



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. 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 was developed to 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 of original (untransformed) and log ratio approaches, 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—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 approaches 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 %), 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, 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. Our systematic comparison helps to elucidate the processing and selection of compositional data in spatial simulation.

1 Abbreviations: psf, soil particle-size fractions; HRB, Heihe River Basin; DSM, digital soil mapping; KNN, 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

1

1. In this total, a total

typo

[Tomislav Hengl]

2. modified the soil sampling...

unclear sentence, you mean maybe "decreased skewness of distributions"?

[Tomislav Hengl]

3. soil psf interpolation.

please add here 1-2 sentences explaining the study area and data set used (number of points, area etc)

[Tomislav Hengl]

4. showed 20 notable consequences.

vague; please rephrase using concrete stat measures

[Tomislav Hengl]

5. Our systematic comparison...

Instead of this sentence (too subjective and not necessary) I would add the main conclusions (is transforming fractions necessary or not? what would you recommend as the best strategy to map PSFs / texture classes?) and implications of this work (to other peoples work).

[Tomislav Hengl]



1 Introduction

Soil texture, classified by ranges of soil particle-size fractions (psf)¹, 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 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 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). 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, which can be obtained from digital elevation models (DEM), remote sensing data, and categorical or geomorphology maps (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

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.

1. (psf),
 better use capital letters
 "PSF" consistently
 [Tomislav Hengl]



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 accuracies 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-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 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 to simulate the spatial variation of soil psf were kriging methods (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 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

1. science have been suggested

You missed an important reference:
<https://link.springer.com/article/10.1007%2Fs11004-018-9769-3>

[Tomislav Hengl]

2. However, most studies using...

Two times however, see next comment.

[Tomislav Hengl]

3. However, few studies combined...

This repeats sentence in L17 on this page. I would recommend reducing the size of these paragraphs.

[Tomislav Hengl]



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-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 the direct soil texture classification, (ii) to evaluate the accuracies of different log ratio approaches and original (untransformed) method applied for soil psf from the perspective of compositional data using machine-learning models, and (iii) to estimate whether the accuracies 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.

1. direct soil texture classification...

please explain for all readers what is "direct soil texture classification"

[Tomislav Hengl]

2. -5-4 °C

make difference between the minus sign and endash sign e.g. "-5 ndash 4"

[Tomislav Hengl]

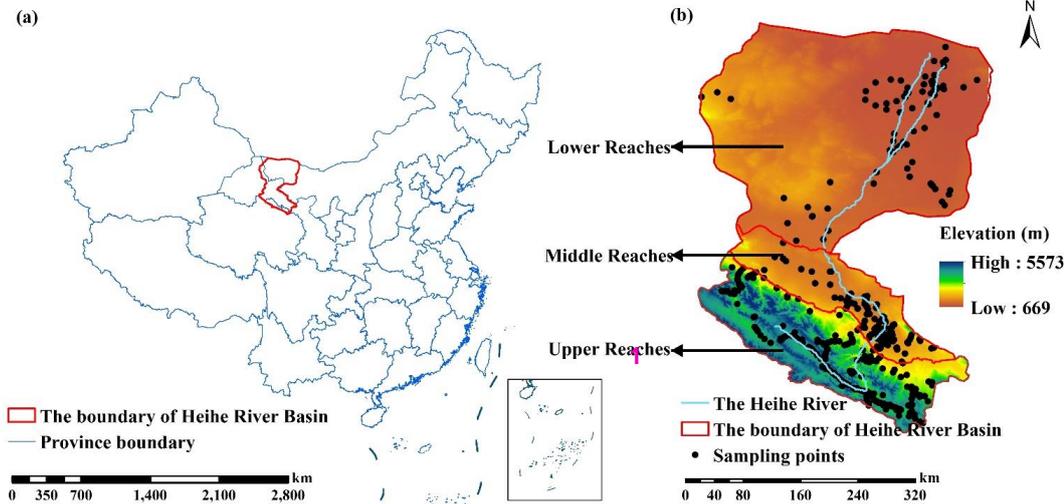


Figure 1. The (a) geographical location, (b) Heihe River, elevation and soil sampling points of Heihe River Basin, China.

