# 1 Compositional balance should be considered in the mapping of soil
# 2 particle-size fractions using hybrid interpolators

3 Mo Zhang[1,2], Wenjiao Shi[1,3]

4 [1]Key Laboratory of Land Surface Pattern and Simulation, Institute of Geographic Sciences and Natural Resources Research,
5 Chinese Academy of Sciences, Beijing 100101, China
6 [2]School of Earth Sciences and Resources, China University of Geosciences, Beijing 100083, China
7 [3]College of Resources and Environment, University of Chinese Academy of Sciences, Beijing 100049, China

8 *Correspondence to:* Wenjiao Shi (shiwj@lreis.ac.cn), Institute of Geographic Sciences and Natural Resources Research,
9 Chinese Academy of Sciences. 11A, Datun Road, Chaoyang District, Beijing 100101, China.

10 **Abstract**. Digital soil mapping of soil particle-size fractions (PSFs) using log-ratio methods is a widely used technique. As a

11 hybrid interpolator, regression kriging (RK) is an alternative way to improve prediction accuracy. However, there is still a lack

12 of comparisons and recommendations when RK is applied for compositional data, and it is not known if the performance based

13 on different balances of isometric log-ratio (ILR) transformation is robust. Here, we compared the generalized linear model

14 (GLM), random forest (RF), and their hybrid patterns (RK) using different transformed data based on three ILR balances, with

15 29 environmental covariables (ECs) for the prediction of soil PSFs in the upper reaches of the Heihe River Basin, China. The

16 results showed that RF performed best, with more accurate predictions, but GLM produced a more unbiased prediction. For

17 the hybrid interpolators, RK was recommended because it widened the data ranges of the prediction values, and modified the

18 bias and accuracy for most models, especially for RF. Moreover, prediction maps generated from RK revealed more details of

19 the soil sampling points. For three ILR balances, different data distributions were produced. Using the most abundant

20 component of the compositional data as the first component of the permutations was not considered the right choice because

21 it produced the worst performance. Compared to the relative abundance of components, we recommend that the focus should

22 be on data distribution. This study provides a reference for the mapping of soil PSFs combined with transformed data at the

23 regional scale.

24 **1 Introduction**

25 Recently, spatial interpolation of soil particle-size fractions (PSFs) has become a focus of soil science researchers. More

26 accurately predicted soil PSFs could contribute to a better understanding of hydrological, physical, and environmental

27 processes (Delbari et al., 2011; Ließ et al., 2012; McBratney et al., 2002).

28 The characteristics of compositional data makes soil PSFs more impressive than other soil properties. Soil PSFs are usually

29 expressed as three components of discrete data – sand, silt, and clay, and carry only relevant percentage information. Soil

30 texture is classified as soil PSFs, which can be demonstrated on a ternary diagram (so-called soil texture triangle). The closure

31 system formed in this triangle is not Euclidean space, but is rather Aitchison space (i.e., the simplex) (Aitchison, 1986). Due

32   to "spurious correlations" (Pawlowsky-Glahn, 1984), traditional statistical methods based on the Euclidean geometry may

33   generate mistakes when dealing directly with soil PSF data (Filzmoser et al., 2009). The requirement for constant sum,

34   nonnegative, unbiased prediction is the key to spatial interpolation (Walvoort and de Gruijter, 2001). Data transformation is

35   crucial for compositional data from the simplex to the real space. Log ratio transformations play a significant role in

36   compositional data analysis, including the additive log-ratio (ALR), centered log-ratio (CLR) (Aitchison, 1986), and isometric

37   log-ratio (ILR) (Egozcue et al., 2003).

38   Although these three log-ratio methods have been widely applied to transform soil PSF data, different study area scales and

39   model selection should be considered when modeling. For local scale study areas, geostatistical models, i.e., ordinary kriging

40   (OK) and compositional kriging, combined with log-ratio transformed data, are sufficient to map spatial patterns, as shown in

41   our previous study (Wang and Shi, 2017). As another perspective, functional compositions combined with the kriging method

42   can also be applied to produce soil particle size curves (PSCs) (Menafoglio et al., 2014), providing an abundance of information.

43   This involves the use of complete and continuous information rather than discrete information, and soil PSFs can be extracted

44   from the predicted soil PSCs (Menafoglio et al., 2016a). Log-ratio transformations can also be combined with functional-

45   compositional data for the stochastic simulation of PSCs (Menafoglio et al., 2016b, Talska et al., 2018). For middle scale study

46   areas, outliers may lead to the overestimation of the variogram, resulting in prediction errors (Lark, 2000). Therefore, the

47   spatial interpolation should take robust variogram estimators into account to improve model performance (Lark, 2003). A

48   previous study proved that applying robust variogram estimators in log-ratio co-kriging significantly improved mapping

49   performance (Wang and Shi, 2018). For large scale study areas, geostatistical models are limited by the number of soil samples

50   and increased spatial variability. An increasing number of studies have concentrated on mapping soil PSFs using different

51   machine-learning models combined with ancillary data (i.e., environmental covariables, ECs) on a broad basin scale (Zhang

52   et al., 2020), national scale (Akpa et al., 2014), and even global scale (Hengl et al., 2017) using log-ratio transformed data.

53   Among these EC-combined models, linear, machine-learning, geostatistical models, and high accuracy surface modeling

54   (Yue et al., 2020) have been commonly used in middle or large-scale studies. Linear models, for example, the generalized

55   linear model (GLM) and multiple linear regression (MLR) have been used in soil PSF predictions with their flexibility and

56   interpretability (Lane, 2002; Buchanan et al., 2012). Many machine-learning models have been applied for the soil PSF

57   interpolation and soil texture classification. For example, tree learners, such as the random forest (RF), have been shown to be

58   advantageous due to their ability to handle overfitting and generate more realistic maps (Zhang et al., 2020). Furthermore,

59   regression kriging (RK), which has been proved to be a powerful and widely accepted method of soil mapping, can not only

60   combine ECs through its regression function, but it also improves model accuracy as a hybrid interpolator for some soil

61   properties, such as topsoil thickness and pH (Hengl et al., 2004; Keskin and Grunwald, 2018). However, the scope of the

62   comparison needs to be expanded to further explore the accuracy and predict compositional data using linear models, machine-

63   learning models, and other models combining RK (hybrid patterns).

64   In log-ratio methods, the ILR method performs better than ALR and CLR in both theory and in practice (Filzmoser and

65   Hron, 2009; Wang and Shi, 2018; Zhang et al., 2020). The ILR method eliminates model collinearity and preserves

66 advantageous properties such as isometry, scale invariance, and sub-compositional coherence, through its use of orthonormal

67 coordinate systems (i.e., balances) using a sequential binary partition (SBP) (Egozcue and Pawlowsky-Glahn, 2005). These

68 choices are not unique, multiple sets of ILR transformed data can be generated by permutations of components (different SBPs)

69 in the compositional data. The choice of an SBP can be based on prior expert knowledge, using a compositional biplot (Lloyd

70 et al., 2012) or variograms and cross-variograms (Molayemat et al., 2018). It has been proven in statistical science that different

71 results are obtained using different choices of ILR balances, and the option of a specific SBP for compositions is crucial for

72 the intended interpretation of coordinates (Fiserova and Hron, 2011). However, most soil science researchers have ignored this

73 point. Martins et al. (2016) reported that clay has been used as the denominator in the ALR method because it is typically the

74 most abundant component of compositions. Few studies have compared the different SBP options from the perspective of

75 accurate assessments and analyzed whether these differences are due to the general characteristics of specific data sets or log-

76 ratio transformations.

