



1 Performance of three machine learning algorithms for predicting soil 2 organic carbon in German agricultural soil

3 Ali Sakhaee¹, Anika Gebauer², Mareike Ließ², Axel Don¹

4 ¹Thünen Institute of Climate Smart Agriculture, Braunschweig, Germany

5 ²Department Soil System Science, Helmholtz Centre for Environmental Research – UFZ, Halle (Saale), Germany

6

7 *Correspondence to:* Ali Sakhaee (a.sakhaee@thuenen.de)

8 **Abstract**

9 Soil organic carbon (SOC), as the largest terrestrial carbon pool, has the potential to influence climate change and
10 mitigation, and consequently SOC monitoring is important in the frameworks of different international treaties.
11 There is therefore a need for high resolution SOC maps. Machine learning (ML) offers new opportunities to do
12 this due to its capability for data mining of large datasets. The aim of this study, therefore, was to test three
13 commonly used algorithms in digital soil mapping – random forest (RF), boosted regression trees (BRT) and
14 support vector machine for regression (SVR) – on the first German Agricultural Soil Inventory to model
15 agricultural topsoil SOC content. Nested cross-validation was implemented for model evaluation and parameter
16 tuning. Moreover, grid search and differential evolution algorithm were applied to ensure that each algorithm was
17 tuned and optimised suitably. The SOC content of the German Agricultural Soil Inventory was highly variable,
18 ranging from 4 g kg⁻¹ to 480 g kg⁻¹. However, only 4% of all soils contained more than 87 g kg⁻¹ SOC and were
19 considered organic or degraded organic soils. The results show that SVR provided the best performance with
20 RMSE of 32 g kg⁻¹ when the algorithms were trained on the full dataset. However, the average RMSE of all
21 algorithms decreased by 34% when mineral and organic soils were modelled separately, with the best result from
22 SVR with RMSE of 21 g kg⁻¹. Model performance is often limited by the size and quality of the available soil
23 dataset for calibration and validation. Therefore, the impact of enlarging the training data was tested by including
24 1223 data points from the European Land Use/Land Cover Area Frame Survey for agricultural sites in Germany.
25 The model performance was enhanced for maximum 1% for mineral soils and 2% for organic soils. Despite the
26 capability of machine learning algorithms in general, and particularly SVR, in modelling SOC on a national scale,
27 the study showed that the most important to improve the model performance was separate modelling of mineral
28 and organic soils.

29 **1 Introduction**

30 Soil organic carbon (SOC) is the largest terrestrial carbon pool (Wang et al., 2020) and plays an essential role in
31 agriculture. Since SOC influences various physical, chemical and biological properties of soil (Reeves, 1997),
32 numerous studies recognise it as a crucial indicator of soil quality (Castaldi et al., 2019; Meersmans et al., 2012;
33 Reeves, 1997). Thus, its decline is identified as a threat that leads to soil degradation (Castaldi et al., 2019; Poeplau
34 et al., 2020). Moreover, when considering carbon sequestration, the SOC pool provides the option for climate
35 change mitigation (Meersmans et al., 2012; Ward et al., 2019). Consequently, SOC monitoring is important in the
36 frameworks of various international treaties such as the European Union Soil Thematic Strategy and the United
37 Nations Framework Convention on Climate Change (Meersmans et al., 2012; Poeplau et al., 2020). There is
38 therefore growing interest in understanding the spatial distribution of SOC at different scales in response to



39 increasing demand for a better assessment of SOC (Minasny et al., 2013). This is particularly important for
40 agricultural land due to its potential for carbon sequestration (Lal, 2004).

41 In digital soil mapping (DSM), a soil attribute is formulated as an empirical quantitative function of seven factors:
42 soil properties, climate, organism, topography, parent material, time and spatial position (McBratney et al., 2003).
43 Therefore, this function, known as the SCORPAN model, can be applied to spatially predict the soil attribute of
44 interest (Minasny et al., 2013). Within this framework, machine learning algorithms aim to automatically extract
45 the information from the data for predictive purposes (Behrens et al., 2005). This is particularly intriguing in view
46 of the expansion of databases at a different scale in soil science and the complexity of the covariates in recent years
47 (McBratney et al., 2003; Wadoux et al., 2020), thus making DSM cost-effective, time-efficient and applicable over
48 large areas with good results (Behrens and Scholten, 2006; Camera et al., 2017).

49 Despite the advantages of DSM, it is crucial to consider that its application requires soil databases of an adequate
50 sample size for training and testing. Furthermore, consistent and quality-checked datasets are a prerequisite for
51 DSM. Several soil inventories and monitoring networks for SOC have been formed on a national scale in countries
52 such as Sweden (Poeplau et al., 2015), France (Belon et al., 2012; Meersmans et al., 2012), Denmark (Taghizadeh-
53 Toosi et al., 2014) and Scotland (Chapman et al., 2013). Nonetheless, the most critical shortcomings of soil
54 inventories in Germany are the lack of a large-scale, high-quality SOC inventory (Wiesmeier et al., 2012) with a
55 periodic and standardised sampling focus on agricultural soils (Prechtel et al., 2009). These issues have now been
56 solved in the first German Agricultural Soil Inventory (Poeplau et al., 2020). This inventory was conducted on a
57 national scale with a sampling depth down to 100 cm at 3104 sampling sites covering agricultural land.
58 Furthermore, on a European scale, the Land Use/Land Cover Area Frame Survey (LUCAS) undertaken in 2009 is
59 the first harmonised topsoil survey with physico-chemical analyses of georeferenced topsoil samples in 23
60 European states (Tóth et al., 2013). Therefore, by taking advantage of DSM and both the German Agricultural Soil
61 Inventory and the LUCAS survey, it is possible to regionalise single-point measurements to obtain complete high-
62 resolution cover soil data and thus provide a baseline for SOC monitoring as well as for environmental and climatic
63 modelling for Germany.

64 Boosted regression trees (BRT), random forest (RF) and support vector machine for regression (SVR) are among
65 the most used algorithms in DSM (Padarian et al., 2020). For example, Martin et al. (2014) predicted topsoil SOC
66 on a national scale for France using the BRT algorithm and compared its results when the same algorithm was
67 coupled with a geostatistical approach. They concluded that since spatial autocorrelation is not feasible in most
68 national inventories, the BRT algorithm alone is sufficient for this purpose. This algorithm was also used on a
69 national scale in China for data from the 1980s and 2010s in order to predict topsoil SOC and its spatial-temporal
70 change, as well as the main drivers of its variability (Wang et al., 2021). Moreover, RF has become more popular
71 in DSM due to its relative simplicity and performance. For example, this algorithm was implemented to map
72 topsoil SOC on a national scale in Madagascar and obtain its main drivers (Ramifehiarivo et al., 2017).
73 Ramifehiarivo et al. (2017) concluded that the uncertainty of the algorithm was lower when compared with the
74 maps formerly generated for the country. Moreover, this algorithm was compared with the Cubist model for
75 mapping SOC at different resolutions on a regional scale in China and could outperformed it (Li et al., 2021).
76 Fewer studies have used SVR to predict SOC than RF. Studies have mainly implemented SVR on a regional scale
77 with a limited number of samples (Forkuor et al., 2017; Were et al., 2015) or on a national scale (Switzerland)
78 with very few samples (150 samples from the European LUCAS survey) (Zhou et al., 2021). However, in a study



79 comparing different algorithms, including SVR and RF, on a continental scale and within each country in Latin
80 America, the results indicated that the best-performing algorithm varied in different countries (Guevara et al.,
81 2018). This difference mainly depended on sample dispersity and also country size which affects the heterogeneity
82 of land use and environmental conditions.

83 Another important consideration when applying machine learning is the impact of the parameter-tuning strategy
84 in algorithm performance. This is particularly crucial when the objective of the study is the comparisons of
85 different machine learning algorithms. Although some algorithms are less sensitive to tuning, this step is more
86 important for others, particularly those with a higher number of parameters (Tziachris et al., 2020; Wadoux et al.,
87 2020). Furthermore, as algorithms differ by the type of their parameters, continuous or discrete, the chosen strategy
88 should be in accordance with this difference. This is particularly more important for algorithms with continuous
89 parameters. For example, it has been shown that that the performance of SVR and BRT is better and more stable
90 when optimised by differential evolution (DE) algorithm than tuned by grid search (Zhang et al., 2011; Gebauer
91 et al., 2020). Despite this importance, in a review of studies that have applied DSM, Wadoux et al. (2020) state
92 that almost half of them implemented parameter tuning, with grid search the most common strategy for this
93 purpose. This finding indicates that the role of parameter tuning and optimisation is unfortunately undermined in
94 DSM. This is particularly evident when the application of machine learning in this field is compared with other
95 fields, where various studies have shown the impact of parameter-tuning strategies on the performance of
96 algorithms such as SVR and BRT (Liang et al., 2011; Santos et al., 2021; Bhadra et al., 2012; Deng et al., 2019).

97 The aim of the present study was therefore i) to address the above-mentioned parameter-tuning issue and
98 consequently provide a true comparison of the performance of BRT, RF, and SVR in modelling the SOC contents
99 of German agricultural topsoil (0-30 cm), ii) to assess the impact of training data size by extending the data of the
100 German Agricultural Soil Inventory with LUCAS data for model calibration, and iii) to develop a two-model
101 approach to address the high variability of SOC in German agricultural soils and compare it with a single-model
102 approach.