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 (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 plimatic 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

1. These dots are beyond the cartographic scale of this visualization and should be removed. It politicizes the topic of soil texture mapping which is unnecessary. The topic of the paper is not the disputed islands but comparison of methods for PSF mapping.

[Tomislav Hengl]



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 topographic attributes from DEM, including slope, aspect, convergence index, 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 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 also address the distributions of soil psf. 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)

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-

1. position information.

this section needs to be extended - how was the sampling plan generated? did you sample per soil horizon or at fixed depths? did you take bulk samples or are the samples close-to-point support? very important that all these things are clarified.

[Tomislav Hengl]

2. Latitude and longitude...

I highly do not recommend using Lat and Lon as covariates. This has shown to lead to artifacts. For more details see:

<https://peerj.com/articles/5518/>

[Tomislav Hengl]

3. geomorphology types, soil...

Please provide more detail - classes etc. Which soil map you used as covariate layer and how? Has this come up as a significant predictor?

[Tomislav Hengl]



nearest neighbors. Weighted KNN is an extended version of KNN that considers the distances of the nearest neighbors; therefore, the parameters of KNN contain the maximum value of k (k_{max}), the distances of the nearest neighbors (distance) and the types of kernel function (kernel). The KNN model is available in the R package “*kkn*” (Schliep and Hechenbichler, 2016).

5 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. 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).

2.4.3 Random forest (RF)

15 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. 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 not necessary because it is insensitive to the range of value. Two parameters should be adjusted for RF model: the number of trees (n_{tree}) and the number of features randomly sampled at each split (m_{try}). The RF model is available in the R package “*randomForest*” (Liaw and Wiener, 2002).

2.4.4 Support vector machines (SVM)

25 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 (Burgess, 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. The advantages of SVMs are that they are effective in high dimensional

1. the most popular multilayer...

does Subasi (2007) claims that it is the most popular? Putting a reference here would be recommended.

[Tomislav Hengl]

2. “randomForest”

Unfortunately not optimized for larger data sets. I recommend using the “*ranger*” package instead.

[Tomislav Hengl]



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)

5 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
 10 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). 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.,
 15 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
 20 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
 25 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
 30

8

1. the other package of “randomF...
 which package?
 [Tomislav Hengl]



because the sum of the dimensions of CLR is 0, and thus the results are collinear. These problems can be overcome by using ILR approach, which transforms all the information into D-1 orthogonal log contrasts (Abdi et al., 2015). The transformation equations for ILR are defined as follows:

$$z = (z_1, \dots, z_{D-1}) = \text{ilr}(x), \quad (8)$$

$$z_i = \sqrt{\frac{D-i}{D-i+1}} \ln \frac{x_i}{\sqrt{\prod_{j=i+1}^D x_j}}, \quad (9)$$

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

$$\text{ilr}(1) = \sqrt{\frac{2}{3}} \ln \frac{\text{sand}}{\sqrt{\text{silt} \times \text{clay}}}, \quad (10)$$

$$\text{ilr}(2) = \sqrt{\frac{1}{2}} \ln \frac{\text{silt}}{\text{clay}}, \quad (11)$$

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 interpolation, only the first two equations were used for ALR and ILR:

$$\text{alr}(3) = \ln \frac{\text{clay}}{\text{sand}}, \quad (12)$$

$$\text{ilr}(3) = \sqrt{\frac{2}{3}} \ln \frac{\text{clay}}{\sqrt{\text{sand} \times \text{silt}}}, \quad (13)$$

The equations for $\text{alr}(1)$, $\text{alr}(2)$, $\text{alr}(3)$ were equivalent to $\text{alr}(\text{sand})$, $\text{alr}(\text{silt})$, $\text{alr}(\text{clay})$ 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).

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%:

$$\text{sand}_s = \frac{\text{sand}}{(\text{sand} + \text{silt} + \text{clay})} \times 100, \quad (14)$$

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

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 (ALR, CLR, ILR). Five machine-learning methods were applied for the direct soil texture classification; additionally, these methods were combined with original (untransformed) and log ratio transformed data for a total of 40 methods for the indirect soil texture classification (20) and soil psf interpolation (20). The data were randomly divided into two sets to guarantee prediction accuracy and persuasion; for instance, one (70 % = 448 soil samples)

1. 2008).