77 Therefore, based on our previous work, the objectives of this study were to: (i) compare the spatial prediction accuracy of

78 soil PSFs using a GLM and RF combined with ECs and ILR transformed data; (ii) determine whether hybrid interpolators

79 (GLMRK and RFRK) can improve the interpolation performance; and (iii) explore the distributions of different transformed

80 data and the variation law of precision based on different choices of SBP.

81 **2 Methods and materials**

82 **2.1 Study area**

83 The study area was the upper reaches of the Heihe River Basin (HRB), which is the source of the Heihe River and the central

84 area of runoff generation in the HRB. The elevation in this area ranges from 1640 to 5573 m (Fig. 1), and the climate is damp

85 and cold, being dominated by the Qilian Mountains. The mean annual rainfall in the study area is 350 mm, and the mean annual

86 temperature is lower than 4°C. Meadow and steppe are the dominant vegetation types. Grassland is the primary land-use type.

87 The main soil classes are frigid calcic soil in the southwest of the study area, with cold desert soil dominating the southeast,

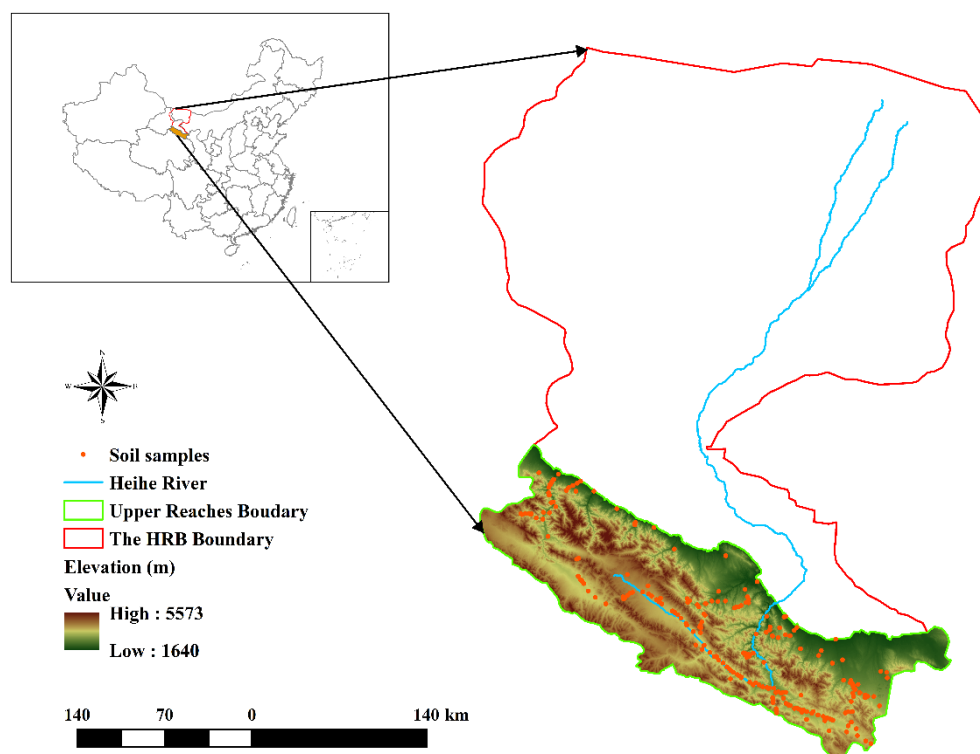88 while Castanozems and Sierozems are distributed in the north of the study area.

**Figure. 1.** The location, elevation, and soil samples on the upper reaches of the Heihe River Basin.

**2.2 Data collection and analysis**

**2.2.1 Soil PSF data**

A total of 262 soil samples were collected in the upper reaches of the HRB based on a purposive sampling strategy and were used to characterize the spatial variability of soil PSFs at the regional scale (Fig. 1). The variability of soil formation factors, such as elevation, soil type, vegetation class, and geomorphology of the upper reaches of the HRB was considered in soil sample collection. The average of three mixed topsoil samples (approximate depth of 0–20 cm) was obtained to reduce the noise of soil sample parameters, and a parallel sample was also measured. Subsequently, about 30 g of each soil sample was air-dried, and chemical and physical analyses were conducted in the laboratory. Soil PSF information was obtained for the soil samples using a Malvern Panalytical Mastersizer 2000, with less than 3% average measurement error.

**2.2.2 The selection of ECs**

There were 29 ECs considered in our study, including both continuous and categorical variables (Table S1.1). They followed the principles of the SCORPAN model (McBratney et al., 2003). The continuous variables included the morphometry and hydrologic characteristics of topographic properties, climatic and vegetative indices, and soil physical and chemical properties (Yi et al., 2015; Song et al., 2016; Yang et al., 2016). The categorical variables included geomorphology, land use types, and

4

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

105    vegetation classes, which were transformed into raster with 1000 m resolution. Due to the intricate patterns of topography in

106    the upper reaches of the HRB, the variable of topographic properties dominated the ECs. The System for Automated

107    Geoscientific Analyses geographic information system (SAGA GIS, Conrad et al., 2015) was applied for a terrain analysis to

108    derive topographic variables using the 30 m resolution digital elevation model (DEM, http://www.gscloud.cn). A collinearity

109    test removed the redundant variables, and the topographic properties were then resampled to 1000 m. More details of the ECs

110    are provided in the Data Availability section.

111    **2.3 ILR transformation and SBP**

112    An orthonormal basis of ILR was chosen to isometrically project the compositions from $S^D$ (the simplex for the Aitchison

113    geometry) to $R^{D-1}$ (real space for the Euclidean geometry). The choice of a specific orthonormal basis for use on $S^D$ can be

114    explained by the SBP for the groups of compositions (Egozcue and Pawlowsky-Glahn, 2005). The choice of the construction

115    of coordinates (i.e., balances) between groups of compositions was calculated as follows:

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

117    where $z_k$ refers to the balance between two groups; $i_1, i_2, \dots, i_{r_k}$ is the $r_k$ part of one group; and $j_1, j_2, \dots, j_{r_k}$ is the $s_k$

118    part of the other group. Therefore, in a stepwise manner, the balances contain all the relevant information of the compositions

119    in two groups. This can also be explained in a tabular form. For soil PSF data (D = 3), all three choices of the balance of SBPs

120    are shown in Table 1. The first component of the ILR contained all the information on soil PSFs, and the main difference in

121    the choice of balances for soil PSFs was the order of the three parts, i.e., the first order of the soil PSF component was used as

122    the numerator of the first ILR equation. In our study, three SBP balances, SBP1, SBP2, and SBP3, were transformed from the

123    original soil PSF data, and the orders of soil PSF data were $(sand, silt, clay)$, $(silt, clay, sand)$, and $(clay, sand, silt)$,

124    respectively. The transformation equations for the ILR can be derived from Eq. (1), and were defined as Eqs. (2) and (3). The

125    inverse equations for ILR were defined as Eqs. (4), (5), (6). The ILR transformation and its inverse were conducted using the

126    R package "compositions" (K. Gerald van den Boogaart and Raimon Tolosana, 2014).

127    $\mathbf{z} = (z_1, \dots z_{D-1}) = ILR(\mathbf{x})$, and for $i = 1, \dots, D-1$ and component $x_i$,    (2)

128    $$z_i = \sqrt{\frac{D-i}{D-i+1}} \ln \frac{x_i}{\sqrt[D-i]{\prod_{j=i+1}^{D} x_j}}.$$    (3)