103 **2 Materials and methods**

104 **2.1 Soil data**

105 The models were built using SOC content data from two soil inventories. The first dataset was from the German
106 Agricultural Soil Inventory, which consists of 3104 sites with a fixed grid of 8x8 km throughout Germany (Poeplau
107 et al., 2020). The sites were sampled and analysed for different soil properties, including SOC content measured
108 via dry combustion, for the upper 30 cm of the soil between 2012 and 2018. The second dataset was the European
109 LUCAS survey that provides SOC content, similarly measured via dry combustion, for all EU countries, with the
110 sampling depth limited to 0-20 cm (Tóth et al., 2013). Therefore, in order to harmonise the depths of both datasets,
111 these were subdivided into mineral and organic soil classes according to a SOC threshold value of 87.0 g kg⁻¹
112 considering all soils above this threshold as organic soils comprising peat soils and disturbed and degraded peat
113 soils (Poeplau et al., 2020). Linear correlation functions between 0-30 cm and 0-20 cm were derived for each soil
114 class of the German Agricultural Soil Inventory separately. These functions were applied to the corresponding soil
115 class from the LUCAS data in order to estimate 0-30 cm topsoil SOC. With a slope of 0.881 for mineral soils and
116 1.02 for organic soils, they changed the mean of LUCAS data by less than 6%. The depth-extrapolated values of



117 mineral and organic soils were then combined to form the complete dataset. The 0-30 cm LUCAS data generated
118 and the original 0-20 cm LUCAS data were then used by each algorithm to check the effect of depth extrapolation.

119 **2.2 Covariates**

120 Covariates from multiple sources were included to approximate the SCORPAN factors throughout Germany. In
121 case of multiple data products for one covariate, the one with the best quality (least artefacts), and the highest
122 spatial resolution was added. These were then resampled in ArcGIS (ESRI, 2013) using the INSPIRE standard
123 grid at 100 m resolution (Eurostat grid generation tool for ArcGIS). The resampling method was either the nearest
124 neighbour for categorical covariates or bilinear interpolation for continuous covariates. The same INSPIRE grid
125 was used to rasterise the vector covariates as well. Finally, they were stacked and overlaid on SOC databases in
126 order to extract the values of sampling points.

127 Following the SCORPAN framework, 24 covariates including x and y spatial positions were compiled. In order to
128 capture climate factors (C factor), precipitation (DWD, 2018c), sunshine duration (DWD, 2017), summer days
129 (DWD, 2018b) and minimum temperature (DWD, 2018a) were used according to the study of Schneider et al.
130 (2021). Using principal component analysis, these four covariates were indicated to be the most important among
131 34 available climate factors for SOC in the German Agricultural Soil Inventory dataset. Moreover, land use is one
132 of the main drivers of SOC variability at a national scale (Poeplau et al., 2020). Thus, the land-use map from the
133 official topographic information system (BKG, 2019) with its corresponding classes according to the German
134 Agricultural Soil Inventory was rasterised and included. This is a categorical covariate, representing the organism
135 factor of SCORPAN (O factor), that distinguishes croplands from grasslands and captures their spatial distribution
136 throughout Germany.

137 The European Digital Elevation Model (EUEM) (European Union Copernicus Land Monitoring Service, 2016)
138 and six covariates derived from this layer were also added to integrate the topography and relief parameters (R
139 factor). Slope, plan curvature and profile curvature, generated on SAGA (Conrad et al., 2015), were included to
140 capture the slope's gradient, convexity-concavity and convergence-divergence. These factors influence the soil
141 distribution throughout the landscape, e.g. affecting flow over the surface, thus impacting SOC and its dynamic
142 (Ritchie et al., 2007). Moreover, north-south and east-west aspects were obtained from EUEM as these influence
143 soil development and subsequently affect SOC (Carter and Ciolkosz, 1991). The Topographic Wetness Index
144 (TWI), generated on SAGA (Conrad et al., 2015), was also added since it captures the soil moisture distribution
145 of the landscape and has a direct correlation with SOC (Pei et al., 2010). A geomorphographic map of Germany
146 containing 25 geomorphic categories was also used to distinguish between four different landscape areas of the
147 country: North German lowlands, highlands, Alpine foothills and the Alps.

148 Continuing with the framework, a large-scale soil landscape unit map ("Bodengrosslandschaft") (Richter et al.,
149 2007) comprising 38 classes was used. This covariate divides Germany by various geo-factors that can be compiled
150 into a map with 12 soil regions representing mainly the parent materials. Similarly, large-scale soil-climate region
151 map ("Bodenklima") (Roßberg, 2007) with 50 classes was added. Moreover, Germany's hydrogeological unit map
152 (BGR and SGD, 2019) provides information about lithology and its hydrological characteristics. These categorical
153 maps were rasterised and applied to the model as the P factor of SCORPAN. Moreover, the soil factor of the
154 framework (S factor) was captured by eight covariates that represent different aspects of its properties: the map of
155 organic soils (Roßkopf et al., 2015) that distinguishes mineral soils from organic ones and explains their spatial



156 distribution throughout the country, as well as the map of nitrogen (Ballabio et al., 2019) and clay content (Ballabio
157 et al., 2016) since they directly correlate with SOC. As nitrogen is a crucial component of soil organic matter,
158 regions with higher total nitrogen have higher SOC (Ballabio et al., 2019). Also for clay content, different studies
159 have shown that coarser soil textures tend to have a lower accumulation of SOC (Zhong et al., 2018; Hoyle et al.,
160 2011). Map of pH (Ballabio et al., 2019) since soil pH directly impacts microbial activities that influence the
161 turnover of soil organic matter, and consequently negatively correlates with SOC (Malik et al., 2018). Furthermore,
162 map of available water capacity (Ballabio et al., 2016) as this soil properties is another interactive factor with SOC
163 through plant productivity and soil texture (Burke et al., 1989; Yu et al., 2021). Soil erosion is also a key factor in
164 the SOC cycle (Li et al., 2019), which was added through the soil erosion map of Europe (Borrelli et al., 2018).

165 **2.3 Boosted Regression Trees**

166 Developed by Friedman et al. (2000), BRT is a tree-based algorithm that applies boosting method to improve
167 accuracy. Boosting method relies on combining several approximate prediction models rather than obtaining one
168 single highly accurate one (Schapire, 2003). Thus, the decision trees are grown sequentially so that each decision
169 tree predicts the residual of the previous one. Consequently, the number of trees influences the performance of the
170 algorithm and requires tuning. However, to incorporate randomness into the model and subsequently increase the
171 robustness of performance, the trees are grown on a randomly selected data subset with no replacement (Friedman,
172 2002). The size of this subset is controlled by a parameter known as a bag fraction. Furthermore, the contribution
173 of each new tree to the final model is regularised by learning rate, also known as shrinkage (Friedman et al., 2009).
174 Finally, the number of splits in each tree that divides the response variable into subsets is optimised by interaction
175 depth. The BRT model was built in R using the “gbm” package (Greenwell et al., 2019).

176 **2.4 Random Forest**

177 Similar to BRT, RF is another tree-based algorithm. RF uses bootstrap sampling of the dataset for growing a
178 decision tree. Subsequently, by aggregating the results of a large number of decision trees, the bias and variance
179 of the final model can be reduced (Breiman, 1999). The method of bootstrapping in conjunction with aggregating,
180 known as bagging, increases the robustness and stability of RF. However, the trees from different bootstraps may
181 form a similar structure if all covariates participate in a split of each node. Thus, the variance cannot be reduced
182 optimally through the bagging process (Kuhn and Johnson, 2013). In order to avoid this tree correlation, a random
183 subset of predictors is selected at each split. The parameter m_{try} defines the number of predictors included in this
184 subset and should be tuned (Kuhn and Johnson, 2013). The RF algorithm was implemented by the “Ranger”
185 package (Wright and Ziegler, 2017) in R.

186 **2.5 Support Vector Regression**

187 SVR is a form of support vector machine adopted for regression. From all possible solutions, i.e. estimation
188 function, for the problem, SVR tries to obtain an estimation function with the maximum ϵ error while minimising
189 model complexity (Smola and Schölkopf, 2004). Thus, a symmetrical tolerance threshold, ϵ -insensitivity zone, is
190 created around the estimation function within which the vectors are not penalised (Awad and Khanna, 2015).
191 However, the vectors that lie on the boundary of the ϵ -insensitivity zone are called support vectors. Therefore, ϵ
192 is an optimisable parameter that controls the width of ϵ -insensitivity, alters the model complexity and impacts the
193 number of support vectors inversely (Cherkassky and Ma, 2004). Moreover, the trade-off between model
194 complexity and tolerance of ϵ deviation is controlled by a parameter named C (Smola and Schölkopf, 2004;



195 Cherkassky and Ma, 2004). Optimising the C parameter has a crucial impact on SVR performance since a high C
196 can lead to overfitting, while a low C can cause under fitting (Kuhn and Johnson, 2013). The use of kernel functions
197 makes SVR a powerful tool for nonlinear problems. By implementing these functions, SVR can map data space
198 from its original dimension to a higher dimensional space where a nonlinear problem can be solved linearly. In
199 this study, the Radial Basis Function (RBF) kernel was used with gamma as its tuneable parameter. This parameter
200 affects the generalisation performance of SVR by controlling the influence of support vectors inversely (Battinini
201 et al., 2019). SVR was implemented from the package e1071 in R (Hornik et al., 2021).