I think the readers would appreciate here is you would explain how are the ALR/ILR numbers back-transformed to original 0-100% scale?
 [Tomislav Hengl]



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):

$$\text{Overall Accuracy} = \frac{TP+TN}{TP+TN+FP+FN} \quad (15)$$

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. 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:

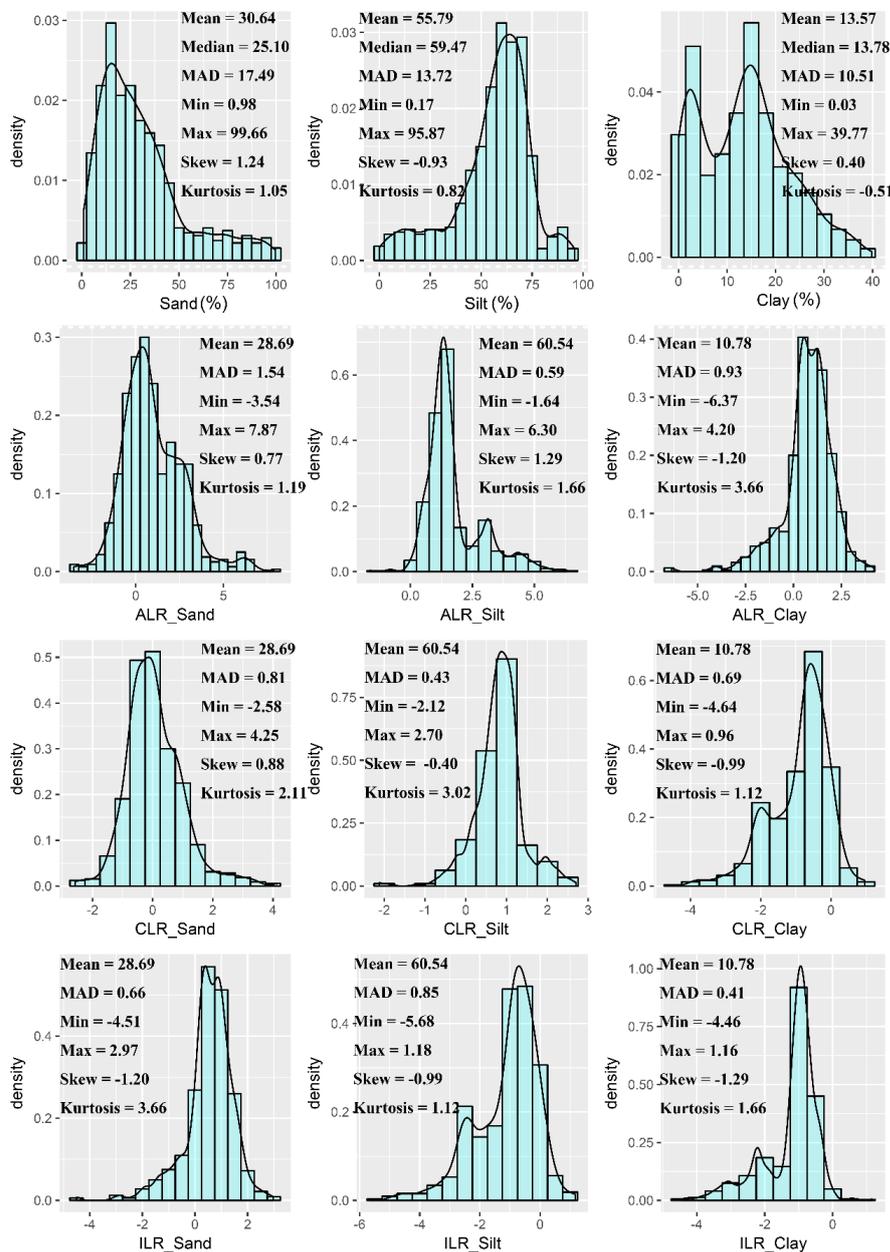
$$\text{kappa} = \frac{p_o - p_e}{1 - p_e} \quad (16)$$

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).