129    $$Y(x_j) = \sum_{j=1}^{D} \frac{ILR(x_j)}{\sqrt{j \times (j+1)}} - \sqrt{\frac{j-1}{j}} \times ILR(x_j),$$    (4)

130    $ILR(x_0) = ILR(x_D) = 0,$    (5)

131    $$\overline{ILR}(x_j) = \frac{exp(Y(x_j))}{\sum_{j=1}^{D} exp(Y(x_j))}.$$    (6)

132    **Table 1** All choices of SBPs for soil PSF data (D = 3), the orders of soil PSFs data are $(sand, silt, clay)$, $(silt, clay, sand)$

133    and $(clay, sand, silt)$ for SBP1, SBP2 and SBP3.

| Groups | Step | Sand | Silt | Clay | r | s | Balance |
|--------|------|------|------|------|---|---|---------|
| SBP1 | 1 | + | - | - | 1 | 2 | Step1: $z_1 = \sqrt{\frac{2}{3}} ln \frac{sand}{\sqrt{silt \times clay}}$ |
|      | 2 | 0 | + | - | 1 | 1 | Step2: $z_2 = \sqrt{\frac{1}{2}} ln \frac{silt}{clay}$ |
| SBP2 | 1 | - | + | - | 1 | 2 | Step1: $z_1 = \sqrt{\frac{2}{3}} ln \frac{silt}{\sqrt{clay \times sand}}$ |
|      | 2 | - | 0 | + | 1 | 1 | Step2: $z_2 = \sqrt{\frac{1}{2}} ln \frac{clay}{sand}$ |
| SBP3 | 1 | - | - | + | 1 | 2 | Step1: $z_1 = \sqrt{\frac{2}{3}} ln \frac{clay}{\sqrt{sand \times silt}}$ |
|      | 2 | + | - | 0 | 1 | 1 | Step2: $z_2 = \sqrt{\frac{1}{2}} ln \frac{sand}{silt}$ |

134

## 2.4 Linear model, machine-learning model, and hybrid patterns

### 2.4.1 GLM

The GLM is an extended version of the linear model, which contains response variables, with non-normal distributions (Nelder and Wedderburn, 1972). The link function is embedded into the GLM to ensure the classical linear model assumptions. The scaled dependent variables and the independent variables can be connected using a link function for the additive combination of model effects, the choice of link function depends on the distribution of response variables (Venables and Dichmont, 2004). A Gaussian distribution with an identity link function was applied in our study, which produced consequences equivalent to that of MLR (Nickel et al., 2014). However, categorical variables can be directly trained in the GLM without setting dummy variables. The Akaike's information criterion (AIC) was applied to choose the best predictors and remove model multicollinearity using a backward stepwise algorithm, and the combinations of ECs for different ILR data were then obtained (Table S2.1).

### 2.4.2 RF

The RF is a non-parametric technique, which combines the bagging method with a selection of random variables as an extended version of a regression tree (RT) (Breiman, 1996, 2001). It can improve model prediction accuracy by producing and aggregating multiple tree models. The principle of the RF is to merge a group of "weak trees" together to generate a "powerful forest." The bootstrap sampling method was applied for each tree, and each predictor was selected randomly from all model predictors. The "out of bag" (OOB) data were applied to produce reliable estimates in an internal validation using a random subset independent of the training tree data. Three parameters needed to be tuned: number of trees ($ntree$); minimum size of terminal nodes ($nodesize$), and number of variables randomly sampled as predictors for each tree ($mtry$) (Liaw and Wiener, 2001). The standard value of the $mtry$ parameter was one-third of the total number of predictors, while $ntree$ and

155  *nodesize* were 500 and 5, respectively. For regression, the mean square errors (MSEs) of predictions were estimated to train

156  the trees. The variable importance of the RF was produced from the OOB data using the "importance" function. One of the

157  benefits of the RF is that the ensembles of trees are used without pruning to ensure that the most significant amount of variance

158  can be expressed. Moreover, the RF can reduce model overfitting and normalization is unnecessary due to the effects on the

159  value range being insensitive. The GLM and RF algorithms and the parameter adjustment of the RF were conducted in the R

160  package "caret" (Max Kuhn, 2018).

161

162  **2.4.3 RK**

163  Regression kriging is a hybrid interpolation technique that combines regression models (e.g., GLM and RF) with the residuals

164  of OK (Odeh et al., 1995). Mathematically, the RK method corresponds to two interpolators, the regression part and the kriging

165  part, which are operated separately (Goovaerts, 1999). One limitation of using only the regression part is that it is usually only

166  useful within the range of values of the training sets (Hengl et al., 2015). The principle of the RK method is that the regression

167  model explains a deterministic component of spatial variability, and the interpolation of regression residuals generated from

168  OK is used to describe the spatial variability (Bishop and McBratney, 2001; Hengl et al., 2004). The residuals create a

169  variogram (e.g., Gaussian, spherical, or exponential) for models based on the MSE from the results of a cross-validation. First,

170  we used the regression part (GLM or RF) to predict soil PSFs, the residual from the fitted model was then calculated by

171  subtracting the regression part from the observations. Subsequently, OK was applied for the whole study area to interpolate

172  the residuals. Finally, the regression prediction and the predicted residuals at the same location were summed. The variograms

173  of the RK method were generated automatically using the "autofitVariogram" function in the R package "automap" (Hiemstra

174  et al., 2009).

175  **2.5 Prediction method system and validation**

176  The method system of spatial interpolation models for soil PSFs is presented in Table 2. We systematically compared 12

177  models: four interpolators, including GLM and RF with or without RK, and three SBPs of the ILR transformation method. For

178  the validation of model performance, the independent data set validation was used to evaluate the prediction bias and accuracy

179  of the models. The sub-training sets (70%) and the sub-testing sets (30%) were randomly selected from data independently,

180  and this process was repeated 30 times.

181  **Table 2.** The method system of spatial interpolation models of soil PSFs.

| Models | GLM | GLMRK | RF | RFRK |
|---|---|---|---|---|
| ILR_SBP1 | GLM_SBP1 | GLMRK_SBP1 | RF_SBP1 | RFRK_SBP1 |
| ILR_SBP2 | GLM_SBP2 | GLMRK_SBP2 | RF_SBP2 | RFRK_SBP2 |
| ILR_SBP3 | GLM_SBP3 | GLMRK_SBP3 | RF_SBP3 | RFRK_SBP3 |

182

183    The mean error (ME), the root mean square error (RMSE), and Aitchison distance (AD) were used to evaluate and compare

184    the prediction performance. The ME and RMSE measure prediction bias and accuracy, respectively (Odeh et al., 1995). The

185    AD is an overall indicator of compositional analysis, which describes the distance between two compositions. Generally, in an

186    accurate, unbiased model all three values will be close to 0. The ME, RMSE, and AD were calculated as follows:

187    $ME = \frac{1}{n}\sum_{i=1}^{n}(M_i - P_i),$ (7)

188    $RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(M_i - P_i)^2},$ (8)

189    $AD = \left[\sum_{i=1}^{D}(\log\frac{M_i}{G(M)} - \log\frac{P_i}{G(P)})^2\right]^{0.5},$ (9)

190    where $M_i$ and $P_i$ are the measured and predicted values at the $i$th position, respectively; $n$ refers to the number of soil

191    samples; $D$ is the number of dimensions of compositions; and $G(M)$ and $G(P)$ denote the geometric mean with the form

192    $G(\mathbf{x}) = (x_1, \ldots, x_D)^{1/D}$ of the measured and predicted values, respectively.