202 2.6 Performance evaluation

203 When training a predictive model, it is important to evaluate its generalisation performance on unseen data of the
204 same type (Hawkins et al., 2003). However, as the number of available samples is usually a limiting factor, the
205 evaluation process is often done by randomly splitting the available dataset into training and testing sets multiple
206 times, i.e. cross-validation (CV). Although this process is effective, it is not entirely immune from biased
207 estimation of error. However, to ensure that the estimated error in model evaluation is as unbiased as possible,
208 every model training step should be performed within the CV. This includes finding the best parameter sets for the
209 chosen algorithm (Varma and Simon, 2006). Thus, the algorithms in this study were applied on a stratified nested
210 CV.

211 First, to ensure that the SOC distribution was represented in the CV scheme, Germany was divided using a 100x100
212 km INSPIRE grid into 50 strata. Random samples from each stratum were then taken and compiled into a fold.
213 This procedure was continued to create five folds and was repeated five times, forming the outer loop of CV used
214 for model evaluation. Large distance between neighboring samples, 8120 m on average, prevents train and test
215 data from being spatially autocorrelated. Since the aim was to tune the parameters of the algorithms, the training
216 set of the outer loop of CV was nested, creating five folds as the inner loop on which the parameter tuning was
217 performed. To evaluate the performance of algorithms, root-mean-squared error (RMSE), Eq. 1, mean absolute
218 error (MAE), Eq. 2, and mean absolute percentage error (MAPE), Eq. 3, were used.

$$219 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (P_i - O_i)^2} \quad (1)$$

$$220 \quad MAE = \frac{1}{n} \sum_{i=1}^n (P_i - O_i) \quad (2)$$

$$221 \quad MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{P_i - O_i}{O_i} \right| \times 100 \quad (3)$$

222 Where n is the number of samples, P_i and O_i are the predicted and observed values, respectively.

223 2.6.1 Parameter tuning

224 As mentioned in Sect.1, choosing a suitable strategy for parameter tuning is a crucial step in machine learning
225 particularly for comparing the algorithms. Therefore, two strategies were applied depending on the algorithm: 1)
226 a grid search for RF and 2) optimisation with the DE algorithm for BRT and SVR. The first strategy was an
227 exhaustive search over a defined space consisting of lower bound, upper bound and n steps in between for the
228 target parameter. Therefore, the target parameters in this strategy should be discrete or discretised beforehand if
229 they are continuous (Probst et al., 2019). This strategy was applied to RF since the tuning parameter is discrete.
230 However, the second strategy is a stochastic approach of searching over a continuous space in order to solve an



231 optimisation problem (Qin et al., 2009) and is described in more detail by Storn & Price (1997). Therefore, SVR
232 and BRT were optimised by this strategy as they have continuous parameters. For the optimisation task in the
233 present study, the R package “DEoptim” was applied (Peterson et al., 2021). Table S1 shows the parameters and
234 their tuning range for each algorithm.

235 **2.6.2 Variable importance**

236 Variable importance was assessed by permutation (Ließ et al., 2021). Therefore, each covariate in the test set was
237 shuffled 10 times and on each occasion the trained model corresponding to that test set was applied. The population
238 of RMSE was averaged and its relative change to the RMSE of the original test set was calculated. Thus, the
239 variable importance of each covariate in terms of percentage relative change in RMSE was obtained.

240 **2.7 Modelling approaches**

241 Three approaches were designed to test the performance of the algorithms. The models were built based on nested
242 CV, while the train and test sets remained identical for the three algorithms to make the results comparable. The
243 first approach (AP1) only used the SOC content from the German Agricultural Soil Inventory and corresponding
244 values from the covariates were used to build the models. Thus, the dataset was cross-validated and used by BRT,
245 RF and SVR to predict the SOC content of German agricultural soils. The results of this approach served as a
246 baseline on which the model improvement for each algorithm in other approaches was assessed.

247 Due to the high variability of SOC in agricultural soils of Germany, two separate models for organic and mineral
248 soils was developed and tested to identify whether it could improve model performance. Accordingly, the German
249 Agricultural Soil Inventory was subdivided by the threshold 87 g kg^{-1} into mineral and organic soils and then were
250 used to train separate models. This approach was named AP2. The same nested CV procedure was applied for both
251 data subsets. The results of BRT, RF and SVR were compared to identify which one had better performance under
252 mineral and organic soils separately. Finally, each algorithm's predicted SOC from two separate models was
253 combined, and the error metrics were calculated for the full data set to identify the impact of AP2 on model
254 performance. The CV folds for this procedure match the one from the AP1 models.

255 The impact of enlarging the training set on model performance was examined for both AP1 and AP2 approaches.
256 Thus, 1223 depth-extrapolated samples of the LUCAS data were added to the training sets of AP1 and named
257 AP1L. Moreover, the same threshold (87 g kg^{-1}) was used to subdivide this dataset and each soil class was included
258 to the training set of the corresponding soil class of AP2 and named AP2L. The test sets of the CV procedure
259 remained the same.

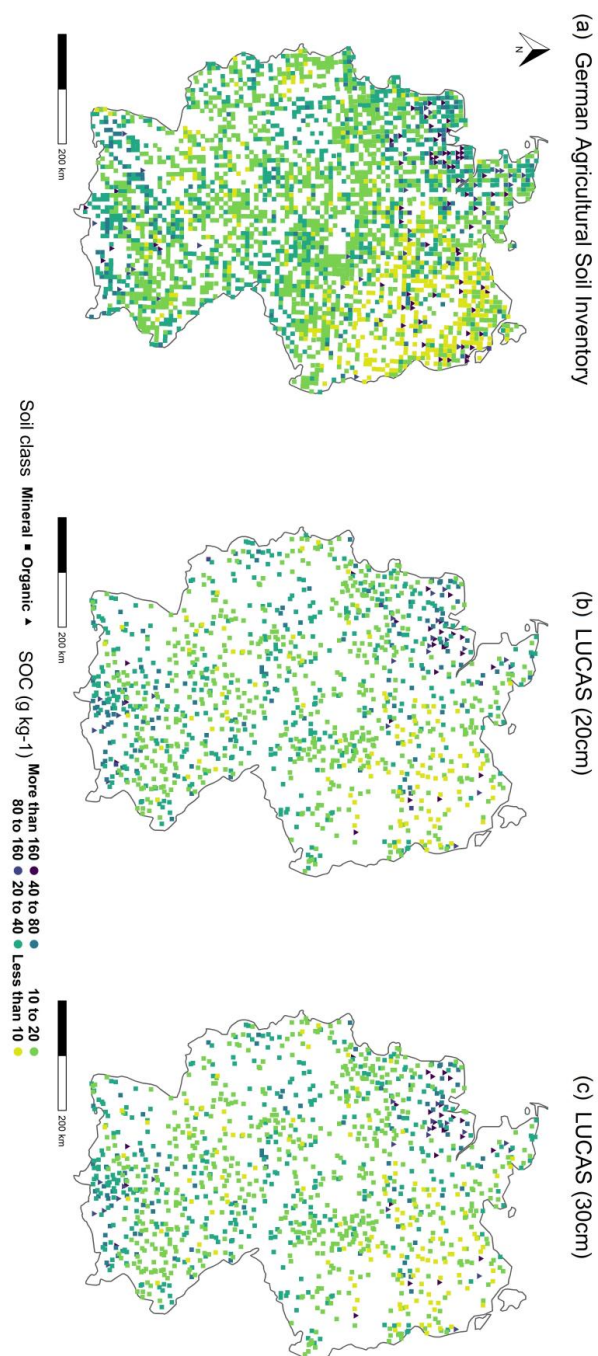
260 **3 Results and Discussion**

261 **3.1 Comparison of algorithms on the data from the German Agricultural Soil Inventory**

262 The range of SOC content of topsoil for the German Agricultural Soil Inventory dataset was 4 g kg^{-1} to 480 g kg^{-1} ,
263 with a mean of 27 g kg^{-1} and median of 16 g kg^{-1} . Figure 1 shows the spatial distribution of the implemented
264 data. The RMSE and MAPE indicate that SVR had a better general performance than the two other algorithms
265 (Fig. 2). In this respect, the RMSE of SVR was 5% lower than that from RF and 4% lower than that from BRT.
266 Furthermore, its MAPE was 3% and 7% lower than that from RF and BRT respectively. However, despite the
267 difference in overall performance, the spatial distribution of relative residuals indicated that all three algorithms
268 were less accurate in the north of Germany compared with the centre and south of the country (Fig. 3A). This can



269 be explained by the characteristics of this region and its higher SOC variability. The northern part of Germany is
270 a lowland dominated by sandy soil texture from pleistocene sedimentation with geomorphological structures such
271 as ground moraines, terminal moraines and aprons (Roßkopf et al., 2015). Despite general geomorphological and
272 pedological similarities throughout the region, 1) organic soils in Germany are mainly located in the north and 2)
273 mineral soils with the lowest and the highest SOC content are also located in the northeast and northwest
274 respectively. Therefore, this region has the highest SOC range on agricultural soils.