1. soil texture classification

which soil.texture system have you used - USDA? Are you aware of the "soil.texture" package?
 [Tomislav Hengl]

2. 0.01- 20 0.20: slight,...
 Kappa depends on number of samples and number of classes - please clarify.
 [Tomislav Hengl]



1

1. This figure could be much improved: (1) put exactly the same scale (0-100% and -8 to 8 for the x-axis) so the distributions can be compared, (2) center all ALR, CLR and ILR plots at 0, (3) remove "Min" and "Max" from the labels (visible from the plots).

[Tomislav Hengl]

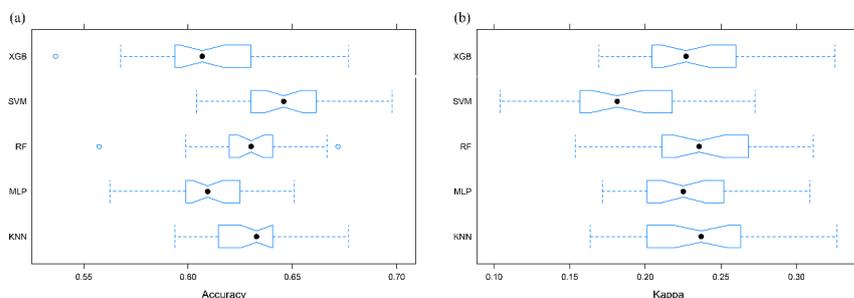


Figure 2. Descriptive statistical analysis for the original (untransformed) and logratio transformed soil sampling data. Note that the means generated from log ratio transformed data were back-transformed to the real space.

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

- 5 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 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).



- 10 **Figure 3.** (a) The overall accuracies 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:

15 RF>SVM>XGB>MLP>KNN.

1. XGB

plot shows "SVM" is the best? please check.
 [Tomislav Hengl]

2. Hmm - SVM is both best and worst performing method based on this plot - please check.

[Tomislav Hengl]



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.

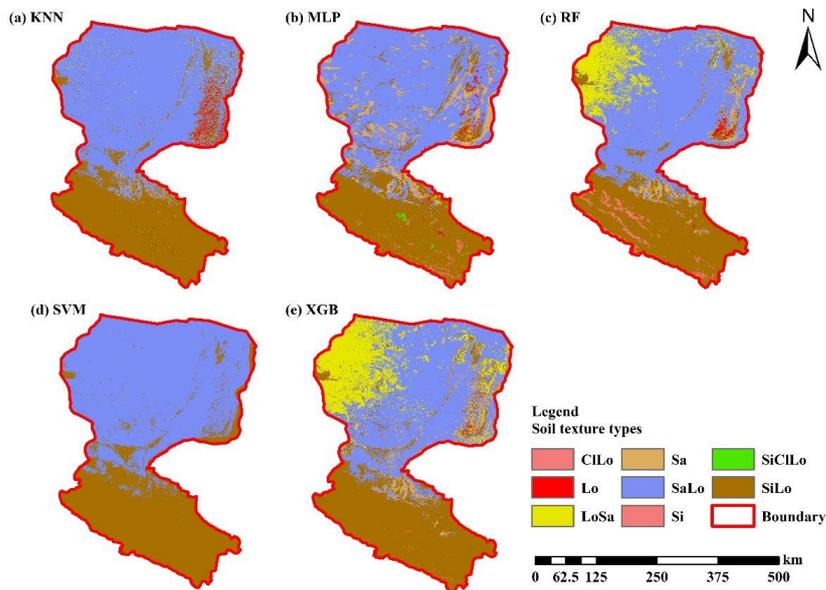


Figure 6. Soil texture classification prediction map of different soil texture types of (a) KNN, (b) MLP, (c) RF, (d) SVM and (e) XGB.

10 3.3 Comparison of the machine-learning models combined with log ratio approaches 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 of STRESS of the methods combined with log ratio transformed data were superior to other approaches using original (untransformed) data (Table 2). With respect to KNN, MLP, RF and XGB, values of RMSE, MAE, R^2 and AD generated from original (untransformed) data outperformed log ratio

15

19

1. (image annotation)

LoSa and SaLo are obviously most confused classes. But they are fairly similar to each other so not a big problem probably. To make that clear to readers I would use as legend the soil texture triangle as in <https://enviromatrix.github.io/PredictiveSoilMapping/soil-variables-chapter.html#converting-texture-by-hand-classes-to-fractions> (so that it is also visible which are the most similar classes).