193

**2.6 Covariance structure analysis**

195    The interpretation of the ILR balances is based on a decomposition of the covariance (COV) structure (Fiserova and Hron,

196    2011). We calculated the variance (VAR), COV, and the corresponding correlation coefficient (CC) of ILR transformed data

197    based on different SBP. The equations for calculating VAR, COV, and CC were derived from Eq. (1) as follows:

198    $VAR(z) = \frac{1}{r+s}\sum_{p=1}^{r}\sum_{q=1}^{s}var(\ln\frac{x_{ip}}{x_{jq}}) - \frac{s}{2r(r+s)}\sum_{p=1}^{r}\sum_{q=1}^{r}var(\ln\frac{x_{ip}}{x_{iq}}) - \frac{r}{2s(r+s)}\sum_{p=1}^{s}\sum_{q=1}^{s}var(\ln\frac{x_{jp}}{x_{jq}}) -$

199    $\frac{r}{2s(r+s)}\sum_{p=1}^{s}\sum_{q=1}^{s}var(\ln\frac{x_{jp}}{x_{jq}})$ (10)

200    $COV(z_1, z_2) = \frac{C}{2r_1s_2}\sum_{p=1}^{r_1}\sum_{q=1}^{s_2}var(\ln\frac{x_{ip}^1}{x_{jq}^2}) + \frac{C}{2r_2s_1}\sum_{p=1}^{r_2}\sum_{q=1}^{s_1}var(\ln\frac{x_{ip}^2}{x_{jq}^1}) - \frac{C}{2r_1r_2}\sum_{p=1}^{r_1}\sum_{q=1}^{r_2}var(\ln\frac{x_{ip}^1}{x_{iq}^2}) -$

201    $\frac{C}{2s_1s_2}\sum_{p=1}^{s_1}\sum_{q=1}^{s_2}var(\ln\frac{x_{jp}^1}{x_{jq}^2}),$ (11)

202    $CC = \frac{COV(z_1,z_2)}{\sqrt{var(z_1)\cdot var(z_2)}}$ (12)

203    For soil PSF data, Eqs. (10), (11), and (12) can be simplified to three dimensions. The relationship between the ratios of soil

204    PSF components and the dominant roles of ILR transformed data were indicated from the covariance structure. All the

205    statistical analyses, such as the descriptive statistics of soil PSF data, calculation and evaluation of indicators, and the spatial

206    prediction mapping, were performed using the R statistical program (R Development Core Team, 2019).

207

## 3 Results

### 3.1 Exploratory data analysis

#### 3.1.1 Descriptive statistics of soil PSF data

From the descriptive statistics of the original (raw) and ILR transformed data, the silt fraction dominated the soil PSFs, accounting for a more substantial amount than the sand and clay fractions. The distributions of the sand and clay fractions were similar (Fig. 2a). The ILR transformed data based on the three SBPs revealed different distributions (Figs. 2b, 2c, and 2d). For example, two ILR components (ILR1 and ILR2) for SBP1 had a symmetric distribution around zero at the $x$-axis (Fig. 2b). In comparison, the distribution of data generated from SBP2 or SBP3 had a mirrored symmetry, with a left-skewed ILR1 of SBP2 and right-skewed ILR2 of SBP3 (Figs. 2c and 2d). The comparison of means and medians demonstrated that the back-transformed means of three sets of ILR transformed data were the same, and the mean ILR of sand was closer to the median compared with the original soil PSF data. In contrast, the opposite patterns were apparent for the silt and clay components (Fig. 2e).



| Center | Sand | Silt | Clay | Sand_ILR | Silt_ILR | Clay_ILR |
|--------|------|------|------|----------|----------|----------|
| Mean | 19.61 | 66.80 | 13.59 | 18.07 | 71.48 | 10.45 |
| Median | 17.03 | 67.44 | 14.53 | — | — | — |

**Figure 2.** Descriptive statistics of original soil PSF and ILR transformed data using different SBPs. Note that means of Sand_ILR, Silt_ILR, and Clay_ILR from different SBPs were back-transformed to the real space.

#### 3.1.2 Covariance structure of ILR transformed data with different balances

The covariance analysis of the transformed data of soil PSFs based on the different SBPs showed that the variance VarILR_1 of SBP3 was the largest, followed by the VarILR_1 of SBP1 and SBP2 (Table 3). The variance of the second component of ILR (VarILR_2) followed the opposite pattern to that of VarILR_1. The COV and corresponding CC followed the same pattern of SBP1 > SBP3 > SBP2. The first ILR equation ($z_1$ in Table 1) contained all information of soil PSFs, while the second one

229   ($z_2$ in Table 1) included only two components. The information of VarILR_1 was therefore more abundant. All of VarILR_1

230   and VarILR_2 values were not 0 (or not nearly 0), indicating that there was no constant (or almost constant) value in any two

231   ratios of soil PSF components. The COV of SBP3 was close to 0, indicating that the proportions of *clay/sand* and *clay/silt*

232   were approximately the same. The same results were generated from the corresponding CC. For the distribution of soil PSFs

233   in a ternary diagram (the United States Department of Agriculture texture triangle, USDA), the main texture class was silt

234   loam (Fig. 3a). The biplot of soil samples demonstrated that the rays of the three components, i.e., sand, silt, and clay, were

235   reasonably well clustered at about 120° in the three groups (Fig. 3b).

236

237   **Table 3.** Covariance structure of soil PSFs based on different SBPs. VarILR_1 and VarILR_2 denote the variance of the first

238   and the second component of ILR, respectively. COV refers to the covariance of ILR1 and ILR2. CC is the correlation

239   coefficient.

| Balances | VarILR_1 | VarILR_2 | COV | CC |
|----------|----------|----------|-----|-----|
| SBP1 | 0.53 | 0.71 | 0.32 | 0.52 |
| SBP2 | 0.39 | 0.86 | -0.24 | -0.41 |
| SBP3 | 0.94 | 0.30 | -0.09 | -0.16 |

240



241

242   **Figure 3.** The distribution of soil PSFs in the USDA triangle (a) and biplot graph (b). The red curve was fitted by loess function.

243   **3.2 Accuracy comparison of different models using ILR data**

244   The first three rows of the boxplots in Figs. 4a, 4b, and 4c indicate the bias of the different models according to their ME

245   values. The ME of sand was closest to 0, followed by the MEs of clay and silt. GLM was more unbiased than RF, with lower

246   ME values. After combining with RK, there was an improvement in the ME for most GLM and RF models (Figs. 4a, 4b, and

247   4c). For the accuracy assessment, the RMSE of silt was higher than for the other two components. The GLMRK did not

248   perform as well as expected in terms of the RMSE, with only the sand component having an improved RMSE (Fig. 4d).

249   However, the RFRK performed better than the GLMRK and improved the accuracy of most parts compared with the RF,

250   except for the RFRK_SBP1 of sand. As an overall indicator, AD showed that the RF (or RFRK) performed better than the

251   GLM (or GLMRK) in terms of both average RMSE values and uncertainties (Fig. 4g). Moreover, the RFRK improved the AD

252   values for the SBP2 and SBP3 methods. For the uncertainty assessment, the RF generated lower uncertainties than the GLM,

253   and the models combined with RK further reduced the uncertainties for most GLM and RF models.



**Figure 4.** Accuracy comparison of GLM, RF, and their RK patterns combined with three ILR balances. The mean values of

256   different model indicators were calculated in their boxes.