275

276 Figure 1: Soil organic carbon content in topsoil of two soil inventories. A) German Agricultural Soil Inventory (0-30
277 cm), B) LUCAS at its original sampling depth (0-20 cm), C) LUCAS after depth extrapolation (0-30 cm)



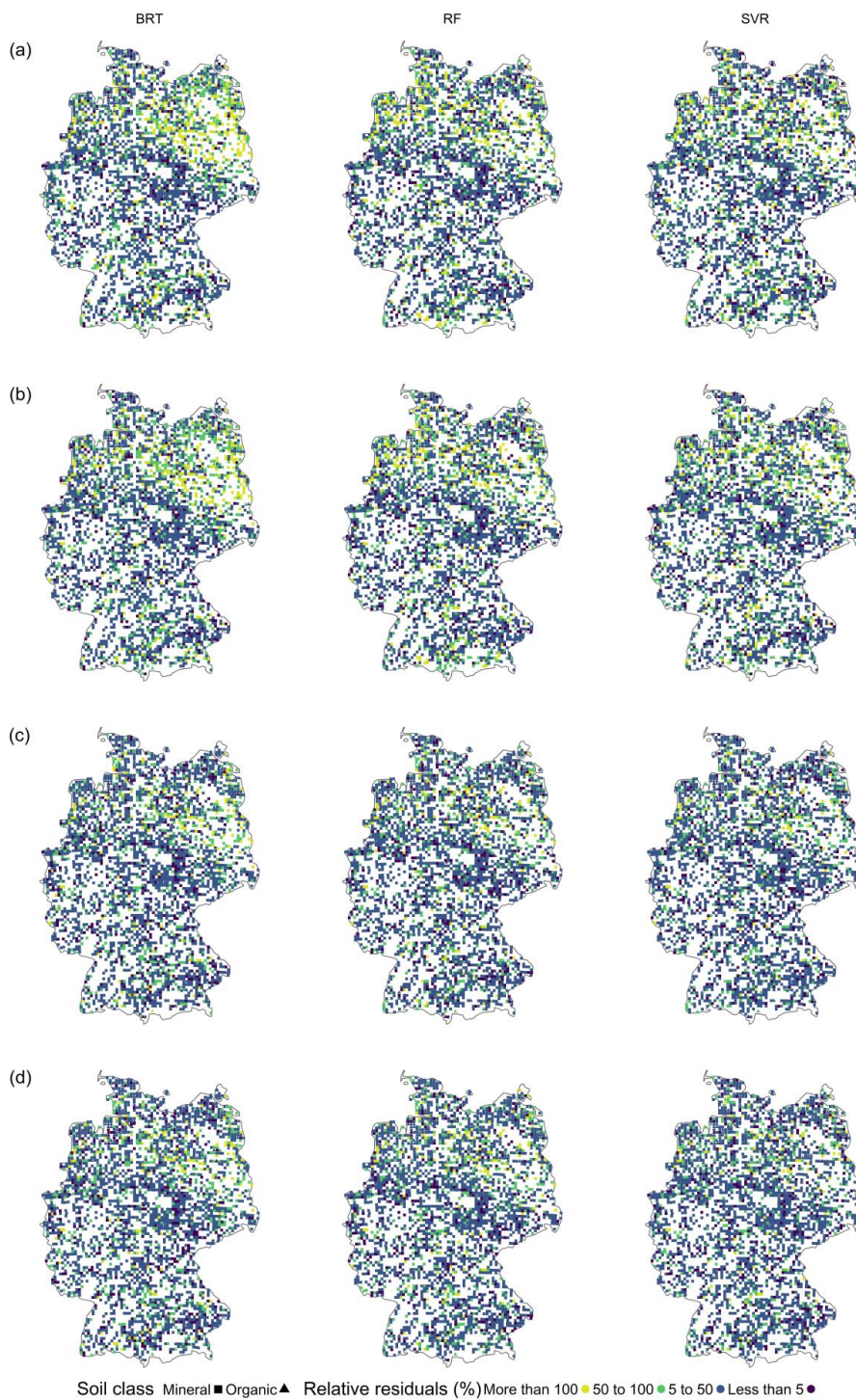
278

279 **Figure 2: Performance indicators of the three algorithms. One-model approach (Without LUCAS data API and with**
280 **LUCAS data AP1L) versus the two-model approach (AP2 and AP2L) for A) RMSE (g kg^{-1}), B) MAE (g kg^{-1}) and C)**
281 **MAPE (%). Please note that the y-axis is shortened for better visibility and does not display a zero.**

282 Consequently, the variable importance (Fig. 4A) indicated that the map of organic soils contains the highest
283 available information among all covariates for the algorithms. The value for variable importance for this covariate
284 was 65% in SVR, 72% in RF and 84% in BRT. These values firstly show the crucial role of the map of organic
285 soils for the algorithms in explaining the variability of SOC and, secondly, how BRT mainly relies on the map of
286 organic soils to predict SOC compared with SVR. Despite the importance of the organic soil map, the scatterplots
287 (Fig. 5A) show that all three algorithms underpredicted the SOC of the organic soils and had similar
288 heteroscedasticity patterns in their residuals. Thus, while most residuals from mineral soils followed the 1:1 line,
289 they became more scattered in soils with a higher SOC content. The underprediction of SOC in organic soils can
290 be explained by their low sample size, resulting in a dataset with a high SOC range and a unimodal distribution
291 that leaves these soils in the tail. Consequently, the organic soils were underrepresented and the results were
292 systematically pulled towards mineral soils, regardless of the choice of algorithm. Different studies have shown
293 that predicting soil properties with mineral and organic soils combined can lead to underprediction or
294 overprediction of one soil class, depending on the distribution of the dataset (Brogniez et al., 2015; Guio Blanco
295 et al., 2018; Mulder et al., 2016).