[Tomislav Hengl]



Table 2. The comparisons of accuracies of different machine-learning models combined with original (untransformed) and transformed data. ¹

	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

1. The comparisons of accuracies...

I would prefer again box-plots as in Fig. 3. Or at least put in bold the best performing methods please.

[Tomislav Hengl]



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.

1. RF and XGB delivered more...

Not substantiated. You mean more contrast / wider output distributions? Please rephrase.

[Tomislav Hengl]



Figure 7. The interpolation maps of sand fraction. ¹

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 and kappa coefficients of models using log ratio transformed data were improved, especially RF and XGB, which combined with all three log ratio approaches were superior to the ORI approach. Table 3 shows 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 approach. In summary, the ILR transformation method of five machine-learning models showed the highest overall accuracy among three log ratio transformation approaches (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 the 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 the highest kappa coefficient among these methods.

Table 3. Overall accuracies and kappa coefficients calculated from soil texture classification by the interpolated maps from five models using original (untransformed) data and log ratio transformed data. ²

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

1. The interpolation maps...

which methods have systematically under-estimated / over-estimated sand content? which are most accurate? please add to the caption.

[Tomislav Hengl]

2. Overall accuracies and...

Overall very small differences in performance. Why is that?

[Tomislav Hengl]

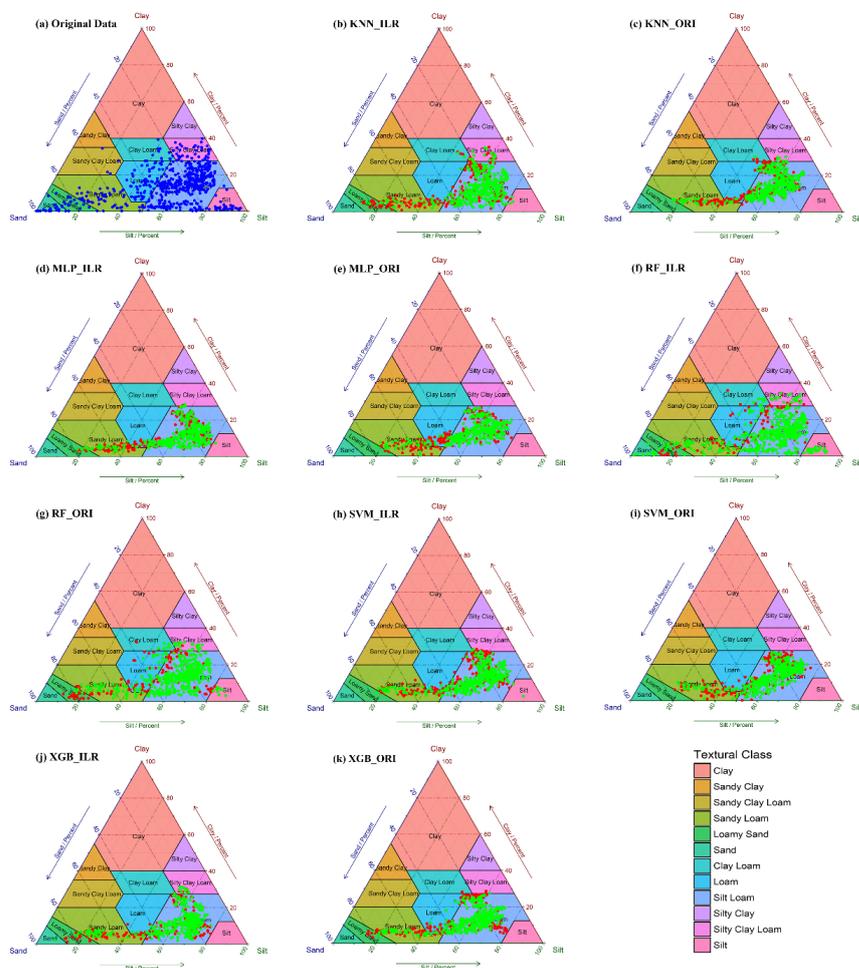
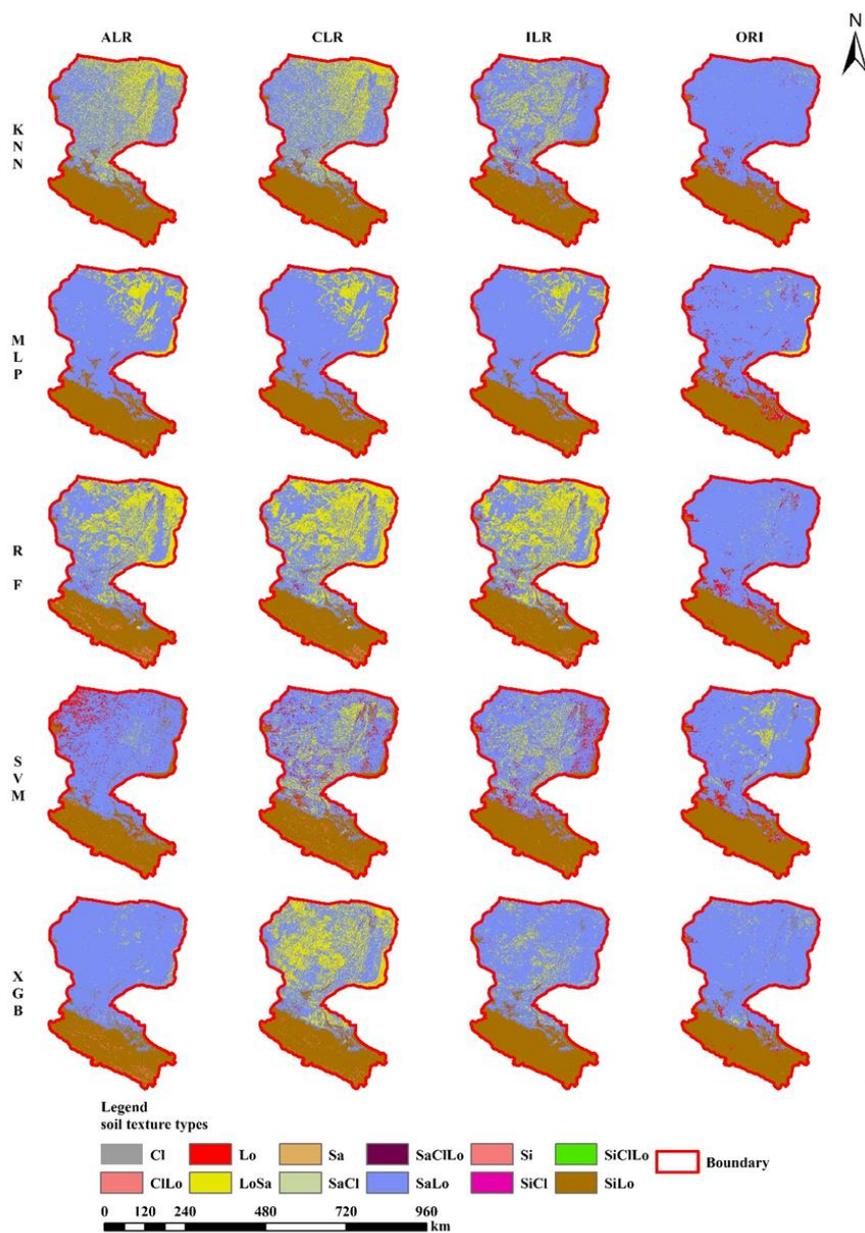


Figure 8. Soil texture types of 640 soil samples shown in USDA texture triangle. The results of soil psf were generated from (a) original (untransformed) data, (b) KNN_ILR (65.0 %), (c) KNN_ORI (65.9 %), (d) MLP_ILR (63.3 %), (e) MLP_ORI (63.6 %), (f) RF_ILR (83.9 %), (g) RF_ORI (81.7 %), (h) SVM_ILR (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

1. Figure 8. Soil

This is an important figure. But it could be much improved: (1) text is too small (unreadable), (2) it is not clear what are the green and red colored points, (3) to each plot add some summary measure of performance (so we know which was best / worst). I would also move the (a) original data and put it into Fig. 1.