258   The model performances were different for the three SBPs. To better evaluate model performance using the different SBP

259   balances, we graded each box from 1 to 3, and the final results are shown in Fig. 5. The results demonstrated that SBP1

260 performed best, with the lowest ME value of all models. For the accuracy comparison there was no apparent pattern, but

261 accuracy could be considered hierarchically: (1) for the GLM, SBP1 performed better than the other two SBP methods, which

262 also performed well when RK was combined (GLMRK); (2) for RF, SBP1 produced the best result. However, the introduction

263 of RK resulted in the Score2 of SBP3 performing best among the three SBPs. However, RFRK of SBP1 performed worst

264 according to the values of Score2 and Score5. Finally, for the comprehensive assessment, SBP1 performed best among three

265 SBPs according to Score6. More details and calculation processes can be found in the Supplementary Material (Table S4.1).



266

267 **Figure 5.** Ranking score of model performance based on three SBPs. Score1 and Score2 are the sum scores of ME and RMSE

268 for each model, respectively; Score3 is the sum scores of ME, RMSE and AD for each model, Score4 and Score5 are the sum

269 scores of ME or RMSE for GLM$_{all}$ (GLM and GLMRK) and RF$_{all}$ (RF and RFRK), Score6 is the sum scores of all indicators.

270 The lower the value, the better the model performance.

271

272 **3.3 Spatial prediction maps of soil PSFs generated from the different models**

273 Prediction maps of soil PSFs made from the different models are shown in Figs. 6, S3.1, and S3.2. For the components of soil

274 PSFs, the maps of the three group maps followed a similar rule. The GLM and GLMRK produced more extensive ranges of

275 predicted values, and their maps were more relevant to the real environment. However, the RF and RFRK predicted a relatively

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

276 narrow range of low values for these components, revealing a smoother distribution than that generated by the GLM and

277 GLMRK. Unlike the regression methods, the RF and RFRK methods produced hot and cold spots on the prediction maps and

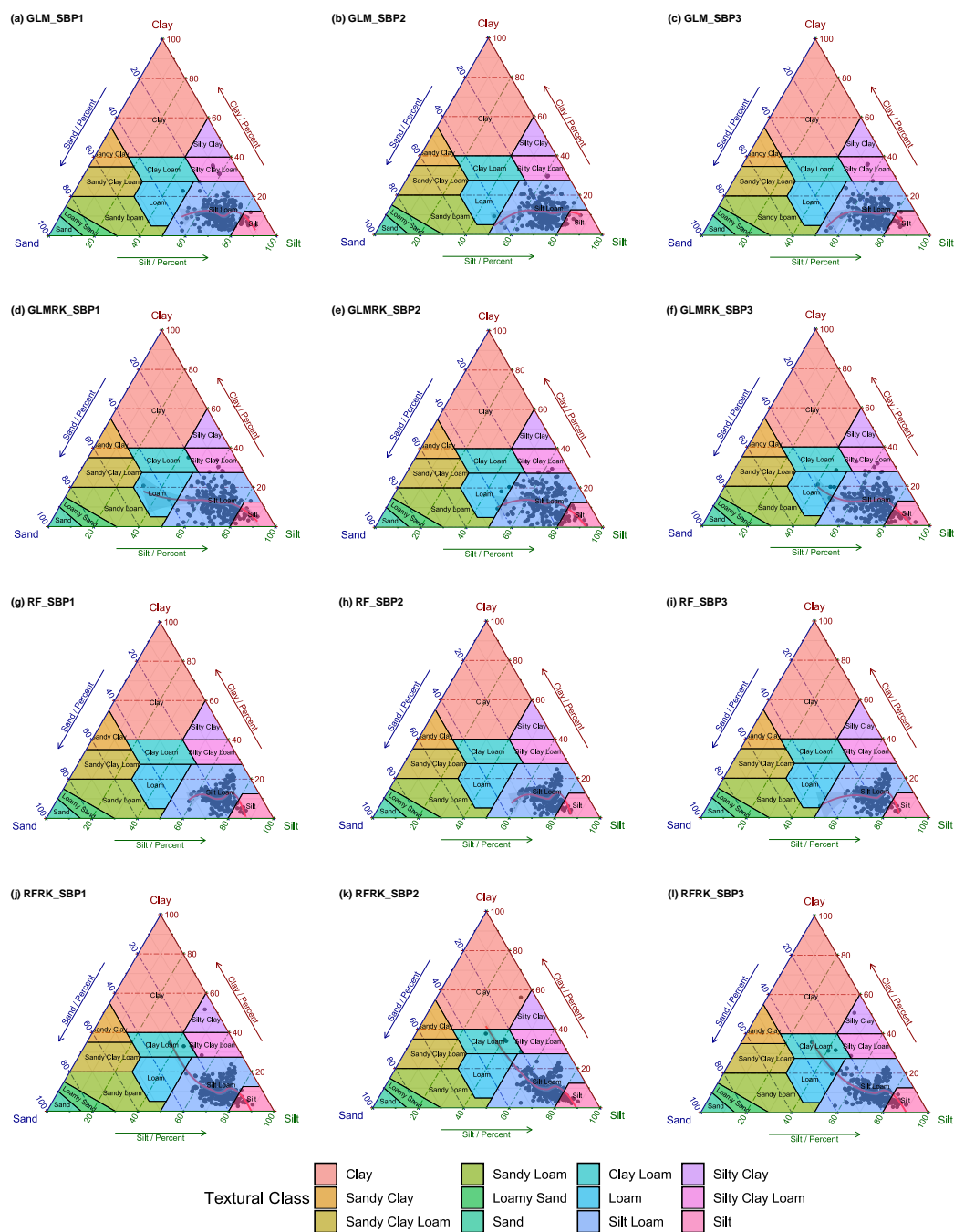278 more details of the soil sampling points were apparent (Fig. S5.1).



279

280 **Figure 6.** Spatial prediction maps of the sand component of the upper reaches of the Heihe River Basin.

281 **3.4 Spatial distribution of soil texture classes in the USDA triangles**

282 The predicted soil textures in the USDA texture triangles (Fig. 7) showed that most predictions fell within the range of observed

283    soil textures (Fig. 3a), and silt loam was the dominant soil texture in all cases. The GLM produced a more discrete distribution

284    than the RF, and the RK method expanded the dispersion. In the trends of the predicted samples, the silt components predicted

285    from all models were overestimated. The pattern fitting curves indicated that the prediction results were closer to the bottom

286    right of the USDA triangle than the soil PSF observations. The GLMRK and RFRK curves were longer than the GLM and RF

287    curves, with a more extensive range of values in triangles. Compared with the GLMRK, the RFRK produced a more upward

288    extension (Figs. 7j, k, l). It was clear that the clay fraction was overestimated and the sand fraction was underestimated.

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

**Figure 7.** Predicted 262 soil samples in USDA texture triangles using (a) GLM_SBP1, (b) GLM_SBP2, (c) GLM_SBP3, (d) GLMRK_SBP1, (e) GLMRK_SBP2, (f) GLMRK_SBP3, (g) RF_SBP1, (h) RF_SBP2, (i) RF_SBP3, (j) RFRK_SBP1, (k) RFRK_SBP2, (l) RFRK_SBP3. Red fitting curves in triangles showed the trends.