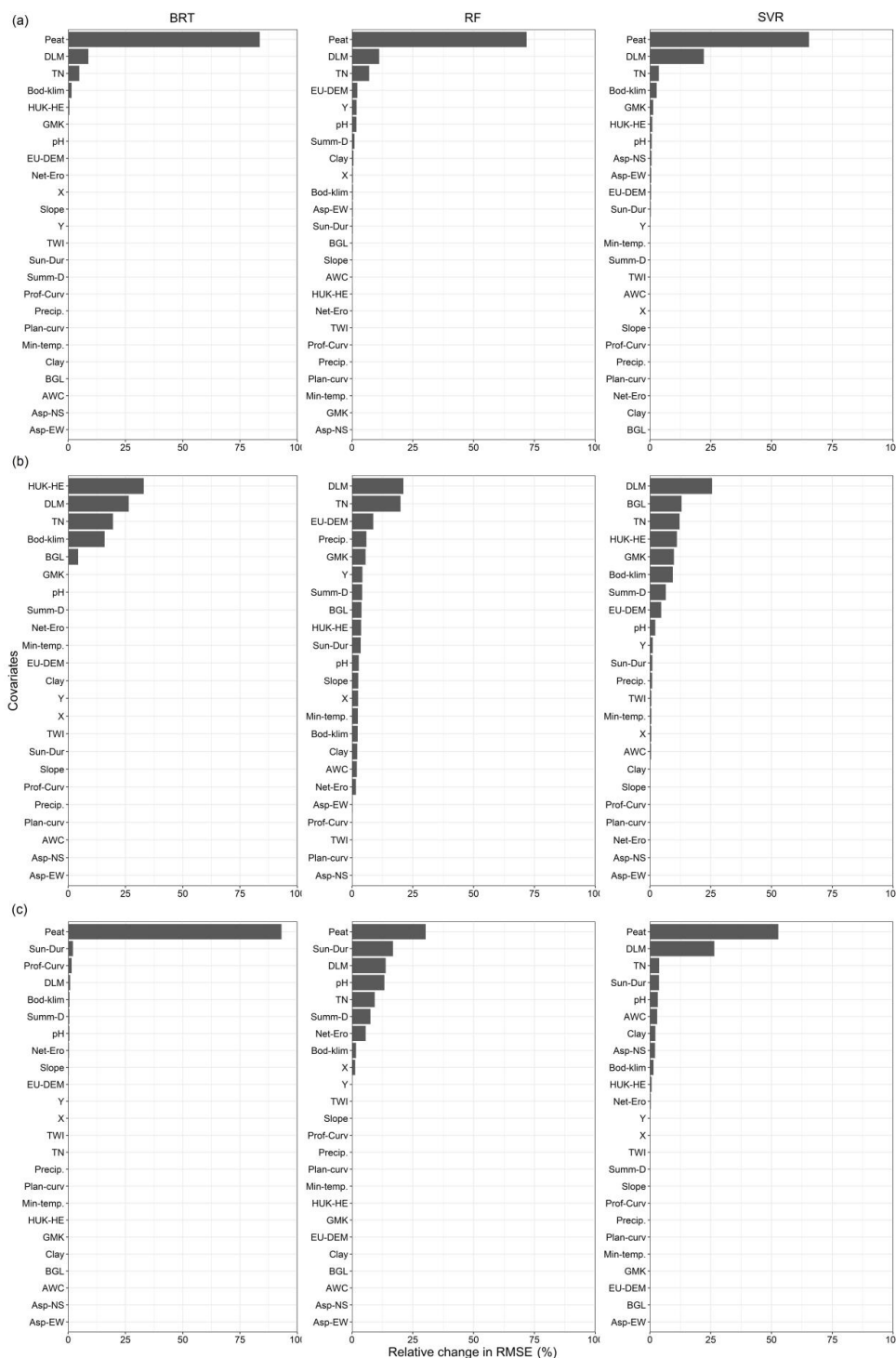


296 Although the map of organic soils was able to distinguish between the two soil classes, i.e. between mineral and
297 organic soil, it could not separate the mineral soils with a low SOC content in the northeast from those with a high
298 SOC content in the northwest. The spatial distribution of the residuals (Fig. 6A) shows that SVR and BRT
299 generally underpredicted the mineral soils in the northwest part of Germany, while RF overpredicted them.
300 Furthermore, unlike RF and SVR, BRT distinctively overpredicted SOC of the north-east's mineral soils with the
301 lowest SOC content ($<10 \text{ g kg}^{-1}$). This result indicates that the algorithms differed in their performance in mineral
302 soils. This difference was mainly due to the information they obtained from land use. As the second most important
303 covariate for all three algorithms (Fig. 4 A), the value for variable importance for this covariate was 22% in SVR,
304 but just 11% in RF and 9% in BRT. Thus, SVR exploits more information from this covariate than RF and
305 particularly BRT. Land use is one of the main drivers of SOC variability on a national scale due to the higher SOC
306 content in grasslands than in croplands (Poeplau et al., 2020). Therefore, this covariate was able to differentiate
307 between the soils of the northeast, which are under cropland, and those in the northwest as they are more under
308 grassland. Consequently, the reliance of BRT on the map of organic soils at the cost of land use could explain why
309 this algorithm overpredicted SOC in croplands in the northeast.



310

311 **Figure 3: Spatial distribution of relative residuals. A) AP1 approach, B) AP1L approach, C) AP2 approach and D)**
312 **AP2L approach.**



313

314 **Figure 4: Variable importance in terms of average relative change (%) in RMSE. A) AP1, B) mineral soil subset of**
 315 **AP2 and C) organic soil subset of AP2. The full name for each abbreviation is presented in Table S3.**



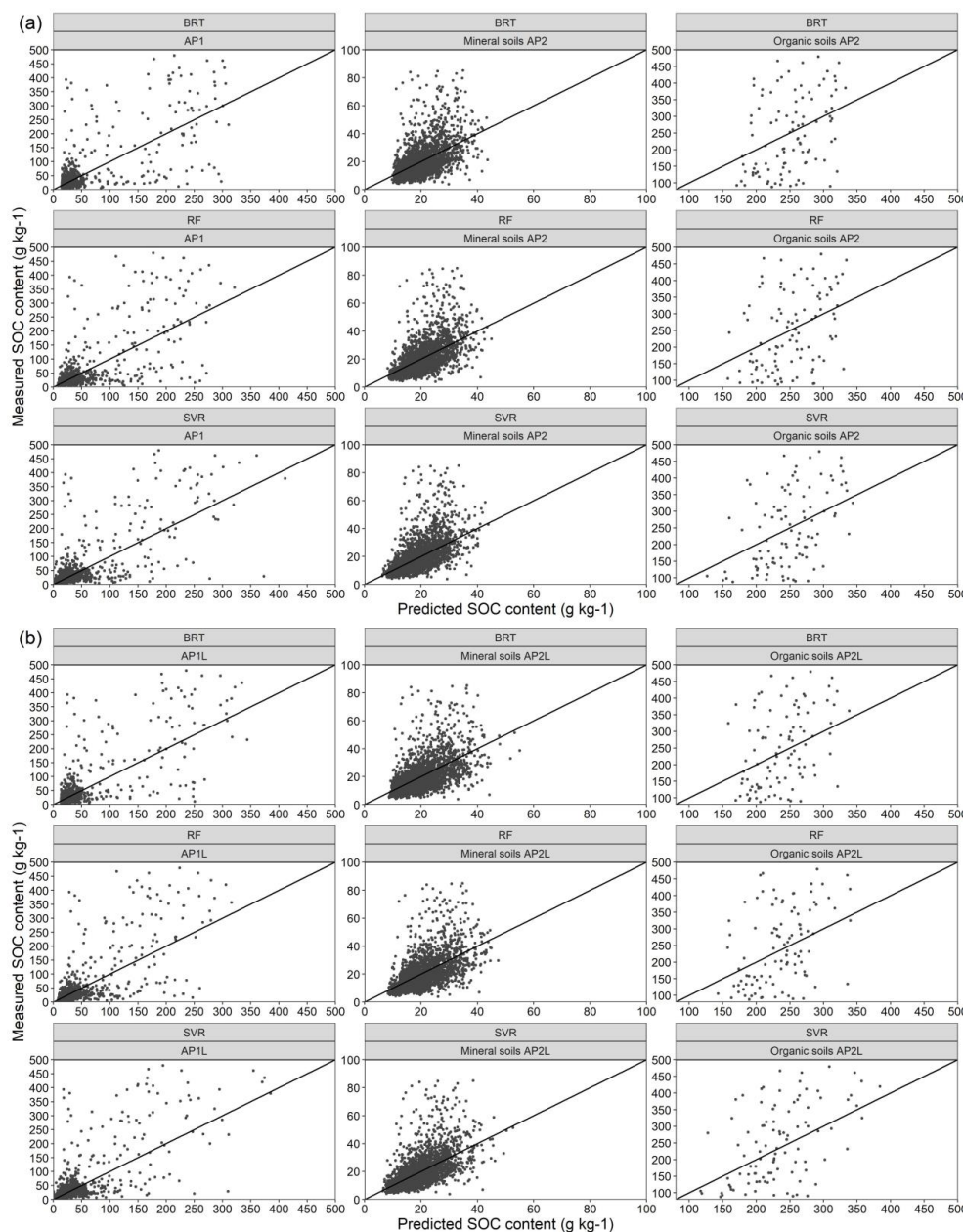
316 **3.2 Enlarging the dataset with additional soil inventories**

317 A larger soil dataset may provide additional information and consequently improve model performance. This
318 possibility was explored in the APIL approach with adding the LUCAS data. The SOC content of LUCAS data at
319 its original depth ranged from 4 g kg⁻¹ to 500 g kg⁻¹ with a mean of 30 g kg⁻¹ and a median of 18 g kg⁻¹. After
320 extrapolating the depth to 30 cm, the new range was from 5 g kg⁻¹ to 512 g kg⁻¹ with a mean of 28 g kg⁻¹ and a
321 median of 17 g kg⁻¹. The spatial distribution of LUCAS data at their original and extrapolated depth is shown in
322 Figure 1.

323 A statistical test was performed on the residuals of models built on LUCAS data with the original and extrapolated
324 depths. That was done to identify whether extrapolating the depth of LUCAS data to that of the German
325 Agricultural Soil Inventory would significantly affect model performance after their inclusion in the training set.
326 With the Shapiro-Wilk test rejecting the normality assumption of residuals of all corresponding algorithms at 20
327 cm and 30 cm, the non-parametric Kruskal-Wallis test showed no significant difference between the residuals at
328 both depths. Thus, the extrapolation of the soil depth had no significant impact on the data quality to regionalize
329 SOC. As a result, any further change in the performance of the algorithms after adding LUCAS data was due to
330 the training set being enlarged. The result of the algorithms at both depths can be found in the supplementary
331 information (Fig. S1).