[Tomislav Hengl]



1. see my comment on Fig. 6.
 [Tomislav Hengl]



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

3.4.4 Comparison of ¹time-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 among ORI and log ratio approaches were similar, 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, it delivered better performance than KNN and SVM in prediction maps of HRB, demonstrating an effective way of dealing with larger data.

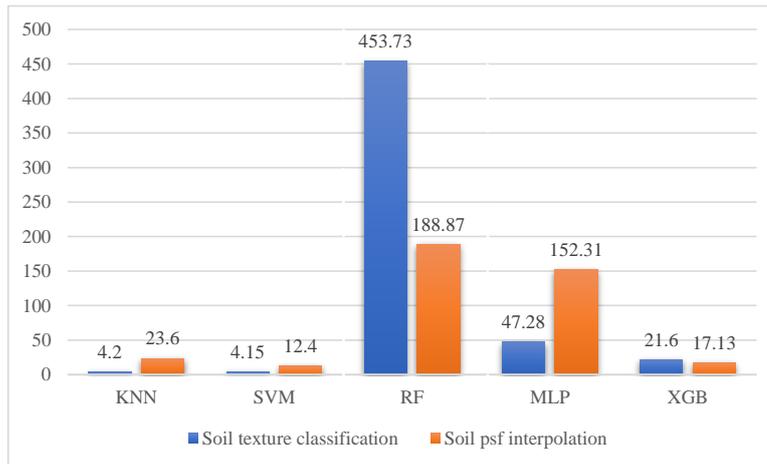


Figure 10. Average time spent running 30 times for KNN, MLP, RF, SVM and XGB of soil texture classification and soil psf interpolation.

1. time-spending
 total computing time
 [Tomislav Hengl]

2. mention how many trees
 you use? If you reduce the
 standard number of trees to
 e.g. 80 you do not lose
 much on accuracy but the
 computing time drops
 significantly.
 [Tomislav Hengl]



than those based on the ORI approach. With respect to the indirect soil texture classification, models using log ratio transformed data improved the overall accuracies 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 approach. Better performance was shown in soil texture classification prediction maps generated from log ratio transformed data. Among the three log ratio approaches, 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 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 approach 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 approach had better performance than ALR and CLR, with the highest R^2 and lowest AD. The CLR approach also performed well due to the lowest RMSE and mean error (ME) among the three log ratio approaches. When comparing the original (untransformed) and log ratio approaches, kriging approaches based on the log ratio delivered slightly decreased accuracies, 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 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

1. however, this improvement...
 please elaborate - what do you mean by "not greatly effective"?
 [Tomislav Hengl]



should be considered, such as the log ratio approaches 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 approaches (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 approaches combined with original (untransformed) and three different log ratio transformed data in the HRB. The results indicate 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 our study, XGB appeared to be more meaningful and more computationally efficient when dealing with large data sets. In addition, the log ratio approaches 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 approaches. The indirect soil texture classification generated preferable consequences in both cases of accuracy indicators and prediction maps. More appropriate environmental covariates, 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 data² (DOIs are: [10.3972/heihe.009.2013.db](https://doi.org/10.3972/heihe.009.2013.db); [10.3972/heihe.009.2013.db](https://doi.org/10.3972/heihe.009.2013.db); [DOI:10.3972/heihe.00135.2016.db](https://doi.org/10.3972/heihe.00135.2016.db); [10.3972/hiwater.147.2013.db](https://doi.org/10.3972/hiwater.147.2013.db); [10.3972/heihe.037.2014.db](https://doi.org/10.3972/heihe.037.2014.db); [10.3972/heihe.0034.2013.db](https://doi.org/10.3972/heihe.0034.2013.db);

1. more symmetric distribution...

I am not sure what you mean here - PSF's are given and can not be manipulated?

[Tomislav Hengl]

2. Data availability. The...

I checked the links and they are (1) made for Chinese users only, (2) require registration / username+pw. I would recommend putting instead the whole project and data on Github + zenodo.org.

[Tomislav Hengl]