Hydrology and
Earth System
Sciences
Discussions

## 4 Discussion

### 4.1 Comparison of the GLM, RF, and RK patterns using ILR data

295 We found RF reveal more accurate results, but with more bias than the GLM, and RK method improved the performance in
296 terms of bias for most models and the accuracy of the RF. Odeh et al. (1995) indicated that RK was superior to the linear
297 models, such as MLR, which was reflected in the prediction results for sand in our study. Scarpone et al. (2016) reported that
298 as a hybrid interpolator, the RFRK outperformed the RF when making soil thickness predictions. We proved that RFRK was
299 also suitable for compositional data and improved model performance when combining with the ILR transformation. In
300 summary, the GLM and RF had both advantages and disadvantages when considering the trade-off between bias and accuracy.

301 The results of GLM and GLMRK should not depend on the ILR basis being chosen, which has been proved by previous
302 studies on the use of linear models and kriging for compositional data (Pawlowsky-Glahn et al, 2015). However, the GLM
303 model used the "glmStepAIC" algorithm (i.e., a stepwise regression) to select the best combination of environmental
304 covariables for each ILR component (Table S2.1). Therefore, the variable inputs are different for these ILR data, and further
305 impact the accuracy assessment and prediction maps. In addition, the difficulty with the use of the GLM is the need for a back-
306 transformation. There is a need to present results on the original untransformed scale after conducting the analysis on a
307 transformed level, which may produce spurious results (Lane, 2002). In our study, we compared the means of ILR transformed
308 data and the original data. We proved the feasibility of the ILR transformation method, especially for meeting the requirements
309 of compositional data. However, the accuracy of the GLM still needs to be improved, which may be because the transformed
310 data did not follow a normal distribution (Fig. 2).

311 Although the RF had the advantage of prediction accuracy, the limited interpretability of the consequences made it difficult
312 to modify the prediction bias – each tree from the model cannot be examined individually (Grimm et al., 2008). Moreover, the
313 ILR transformation before modeling increased the difficulty of interpretation for not only the predicted values on the ILR scale
314 but also the residuals. The back-transformation of the optimal estimate of log-ratio variables does not generate the optimal
315 estimation of compositional data (Lark and Bishop, 2007), which should also be considered.

### 4.2 Comparison of three SBPs of ILR transformation

317 For the comparison of the three SBPs, the ME and RMSE performed better when using SBP1 for ILR transformed data,
318 which may be interpreted as the distributions of the ILR1 and ILR2 of SBP1 being more symmetric (Fig. 2b). In contrast, the
319 performance of SBP2 was worse than that of SBP1 and SBP3 because the ILR_1 component, including all the soil PSF
320 information, was left-skewed (Fig. 2c). This result was especially apparent for the GLM and GLMRK, because the data in a
321 linear model needs to be normally distributed (Lane, 2002).

322 The negligible difference among the three SBP balances revealed a triangular shape with a cluster at about 120° (Fig. 3b).
323 This could be interpreted as the three soil PSFs having a mixed pattern, with each component dominated by the components
324 in one cluster (Tolosana-Delgado et al., 2005). Although the silt component dominated the soil PSFs (Fig. 2a), sand and clay

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

325   also played important roles in soil compositions. Taking either the most abundant component of the compositional data as the

326   denominator (Martins et al., 2016) or the first component of the permutations did not provide convincing evidence. Because

327   using the most abundant component of the compositional data as the primary component of the alterations, i.e., SBP2, resulted

328   in a relatively poor performance compared to the other SBPs. Thus, we recommend that the focus should be on data distribution.

329   Furthermore, the choice of balance and combination of RK are also the key to improving model accuracy, as shown by the

330   result of the RFRK-SBP3 model (Fig. 4).

**4.3 Limitations**

332   Firstly, the scope of this study is limited to independent modeling. Each ILR component was modeled separately, which may

333   suboptimal because they cannot further consider the cross correlations among ILR coordinates. However, the study

334   demonstrated the relation of the raw data (sand, silt, and clay), and has confirmed that the currently used prediction models are

335   suitable. In our pervious study, we have used compositional kriging (CK) for the spatial prediction of soil PSFs (Wang and Shi,

336   2017), and the cross correlations of ILRs can be taken into account using CK. Although it is optimal, it cannot consider different

337   balances of ILR, nor can it be combined with the hybrid interpolator (e.g., RK). Moreover, predicting each ILR component

338   separately was a more suitable approach for the spatial prediction models currently used (such as the GLM and RF). Therefore,

339   more alternative spatial prediction models combined with interpretation of ILR balances for compositional data should be

340   considered in the future. For example, CK and high accuracy surface modelling (HASM; Yue et al., 2016) can be applied for

341   small scale study areas. For large scale study areas, multivariate RF (Segal and Xiao, 2011) can be combined with a log-ratio

342   transformation and hybrid interpolation method, enabling the cross correlations among ILR coordinates to be better interpreted.

343   Secondly, the weighting problem was not considered in this study, because the ILR method can be qualified as an unweighted

344   log-ratio transformation, giving all parts the same weight for both the definition of the total variance and the reduction of

345   dimension. This may enlarge the ratios generated from the rare parts, which would dominate the analysis (Greenacre and Lewi,

346   2009). The pairwise log-ratio can be used to set weights by their proportions when there is no additional knowledge about the

347   component measurement errors (Greenacre, 2019). Nevertheless, all three parts of the soil PSF data dominated the biplot

348   diagram, without the influence of rare elements and with no redundancy; thus, none of the shortcomings mentioned above

349   were apparent. Accuracy assessments using a pairwise log-ratio transformation require further study in the future.

**5 Conclusions**

351   We evaluated and compared the performance of the GLM, RF, and their hybrid pattern (i.e., GLMRK and RFRK) using

352   different balances of ILR transformed data. The bias of the GLM was lower than that of the RF; however, the accuracy of the

353   GLM was relatively low. More discrete distributions and broader ranges of prediction value distributions were produced from

354   GLMs in the USDA soil texture triangles. In other words, different predicted data sets were generated from the use of the GLM

355   and RF, with unbiased and inaccurate predictions for the GLM and biased and more accurate predictions for the RF.

356   The hybrid patterns, GLMRK and RFRK, were found to be the best solution because it produced a relatively high prediction

357   accuracy and strong correlations with ECs, providing more details about the soil sampling points (hot spots and cold spots)

358 compared with only the regression model. However, the non-normal distribution of ILR data and its residuals, and more data

359 transformation and inverse transformation processes make models further difficult to interpreted and improve.

360     For the different SBPs, the three SBP-based data generated different distributions, but no pattern was apparent. This could

361 be explained by the angle of the biplot diagram, with three rays of soil PSF components clustered into three modes, and each

362 part dominating its cluster. Using the most abundant component of the compositional data as the first component of the

363 permutations was not considered the right choice because SBP2 produced the worst performance. Thus, we recommend that

364 the focus should be on data distribution. This study can provide a reference for the spatial simulation of soil PSFs combined

365 with ECs at the regional scale, and how to choose the balances of ILR transformed data.

366

379

382

384

392

# References

394  Aitchison, J.: The statistical analysis of compositional data, Chapman and Hall, Ltd., 416 pp., 1986.

395  Akpa, S. I. C., Odeh, I. O. A., Bishop, T. F. A., and Hartemink, A. E.: Digital mapping of soil particle-size fractions for Nigeria,
396       Soil Sci. Soc. Am. J., 78, 1953-1966, https://doi.org/10.2136/sssaj2014.05.0202, 2014.

397  Bishop, T. F. A., and McBratney, A. B.: A comparison of prediction methods for the creation of field-extent soil property maps,
398       Geoderma, 103, 149-160, https://doi.org/10.1016/S0016-7061(01)00074-X, 2001.