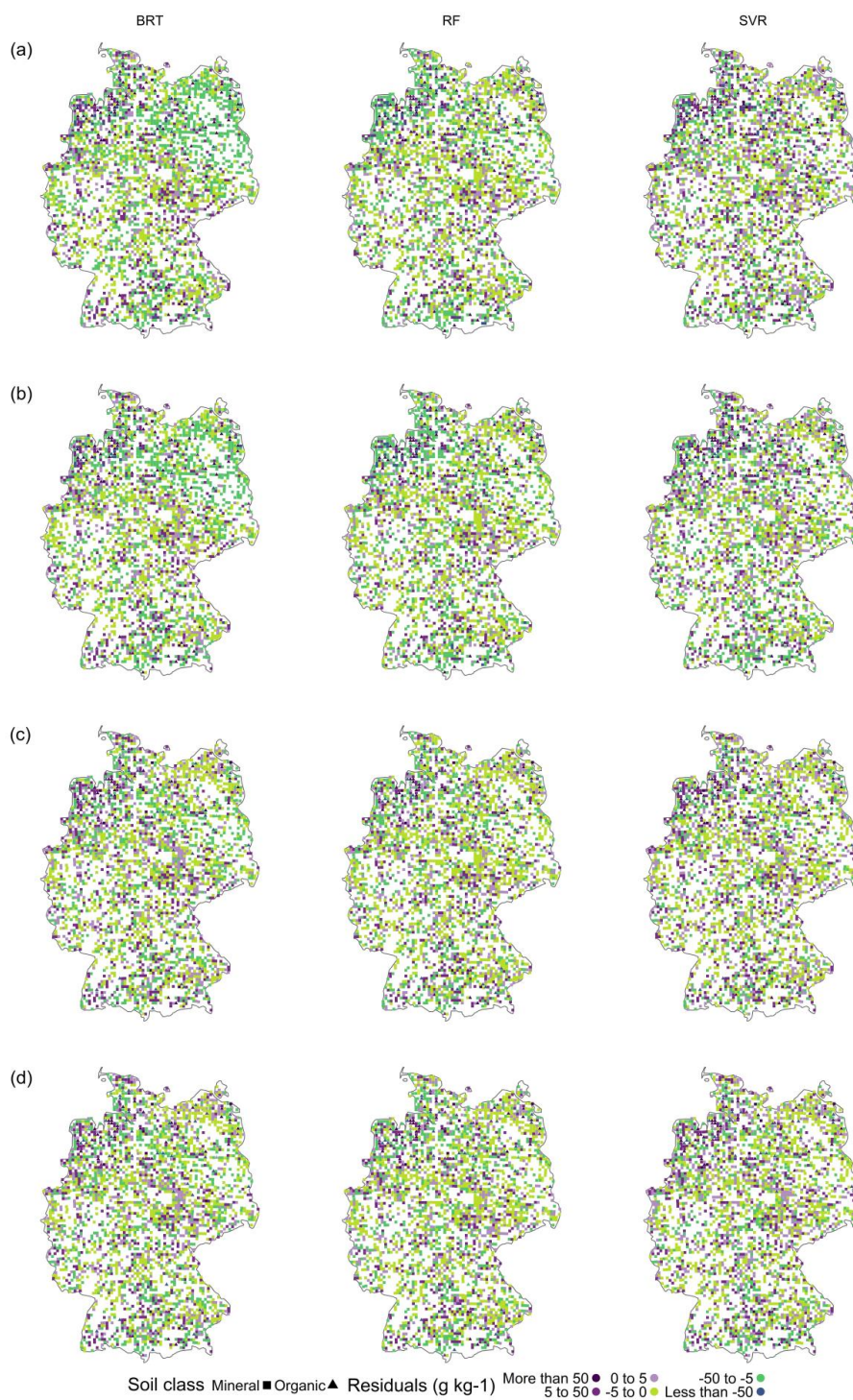
332 After enlarging the training set from 2278 to 3501 sampling points, BRT obtained the lowest RMSE and MAE
333 among the algorithms (Fig. 2). A comparison of the error metrics of corresponding algorithms from the AP1
334 approach with those from the APIL approach showed that BRT had the highest error reduction at 7% in MAPE
335 and 5% in RMSE and MAE. Furthermore, although the error metrics of RF did not improve as much as those of
336 BRT, additional training points were still beneficial for this algorithm. However, SVR did not follow any
337 systematic change under the APIL. Despite a 2% decrease in MAPE, RMSE increased by 3% and MAE remained
338 unchanged. To explore the potential explanation for this behaviour by SVR, the residuals of mineral soils were
339 separated from those of organic soils. Additional samples reduced the RMSE in mineral soils for all algorithms
340 between 9% and 13%. However, this error increased by 9% in the organic subset for SVR, while it increased by
341 just 1% for RF and even decreased by 1% for BRT. This indicated that enlarging the training set by data with
342 similar characteristics had a greater influence on systematic error of the underrepresented soil class in SVR. This
343 influence is understandable when considering the higher optimised ϵ in the APIL approach compared with that of
344 the AP1 approach. The higher value of ϵ means that the hyperplane for the training set is less complex (Cherkassky
345 and Ma, 2004) and more suitable for predicting most soil samples, i.e. mineral soils. Thus, when this hyperplane
346 was fitted to the test set identical to the AP1, the generalisation performance was hindered because it could not
347 capture the variability of samples with higher SOC values, i.e. organic soils.

348 Further evaluation revealed that regardless of the change in error metrics, the relative residuals of the three
349 algorithms had a similar spatial pattern to their counterpart from the AP1. Thus, they all showed lower accuracy
350 in the northern region of Germany for similar reasons (Fig. 3B). Moreover, the scatterplots had a similar pattern
351 with underpredicted organic soils (Fig. 5B). This confirms that when organic soils are modelled with mineral soils,
352 enlarging the training set does not provide enough information for BRT or RF to capture the high variability of
353 SOC, particularly in the north of Germany.



354

355 **Figure 5: Scatterplot of residuals. A) AP1 approach and mineral and organic soils of AP2 and B) AP1L approach and**
356 **mineral and organic soils of AP2L.**



357

358
359

Figure 6: Spatial distribution of residuals. A) AP1 approach, B) AP1L approach, C) AP2 approach and D) AP2L approach.



360 3.3 Subdividing soil inventories into mineral and organic subsets

361 As presented in the sections above, the modelling of SOC content when mineral and organic soils were combined
362 led to a systematic underprediction of soils with higher SOC values by all three algorithms, regardless of the
363 number of training samples. Therefore, by implementing the AP2 approach with two models one for mineral soils
364 and one for organic soils, a noticeable improvement in the performance of all algorithms was observed, with SVR
365 showing the best error metrics (Fig. 2). This meant 34% lower RMSE, 30% lower MAE, and 32% lower MAPE
366 than when this algorithm was trained under the AP1 approach with one model for all soils. As the high variability
367 of SOC was initially hard to capture, the subdivision of the dataset provided a range that better represented each
368 soil class. This was particularly beneficial for mineral soils (ranging from 4 g kg⁻¹ to 85 g kg⁻¹) since the number
369 of samples did not reduce drastically (only by 99 samples). Thus, the algorithms could better capture the relation
370 between SOC and covariates. Consequently, the overall performance improved when the underrepresented soil
371 class was modelled separately. This is in line with the study of Rawlins et al. (2009) which recommends the
372 separate modelling of mineral and organic soils.

373 Nonetheless, following the AP2L approach with additional data, the RMSE and MAPE of the algorithms improved
374 by less than 2% compared with AP2. However, the greatest change was observed in the MAE of SVR with a 2%
375 improvement. Therefore, additional training samples did not considerably influence the performance since the
376 majority of these samples were in mineral soils, while the limiting factor was the high variability of organic soils
377 combined with its low number of samples. Nevertheless, an improvement was noted in relation to the all error
378 metrics of SVR in the AP2L approach. This was in contrast to when the training set was enlarged without
379 subdividing the data, i.e. AP1L. Therefore, it further confirmed that it is more important for SVR than BRT and
380 RF to model the soil classes separately when its training set is enlarged by datasets with similar characteristics.

381 Furthermore, the improvement of the algorithms in AP2 and AP2L was particularly noticeable in their relative
382 residuals. By comparing these results with those from AP1 and AP1L, it was evident that the greatest improvement
383 was observed in the northern region and the spatial distribution of relative residuals was more homogenous
384 throughout the country for all algorithms, but particularly for RF and SVR (Fig. 3 C and D). This is understandable
385 since by subdividing the data, the algorithms can no longer exploit any information from the map of organic soil
386 for spatial variability of SOC in mineral soils. Thus, they obtain information from other covariates for this soil
387 class (Fig. 4 B). Although land use and total nitrogen were still among the most important variables for the
388 algorithms in mineral soils, the importance of the predictors representing the SCORPAN C and P factors increased
389 in the absence of a soil organic map. This could be expected because the north-east of Germany, for example, has
390 continental climate (Roßkopf et al., 2015) and young moraine landscapes, while the north-west has a more oceanic
391 climate (Roßkopf et al., 2015) with old moraine landscapes.

392 It is unsurprising that all the algorithms still relied on the map of organic soil to explain SOC in organic soil class.
393 However, while SVR and RF still obtained information from other covariates, the value for variable importance
394 of this map alone is 93% in BRT (Fig. 4 C). That makes this algorithm prone to greater errors, as can be seen in
395 its error metrics (Table S2). Similar to mineral soils, the order of covariates was different between the algorithms
396 in organic soils. In other words, in AP1 the three algorithms obtained almost all information from the map of
397 organic soil, land-use and total nitrogen with similar order. In contrast, after subdividing the data, the algorithms
398 differentiated from each other by the order of covariates in their variable importance (Figure 4).



399 A comparison of the error metrics of each soil class in AP2 with its counterpart in AP2L revealed that the additional
400 1177 samples had a minor influence on the performance (from zero to a maximum of 2%) of the algorithms in
401 mineral soils (Table S2). These results indicated that the German Agricultural Soil Inventory offers a good
402 representation of the spatial variability of SOC in mineral soil under agricultural use throughout the country and
403 including more sample points do not provide additional information about SOC variability in this soil class.

