22, <https://doi.org/10.1016/j.geoderma.2015.06.024>, 2016.
- Subasi, A.: Eeg signal classification using wavelet feature extraction and a mixture of expert model, *Expert Systems with Applications*, 32, 1084-1093, <https://doi.org/10.1016/j.eswa.2006.02.005>, 2007.
- 15 Sun, X.-L., Wu, Y.-J., Wang, H.-L., Zhao, Y.-G., and Zhang, G.-L.: Mapping soil particle size fractions using compositional kriging, cokriging and additive log-ratio cokriging in two case studies, *Mathematical Geosciences*, 46, 429-443, <https://doi.org/10.1007/s11004-013-9512-z>, 2014.
- Taalab, K., Corstanje, R., Zawadzka, J., Mayr, T., Whelan, M. J., Hannam, J. A., and Creamer, R.: On the application of bayesian networks in digital soil mapping, *Geoderma*, 259, 134-148, <https://doi.org/10.1016/j.geoderma.2015.05.014>, 2015.
- 20 Taghizadeh-Mehrjardi, R., Nabiollahi, K., Minasny, B., and Triantafyllis, J.: Comparing data mining classifiers to predict spatial distribution of usda-family soil groups in baneh region, iran, *Geoderma*, 253, 67-77, <https://doi.org/10.1016/j.geoderma.2015.04.008>, 2015.
- Thompson, J. A., Roecker, S., Grunwald, S., and Owens, P. R.: Chapter 21 - digital soil mapping: Interactions with and applications for hydropedology, in: *Hydropedology*, edited by: Lin, H., Academic Press, Boston, 665-709, 2012.
- 25 van den Boogaart, K. G., and Tolosana-Delgado, R.: "Compositions": A unified r package to analyze compositional data, *Computers & Geosciences*, 34, 320-338, <https://doi.org/10.1016/j.cageo.2006.11.017>, 2008.
- Vapnik, V.: The support vector method of function estimation, *Nonlinear modeling: Advanced black-box techniques*, edited by: Suykens, J. A. K., and Vandewalle, J., 55-85 pp., 1998.
- 30 Walvoort, D. J. J., and de Gruijter, J. J.: Compositional kriging: A spatial interpolation method for compositional data, *Mathematical Geology*, 33, 951-966, <https://doi.org/10.1023/a:1012250107121>, 2001.
- Wang, Z., and Shi, W.: Mapping soil particle-size fractions: A comparison of compositional kriging and log-ratio kriging, *Journal of Hydrology*, 546, 526-541, <https://doi.org/10.1016/j.jhydrol.2017.01.029>, 2017.

- Wang, Z., and Shi, W.: Robust variogram estimation combined with isometric log-ratio transformation for improved accuracy of soil particle-size fraction mapping, *Geoderma*, 324, 56-66, <https://doi.org/10.1016/j.geoderma.2018.03.007>, 2018.
- Wu, B., Yan, N., Xiong, J., Bastiaanssen, W. G. M., Zhu, W., and Stein, A.: Validation of etwatch using field measurements at diverse landscapes: A case study in hai basin of china, *Journal of Hydrology*, 436, 67-80, <https://doi.org/10.1016/j.jhydrol.2012.02.043>, 2012.
- Wu, W., Li, A.-D., He, X.-H., Ma, R., Liu, H.-B., and Lv, J.-K.: A comparison of support vector machines, artificial neural network and classification tree for identifying soil texture classes in southwest china, *Computers and Electronics in Agriculture*, 144, 86-93, <https://doi.org/10.1016/j.compag.2017.11.037>, 2018.
- Yang, R.-M., Zhang, G.-L., Liu, F., Lu, Y.-Y., Yang, F., Yang, F., Yang, M., Zhao, Y.-G., and Li, D.-C.: Comparison of boosted regression tree and random forest models for mapping topsoil organic carbon concentration in an alpine ecosystem, *Ecological Indicators*, 60, 870-878, <https://doi.org/10.1016/j.ecolind.2015.08.036>, 2016.
- Yi, C., Li, D., Zhang, G., Zhao, Y., Yang, J., Liu, F., and Song, X.: Criteria for partition of soil thickness and case studies, *Acta Pedologica Sinica*, 52, 220-227, 2015.
- Yoo, K., Amundson, R., Heimsath, A. M., and Dietrich, W. E.: Spatial patterns of soil organic carbon on hillslopes: Integrating geomorphic processes and the biological c cycle, *Geoderma*, 130, 47-65, <https://doi.org/10.1016/j.geoderma.2005.01.008>, 2006.
- Zeraatpisheh, M., Ayoubi, S., Jafari, A., and Finke, P.: Comparing the efficiency of digital and conventional soil mapping to predict soil types in a semi-arid region in iran, *Geomorphology*, 285, 186-204, <https://doi.org/10.1016/j.geomorph.2017.02.015>, 2017.
- Zeraatpisheh, M., Ayoubi, S., Jafari, A., Tajik, S., and Finke, P.: Digital mapping of soil properties using multiple machine learning in a semi-arid region, central iran, *Geoderma*, <https://doi.org/10.1016/j.geoderma.2018.09.006>, 2018.
- Zhang, S.-w., Shen, C.-y., Chen, X.-y., Ye, H.-c., Huang, Y.-f., and Lai, S.: Spatial interpolation of soil texture using compositional kriging and regression kriging with consideration of the characteristics of compositional data and environment variables, *Journal of Integrative Agriculture*, 12, 1673-1683, [https://doi.org/10.1016/s2095-3119\(13\)60395-0](https://doi.org/10.1016/s2095-3119(13)60395-0), 2013.
- Zhao, Z., Chow, T. L., Rees, H. W., Yang, Q., Xing, Z., and Meng, F.-R.: Predict soil texture distributions using an artificial neural network model, *Computers and Electronics in Agriculture*, 65, 36-48, <https://doi.org/10.1016/j.compag.2008.07.008>, 2009.