399  Breiman, L.: Bagging predictors, Machine Learning, 24, 123-140, https://doi.org/10.1023/a:1018054314350, 1996.

400  Breiman, L.: Random forests, Machine Learning, 45, 5-32, https://doi.org/10.1023/a:1010933404324, 2001.

401  Buchanan, S., Triantafilis, J., Odeh, I. O. A., and Subansinghe, R.: Digital soil mapping of compositional particle-size fractions
402       using proximal and remotely sensed ancillary data, Geophysics, 77, WB201-WB211, https://doi.org/10.1190/geo2012-
403       0053.1, 2012.

404  Conrad, O., Bechtel, B., Bock, M., Dietrich, H., Fischer, E., Gerlitz, L., Wehberg, J., Wichmann, V., and Böhner, J.: System
405       for automated geoscientific analyses (SAGA) v. 2.1.4, Geosci. Model Dev., 8, 1991-2007, https://doi.org/10.5194/gmd-8-
406       1991-2015, 2015.

407  Delbari, M., Afrasiab, P., and Loiskandl, W.: Geostatistical analysis of soil texture fractions on the field scale, Soil and Water
408       Research, 6, 173-189, https://doi.org/10.17221/9/2010-SWR, 2011.

409  Egozcue, J. J., Pawlowsky-Glahn, V., Mateu-Figueras, G., and Barcelo-Vidal, C.: Isometric logratio transformations for
410       compositional data analysis, Mathematical Geology, 35, 279-300, https://doi.org/10.1023/a:1023818214614, 2003.

411  Egozcue, J. J., and Pawlowsky-Glahn, V.: Groups of parts and their balances in compositional data analysis, Mathematical
412       Geology, 37, 795-828, https://doi.org/10.1007/s11004-005-7381-9, 2005.

413  Filzmoser, P., and Hron, K.: Correlation analysis for compositional data, Math Geosci., 41, 905-919,
414       https://doi.org/10.1007/s11004-008-9196-y, 2009.

415  Filzmoser, P., Hron, K., and Reimann, C.: Univariate statistical analysis of environmental (compositional) data: Problems and
416       possibilities, Sci. Total Environ., 407, 6100-6108, https://doi.org/10.1016/j.scitotenv.2009.08.008, 2009.

417  Fiserova, E., and Hron, K.: On the interpretation of orthonormal coordinates for compositional data, Math Geosci., 43, 455-
418       468, https://doi.org/10.1007/s11004-011-9333-x, 2011.

419  Goovaerts, P.: Geostatistics in soil science: state-of-the-art and perspectives, Geoderma, 89, 1-45,
420       https://doi.org/10.1016/S0016-7061(98)00078-0, 1999.

421  Greenacre, M., and Lewi, P.: Distributional equivalence and subcompositional coherence in the analysis of compositional data,
422       contingency tables and ratio-scale measurements, Journal of Classification, 26, 29-54, https://doi.org/10.1007/s00357-
423       009-9027-y, 2009.

424  Greenacre, M.: variable selection in compositional data analysis using pairwise logratios, Math Geosci., 51, 649-682,

425      https://doi.org/10.1007/s11004-018-9754-x, 2019.

426  Grimm, R., Behrens, T., Märker, M., and Elsenbeer, H.: Soil organic carbon concentrations and stocks on Barro Colorado

427      Island – Digital soil mapping using Random Forests analysis, Geoderma, 146, 102-113,

428      https://doi.org/10.1016/j.geoderma.2008.05.008, 2008.

429  Hengl, T., Heuvelink, G. B. M., and Stein, A.: A generic framework for spatial prediction of soil variables based on regression-

430      kriging, Geoderma, 120, 75-93, https://doi.org/10.1016/j.geoderma.2003.08.018, 2004.

431  Hengl, T., Heuvelink, G. B. M., Kempen, B., Leenaars, J. G. B., Walsh, M. G., Shepherd, K. D., Sila, A., MacMillan, R. A., de

432      Jesus, J. M., Tamene, L., and Tondoh, J. E.: Mapping soil properties of Africa at 250 m resolution: random forests

433      significantly improve current predictions, Plos One, 10, 26, https://doi.org/10.1371/journal.pone.0125814, 2015.

434  Hengl, T., Mendes de Jesus, J., Heuvelink, G. B. M., Ruiperez Gonzalez, M., Kilibarda, M., Blagotić, A., Shangguan, W.,

435      Wright, M. N., Geng, X., Bauer-Marschallinger, B., Guevara, M. A., Vargas, R., MacMillan, R. A., Batjes, N. H., Leenaars,

436      J. G. B., Ribeiro, E., Wheeler, I., Mantel, S., and Kempen, B.: SoilGrids250m: Global gridded soil information based on

437      machine learning, Plos One, 12, e0169748, https://doi.org/10.1371/journal.pone.0169748, 2017.

438  Hiemstra, P. H., Pebesma, E. J., Twenhöfel, C. J. W., and Heuvelink, G. B. M.: Real-time automatic interpolation of ambient

439      gamma dose rates from the Dutch radioactivity monitoring network, Computers & Geosciences, 35, 1711-1721,

440      https://doi.org/10.1016/j.cageo.2008.10.011, 2009.

441  Keskin, H., and Grunwald, S.: Regression kriging as a workhorse in the digital soil mapper's toolbox, Geoderma, 326, 22-41,

442      https://doi.org/10.1016/j.geoderma.2018.04.004, 2018.

443  K. Gerald van den Boogaart, and Raimon Tolosana, M. B.: Compositions: compositional data analysis, R package version

444      1.40-1 ed., available at: https://cran.rstudio.com/web/packages/compositions/index.html (last access: 14 July 2020), 2014.

445  Lane, P. W.: Generalized linear models in soil science, European Journal of Soil Science, 53, 241-251,

446      https://doi.org/10.1046/j.1365-2389.2002.00440.x, 2002.

447  Lark, R. M.: A comparison of some robust estimators of the variogram for use in soil survey, European Journal of Soil Science,

448      51, 137-157, https://doi.org/10.1046/j.1365-2389.2000.00280.x, 2000.

449  Lark, R. M.: Two robust estimators of the cross-variogram for multivariate geostatistical analysis of soil properties, European

450      Journal of Soil Science, 54, 187-201, https://doi.org/10.1046/j.1365-2389.2003.00506.x, 2003.

451  Lark, R. M., and Bishop, T. F. A.: Cokriging particle size fractions of the soil, Eur. J. Soil Sci., 58, 763-774,

452      https://doi.org/10.1111/j.1365-2389.2006.00866.x, 2007.

453  Liaw, A., and Wiener, M.: Classification and regression by random forest, 23, available at: https://cran.r-

454      project.org/web/packages/randomForest/index.html (last access: 14 July 2020), 2001.

455  Ließ, M., Glaser, B., and Huwe, B.: Uncertainty in the spatial prediction of soil texture: Comparison of regression tree and

456      Random Forest models, Geoderma, 170, 70–79, https://doi.org/10.1016/j.geoderma.2011.10.010, 2012.

457  Lloyd, C. D., Pawlowsky-Glahn, V., and Jose Egozcue, J.: Compositional data analysis in population studies, Annals of the

458      Association of American Geographers, 102, 1251-1266, https://doi.org/10.1080/00045608.2011.652855, 2012.

459 Martins, A. B. T., Bonat, W. H., and Ribeiro, P. J.: Likelihood analysis for a class of spatial geostatistical compositional models,
460      Spat. Stat., 17, 121-130, https://doi.org/10.1016/j.spasta.2016.06.008, 2016.