404 However, 46 additional organic soil samples from the LUCAS dataset improved MAPE and MAE by 12% and
405 6% for SVR, by 10%, and 4% for RF, and by 7% and 2% for BRT, respectively, but the RMSE of the three
406 algorithms was improved by less than 2%. Thus, additional organic samples mainly influenced the average
407 magnitude of the error. This could be explained by organic soils having a wide range of SOC and the number of
408 samples was limited. Thus, the addition of LUCAS data to the training set offered the algorithms more information
409 about spatial variability of SOC in this soil class. Despite this limitation, SVR had the best overall performance
410 among the algorithms in AP2 and AP2L. It should be noted that training samples must span the complexity of the
411 parameter space in order for the model to be able to effectively match the training data and to generalize unseen
412 data. Small sample size can therefore negatively influence the predictive power of the algorithms. This complexity
413 can be addressed by structural risk minimisation (SRM) (Al-Anazi and Gates, 2012). Implementation of SRM
414 makes SVR capable of performing well in such datasets. Other studies have compared the performance of
415 algorithms on different sample sizes for predicting soil properties and shown that SVR is one of the best choices,
416 if not the best, when the number of samples is a limiting factor (Al-Anazi and Gates, 2012; Khaledian and Miller,
417 2020). In contrast, in a study by Zhou et al. (2021), 150 samples with different sets of covariates at different
418 resolutions were used to compare RF, BRT and SVR to predict SOC content in Switzerland. Their results showed
419 that the best-performing algorithm varied depending on the resolution and covariates. However, the best
420 performance throughout all scenarios was obtained by BRT. The discrepancy between their results and the results
421 of the present study may be due to the parameter-tuning method of the algorithms, as they only used grid search,
422 or other factors, including the spatial distribution of samples or the chosen set of covariates.

423 **4 Conclusions**

424 The three most commonly used algorithms in DSM were implemented to predict the SOC content of German
425 agricultural soils under different approaches. Suitable tuning strategies for each algorithm ensured optimum
426 parameter tuning and made their performance truly comparable. Machine learning algorithms was shown to be
427 powerful in modelling SOC on a national scale. However, the study showed that separate modelling of mineral
428 and organic soils was a better approach for modelling SOC compared to using one model. Thus, this approach has
429 priority to the choice of algorithm and number of training samples. We recommend this approach to be further
430 tested in countries and regions that cover both of these soil classes. Nonetheless, SVR had better performance than
431 RF and BRT except when the number of samples in training was increased. This was disadvantageous for SVR
432 and advantageous for BRT unless mineral and organic soils were modelled separately. Therefore, this approach
433 should be done with consideration of the algorithm and the characteristics of the data. Furthermore, better
434 performance of SVR over RF and BRT was particularly highlighted when predicting SOC in organic soils. Thus,
435 this algorithm should therefore be taken into greater account in DSM when the number of samples is limited.



436 **Data availability**

437 The soil data used in this study are publicly available via: <https://doi.org/10.3220/DATA20200203151139> and
438 <https://esdac.jrc.ec.europa.eu/content/lucas-2009-topsoil-data>

439 **Author contribution**

440 AS and AD conceptualised and developed the methodology of the presented work, with input from ML.AS
441 gathered the predictors with contribution from AD. AS executed programming, testing of existing code
442 components, formal analysis and visualization. AG contributed to programming. The preparation of the paper was
443 done by all authors.

444 **Competing interests**

445 The authors declare that they have no conflict of interest except author AD is a member of the editorial board of
446 the journal.

447 **Acknowledgement**

448 This work is part of the SoilSpace3D-DE project. The LUCAS topsoil dataset used in this work was made available
449 by the European Commission through the European Soil Data Centre managed by the Joint Research Centre (JRC),
450 <http://esdac.jrc.ec.europa.eu/>.

451

452



453 **Reference**

- 454 Al-Anazi, A. F. and Gates, I. D.: Support vector regression to predict porosity and permeability: Effect of sample
455 size, *Comput. Geosci.*, 39, 64–76, <https://doi.org/10.1016/j.cageo.2011.06.011>, 2012.
- 456 Awad, M. and Khanna, R.: Support Vector Regression, in: *Efficient Learning Machines*, Apress, Berkeley, CA,
457 67–80, https://doi.org/10.1007/978-1-4302-5990-9_4, 2015.
- 458 Ballabio, C., Panagos, P., and Monatanarella, L.: Mapping topsoil physical properties at European scale using
459 the LUCAS database, *Geoderma*, 261, 110–123, <https://doi.org/10.1016/j.geoderma.2015.07.006>, 2016.
- 460 Ballabio, C., Lugato, E., Fernández-Ugalde, O., Orgiazzi, A., Jones, A., Borrelli, P., Montanarella, L., and
461 Panagos, P.: Mapping LUCAS topsoil chemical properties at European scale using Gaussian process regression,
462 *Geoderma*, 355, 113912, <https://doi.org/10.1016/j.geoderma.2019.113912>, 2019.
- 463 Battineni, G., Chintalapudi, N., and Amenta, F.: Machine learning in medicine: Performance calculation of
464 dementia prediction by support vector machines (SVM), *Informatics Med. Unlocked*, 16, 100200,
465 <https://doi.org/10.1016/j.imu.2019.100200>, 2019.
- 466 Behrens, T. and Scholten, T.: Digital soil mapping in Germany - A review, *J. Plant Nutr. Soil Sci.*, 169, 434–
467 443, <https://doi.org/10.1002/jpln.200521962>, 2006.
- 468 Behrens, T., Förster, H., Scholten, T., Steinrücken, U., Spies, E. D., and Goldschmitt, M.: Digital soil mapping
469 using artificial neural networks, *J. Plant Nutr. Soil Sci.*, 168, 21–33, <https://doi.org/10.1002/jpln.200421414>,
470 2005.
- 471 Belon, E., Boisson, M., Deportes, I. Z., Eglin, T. K., Feix, I., Bispo, A. O., Galsomies, L., Leblond, S., and
472 Guellier, C. R.: An inventory of trace elements inputs to French agricultural soils, *Sci. Total Environ.*, 439, 87–
473 95, <https://doi.org/10.1016/j.scitotenv.2012.09.011>, 2012.
- 474 Bhadra, T., Bandyopadhyay, S., and Maulik, U.: Differential Evolution Based Optimization of SVM Parameters
475 for Meta Classifier Design, *Procedia Technol.*, 4, 50–57, <https://doi.org/10.1016/j.protecy.2012.05.006>, 2012.
- 476 BKG: *Digitales Basis-Landschaftsmodell (Basis-DLM)*, Leipzig, 2019.
- 477 Borrelli, P., Van Oost, K., Meusburger, K., Alewell, C., Lugato, E., and Panagos, P.: A step towards a holistic
478 assessment of soil degradation in Europe: Coupling on-site erosion with sediment transfer and carbon fluxes,
479 *Environ. Res.*, 161, 291–298, <https://doi.org/10.1016/j.envres.2017.11.009>, 2018.
- 480 Breiman, L.: *Randon Forests*, 1–35, 1999.
- 481 de Brogniez, D., Ballabio, C., Stevens, A., Jones, R. J. A., Montanarella, L., and van Wesemael, B.: A map of
482 the topsoil organic carbon content of Europe generated by a generalized additive model, *Eur. J. Soil Sci.*, 66,
483 121–134, <https://doi.org/10.1111/ejss.12193>, 2015.
- 484 Burke, I. C., Yonker, C. M., Parton, W. J., Cole, C. V., Flach, K., and Schimel, D. S.: Texture, Climate, and
485 Cultivation Effects on Soil Organic Matter Content in U.S. Grassland Soils, *Soil Sci. Soc. Am. J.*, 53, 800–805,
486 <https://doi.org/10.2136/sssaj1989.03615995005300030029x>, 1989.
- 487 Camera, C., Zomeni, Z., Noller, J. S., Zissimos, A. M., Christoforou, I. C., and Bruggeman, A.: A high
488 resolution map of soil types and physical properties for Cyprus: A digital soil mapping optimization, *Geoderma*,
489 285, 35–49, <https://doi.org/10.1016/j.geoderma.2016.09.019>, 2017.
- 490 Carter, B. J. and Ciolkosz, E. J.: Slope gradient and aspect effects on soils developed from sandstone in
491 Pennsylvania, *Geoderma*, 49, 199–213, [https://doi.org/10.1016/0016-7061\(91\)90076-6](https://doi.org/10.1016/0016-7061(91)90076-6), 1991.
- 492 Castaldi, F., Hueni, A., Chabrilat, S., Ward, K., Buttafuoco, G., Bomans, B., Vreys, K., Brell, M., and van
493 Wesemael, B.: Evaluating the capability of the Sentinel 2 data for soil organic carbon prediction in croplands,
494 *ISPRS J. Photogramm. Remote Sens.*, 147, 267–282, <https://doi.org/10.1016/j.isprsjprs.2018.11.026>, 2019.
- 495 Chapman, S. J., Bell, J. S., Campbell, C. D., Hudson, G., Lilly, A., Nolan, A. J., Robertson, A. H. J., Potts, J. M.,
496 and Towers, W.: Comparison of soil carbon stocks in Scottish soils between 1978 and 2009, *Eur. J. Soil Sci.*, 64,
497 455–465, <https://doi.org/10.1111/ejss.12041>, 2013.
- 498 Cherkassky, V. and Ma, Y.: *FL Methods New Genetic Technology.pdf*, 17, 113–126, 2004.
- 499 Conrad, O., Bechtel, B., Bock, M., Dietrich, H., Fischer, E., Gerlitz, L., Wehberg, J., Wichmann, V., and
500 Böhner, J.: System for Automated Geoscientific Analyses (SAGA) v. 2.1.4, *Geosci. Model Dev.*, 8, 1991–2007,