461 Max Kuhn: Caret: Classification and regression training, R package version 6.0-80 ed., available at: https://cran.r-
462      project.org/web/packages/caret/index.html (last access: 14 July 2020), 2018.

463 McBratney, A. B., Minasny, B., Cattle, S. R., and Vervoort, R. W.: From pedotransfer functions to soil inference systems,
464      Geoderma, 109, 41-73, https://doi.org/10.1016/S0016-7061(02)00139-8, 2002.

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

467 Menafoglio, A., Guadagnini, A., and Secchi, P.: A kriging approach based on Aitchison geometry for the characterization of
468      particle-size curves in heterogeneous aquifers, Stochastic Environmental Research and Risk Assessment, 28, 1835-1851,
469      https://doi.org/10.1007/s00477-014-0849-8, 2014.

470 Menafoglio, A., Guadagnini, A., and Secchi, P.: A kriging approach based on Aitchison geometry for the characterization of
471      particle-size   curves   in   heterogeneous   aquifers,   Stoch.   Environ.   Res.   Risk   Assess.,   28,   1835-1851,
472      https://doi.org/10.1007/s00477-014-0849-8, 2014.

473 Menafoglio, A., Secchi, P., and Guadagnini, A.: A class-kriging predictor for functional compositions with application to
474      particle-size curves in heterogeneous aquifers, Math Geosci., 48, 463-485, https://doi.org/10.1007/s11004-015-9625-7,
475      2016a.

476 Menafoglio, A., Guadagnini, A., and Secchi, P.: Stochastic simulation of soil particle-size curves in heterogeneous aquifer
477      systems   through   a   Bayes   space   approach,   Water   Resources   Research,   52,   5708-5726,
478      https://doi.org/10.1002/2015wr018369, 2016b.

479 Molayemat, H., Torab, F. M., Pawlowsky-Glahn, V., Morshedy, A. H., and Jose Egozcue, J.: The impact of the compositional
480      nature of data on coal reserve evaluation, a case study in Parvadeh IV coal deposit, Central Iran, International Journal of
481      Coal Geology, 188, 94-111, https://doi.org/10.1016/j.coal.2018.02.003, 2018.

482 Nelder, J. A., and Wedderburn, R. W. M.: Generalized linear models, Journal of the Royal Statistical Society. Series A (General),
483      135, 370-384, https://doi.org/10.2307/2344614, 1972.

484 Nickel, S., Hertel, A., Pesch, R., Schroeder, W., Steinnes, E., and Uggerud, H. T.: Modelling and mapping spatio-temporal
485      trends of heavy metal accumulation in moss and natural surface soil monitored 1990-2010 throughout Norway by
486      multivariate   generalized   linear   models   and   geostatistics,   Atmospheric   Environment,   99,   85-93,
487      https://doi.org/10.1016/j.atmosenv.2014.09.059, 2014.

488 Odeh, I. O. A., McBratney, A. B., and Chittleborough, D. J.: Further results on prediction of soil properties from terrain
489      attributes: heterotopic cokriging and regression-kriging, Geoderma, 67, 215-226, https://doi.org/10.1016/0016-
490      7061(95)00007-B, 1995.

491 Pawlowsky-Glahn, V.: On spurious spatial covariance between variables of constant sum, 107-113 pp., 1984.

492 Pawlowsky-Glahn V, Egozcue JJ, Tolosana-Delgado R.: Modeling and analysis of compositional data. John Wiley & Sons,

493     Ltd, 2015.

494 R Development Core Team: R: A language and environment for statistical computing, in, R Foundation for Statistical
495     Computing, Vienna, Austria, 2019.

496 Scarpone, C., Schmidt, M. G., Bulmer, C. E., and Knudby, A.: Modelling soil thickness in the critical zone for Southern British
497     Columbia, Geoderma, 282, 59-69, https://doi.org/10.1016/j.geoderma.2016.07.012, 2016.

498 Segal, M. and Xiao, Y. Y.: Multivariate random forests,Wiley Interdisciplinary Reviews-Data Mining and Knowledge
499     Discovery, 1, 80–87, https://doi.org/10.1002/widm.12, 2011.

500 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
501     by geographically weighted regression: A case study in the Heihe River Basin, China, Geoderma, 261, 11-22,
502     https://doi.org/10.1016/j.geoderma.2015.06.024, 2016.

503 Talska, R., Menafoglio, A., Machalova, J., Hron, K., and Fiserova, E.: Compositional regression with functional response,
504     Computational Statistics & Data Analysis, 123, 66-85, 10.1016/j.csda.2018.01.018, 2018.

505 Tolosana-Delgado, R., Otero, N., Pawlowsky-Glahn, V., and Soler, A.: Latent compositional factors in the Llobregat River
506     Basin (Spain) hydrogeochemistry, Mathematical Geology, 37, 681-702, https://doi.org/10.1007/s11004-005-7375-7, 2005.

507 Venables, W. N., and Dichmont, C. M.: GLMs, GAMs and GLMMs: an overview of theory for applications in fisheries research,
508     Fisheries Research, 70, 319-337, https://doi.org/10.1016/j.fishres.2004.08.011, 2004.

509 Walvoort, D. J. J., and de Gruijter, J. J.: Compositional Kriging: A spatial interpolation method for compositional data,
510     Mathematical Geology, 33, 951-966, https://doi.org/10.1023/a:1012250107121, 2001.

511 Wang, Z., and Shi, W. J.: Mapping soil particle-size fractions: A comparison of compositional kriging and log-ratio kriging, J.
512     Hydrol., 546, 526-541, https://doi.org/10.1016/j.jhydrol.2017.01.029, 2017.

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

515 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
516     regression tree and random forest models for mapping topsoil organic carbon concentration in an alpine ecosystem,
517     Ecological Indicators, 60, 870-878, https://doi.org/10.1016/j.ecolind.2015.08.036, 2016.

518 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
519     Pedologica Sinica, 52, 220-227, 2015.

520 Yue, T., Liu, Y., Zhao, M., Du, Z., and Zhao, N.: A fundamental theorem of Earth's surface modelling, Environ. Earth Sci., 75,
521     751, https://doi.org/10.1007/s12665-016-5310-5, 2016.

522 Yue, T., Zhao, N., Liu, Y., Wang, Y., Zhang, B., Du, Z., Fan, Z., Shi, W., Chen, C., Zhao, M., Song, D., Wang, S., Song, Y.,
523     Yan, C., Li, Q., Sun, X., Zhang, L., Tian, Y., Wang, W., Wang, Y. a., Ma, S., Huang, H., Lu, Y., Wang, Q., Wang, C., Wang,
524     Y., Lu, M., Zhou, W., Liu, Y., Wang, Z., Bao, Z., Zhao, M., Zhao, Y., Rao, Y., Naseer, U., Fan, B., Li, S., Yang, Y., and
525     Wilson, J. P.: A fundamental theorem for eco-environmental surface modelling and its applications, Science China-Earth
526     Sciences, 63, 1092-1112, https://doi.org/10.1007/s11430-019-9594-3, 2020.

527    Zhang, M., Shi, W., and Xu, Z.: Systematic comparison of five machine-learning models in classification and interpolation of

528        soil particle size fractions using different transformed data, Hydrol. Earth Syst. Sci., 24, 2505-2526,

529        https://doi.org/10.5194/hess-24-2505-2020, 2020.