- 501 <https://doi.org/10.5194/gmd-8-1991-2015>, 2015.
- 502 Deng, S., Wang, C., Wang, M., and Sun, Z.: A gradient boosting decision tree approach for insider trading
503 identification: An empirical model evaluation of China stock market, *Appl. Soft Comput. J.*, 83, 105652,
504 <https://doi.org/10.1016/j.asoc.2019.105652>, 2019.
- 505 DWD Climate Data Center (CDC): Multi-annual grids of annual sunshine duration over Germany 1981-
506 2010, version v1.0, 2017.
- 507 DWD Climate Data Center (CDC): Multi-annual grids of monthly averaged daily minimum air temperature (2m)
508 over Germany, version v1.0, 2018a.
- 509 DWD Climate Data Center (CDC): Multi-annual grids of number of summer days over Germany, version v1.0,
510 2018b.
- 511 DWD Climate Data Center (CDC): Multi-annual grids of precipitation height over Germany 1981-2010, version
512 v1.0, 2018c.
- 513 Environmental Systems Research Institute (ESRI): ArcGIS 10.2 for Desktop, 2013.
- 514 European Union Copernicus Land Monitoring Service, and E. E. A. (EEA): European Digital Elevation Model
515 (EU-DEM), Version 1.1, 2016.
- 516 Eurostat grid generation tool for ArcGIS: [https://www.efgs.info/information-base/best-practices/tools/eurostat-
517 grid-generation-tool-arcgis/](https://www.efgs.info/information-base/best-practices/tools/eurostat-grid-generation-tool-arcgis/).
- 518 Federal Institute for Geosciences and Natural Resources (BGR): Geomorphographic Map of Germany
519 1:1,000,000, 2006.
- 520 Federal Institute for Geosciences and Natural Resources (BGR) and German State Geological Surveys (SGD):
521 Hydrogeological Map of Germany 1:250,000 (HÜK250), 2019.
- 522 Forkuor, G., Hounkpatin, O. K. L., Welp, G., and Thiel, M.: High resolution mapping of soil properties using
523 Remote Sensing variables in south-western Burkina Faso: A comparison of machine learning and multiple linear
524 regression models, *PLoS One*, 12, 1–21, <https://doi.org/10.1371/journal.pone.0170478>, 2017.
- 525 Friedman, J., Tibshirani, R., and Hastie, T.: Additive logistic regression: a statistical view of boosting (With
526 discussion and a rejoinder by the authors), *Ann. Stat.*, 28, 337–407, <https://doi.org/10.1214/aos/1016120463>,
527 2000.
- 528 Friedman, J., Hastie, T., and Tibshirani, R.: *The Elements of Statistical Learning*, second., Springer New York,
529 New York, NY, 158-61 pp., <https://doi.org/10.1007/b94608>, 2009.
- 530 Friedman, J. H.: Stochastic gradient boosting, *Comput. Stat. Data Anal.*, 38, 367–378,
531 [https://doi.org/10.1016/S0167-9473\(01\)00065-2](https://doi.org/10.1016/S0167-9473(01)00065-2), 2002.
- 532 Gebauer, A., Ellinger, M., M. Brito Gomez, V., and Ließ, M.: Development of pedotransfer functions for water
533 retention in tropical mountain soil landscapes: Spotlight on parameter tuning in machine learning, 6, 215–229,
534 <https://doi.org/10.5194/soil-6-215-2020>, 2020.
- 535 Greenwell, B., Boehmke, B., and Cunningham, J.: Package “gbm” - Generalized Boosted Regression Models,
536 CRAN Repos., 39, 2019.
- 537 Guevara, M., Olmedo, G. F., Stell, E., Yigini, Y., Aguilar Duarte, Y., Arellano Hernández, C., Arévalo, G.,
538 Arroyo-Cruz, C. E., Bolivar, A., Bunning, S., Bustamante Cañas, N., Cruz-Gaistardo, C. O., Davila, F., Dell
539 Acqua, M., Encina, A., Figueredo Tacona, H., Fontes, F., Hernández Herrera, J. A., Ibelle Navarro, A. R.,
540 Loayza, V., Manueles, A., Mendoza Jara, F., Olivera, C., Osorio Herмосilla, R., Pereira, G., Prieto, P., Alexis
541 Ramos, I., Rey Brina, J. C., Rivera, R., Rodríguez-Rodríguez, J., Roopnarine, R., Rosales Ibarra, A., Rosales
542 Riveiro, K. A., Schulz, G. A., Spence, A., Vasques, G., Vargas, R., and Vargas, R.: No Silver Bullet for Digital
543 Soil Mapping: Country-specific Soil Organic Carbon Estimates across Latin America, *SOIL Discuss.*, 1–20,
544 <https://doi.org/10.5194/soil-2017-40>, 2018.
- 545 Guio Blanco, C. M., Brito Gomez, V. M., Crespo, P., and Ließ, M.: Spatial prediction of soil water retention in a
546 Páramo landscape: Methodological insight into machine learning using random forest, *Geoderma*, 316, 100–114,
547 <https://doi.org/10.1016/j.geoderma.2017.12.002>, 2018.
- 548 Hawkins, D. M., Basak, S. C., and Mills, D.: Assessing model fit by cross-validation, *J. Chem. Inf. Comput. Sci.*,



- 549 43, 579–586, <https://doi.org/10.1021/ci025626i>, 2003.
- 550 Hornik, K., Weingessel, A., Leisch, F., and Davidmeyer-projectorg, M. D. M.: Package ‘e1071’, 2021.
- 551 Hoyle, F. C., Baldock, J. A., and Murphy, D. V.: Soil Organic Carbon – Role in Rainfed Farming Systems,
552 Rainfed Farming Syst., <https://doi.org/10.1007/978-1-4020-9132-2>, 2011.
- 553 Khaleedian, Y. and Miller, B. A.: Selecting appropriate machine learning methods for digital soil mapping, Appl.
554 Math. Model., 81, 401–418, <https://doi.org/10.1016/j.apm.2019.12.016>, 2020.
- 555 Kuhn, M. and Johnson, K.: Applied predictive modeling, 1st ed., Springer-Verlag, New York, 1–600 pp.,
556 <https://doi.org/10.1007/978-1-4614-6849-3>, 2013.
- 557 Lal, R.: Soil carbon sequestration impacts on global climate change and food security, Science (80-.), 304,
558 1623–1627, <https://doi.org/10.1126/science.1097396>, 2004.
- 559 Li, T., Zhang, H., Wang, X., Cheng, S., Fang, H., Liu, G., and Yuan, W.: Soil erosion affects variations of soil
560 organic carbon and soil respiration along a slope in Northeast China, Ecol. Process., 8,
561 <https://doi.org/10.1186/s13717-019-0184-6>, 2019.
- 562 Li, X., Ding, J., Liu, J., Ge, X., and Zhang, J.: Digital mapping of soil organic carbon using sentinel series data:
563 A case study of the ebinur lake watershed in xinjiang, Remote Sens., 13, 1–19,
564 <https://doi.org/10.3390/rs13040769>, 2021.
- 565 Liang, W., Zhang, L., and Wang, M.: The chaos differential evolution optimization algorithm and its application
566 to support vector regression machine, J. Softw., 6, 1297–1304, <https://doi.org/10.4304/jsw.6.7.1297-1304>, 2011.
- 567 Ließ, M., Gebauer, A., and Don, A.: Machine Learning With GA Optimization to Model the Agricultural Soil-
568 Landscape of Germany: An Approach Involving Soil Functional Types With Their Multivariate Parameter
569 Distributions Along the Depth Profile, Front. Environ. Sci., 9, 1–24, <https://doi.org/10.3389/fenvs.2021.692959>,
570 2021.
- 571 Malik, A. A., Puissant, J., Buckeridge, K. M., Goodall, T., Jehmlich, N., Chowdhury, S., Gweon, H. S., Peyton,
572 J. M., Mason, K. E., van Agtmaal, M., Bland, A., Clark, I. M., Whitaker, J., Pywell, R. F., Ostle, N., Gleixner,
573 G., and Griffiths, R. I.: Land use driven change in soil pH affects microbial carbon cycling processes, Nat.
574 Commun., 9, 1–10, <https://doi.org/10.1038/s41467-018-05980-1>, 2018.
- 575 Martin, M. P., Orton, T. G., Lacarce, E., Meersmans, J., Saby, N. P. A., Paroissien, J. B., Jolivet, C., Boulonne,
576 L., and Arrouays, D.: Evaluation of modelling approaches for predicting the spatial distribution of soil organic
577 carbon stocks at the national scale, Geoderma, 223–225, 97–107,
578 <https://doi.org/10.1016/j.geoderma.2014.01.005>, 2014.
- 579 McBratney, A. B., Mendonça Santos, M. L., and Minasny, B.: On digital soil mapping, 3–52 pp.,
580 [https://doi.org/10.1016/S0016-7061\(03\)00223-4](https://doi.org/10.1016/S0016-7061(03)00223-4), 2003.
- 581 Meersmans, J., Martin, M. P., Lacarce, E., De Baets, S., Jolivet, C., Boulonne, L., Lehmann, S., Saby, N. P. A.,
582 Bispo, A., and Arrouays, D.: A high resolution map of French soil organic carbon, Agron. Sustain. Dev., 32,
583 841–851, <https://doi.org/10.1007/s13593-012-0086-9>, 2012a.
- 584 Meersmans, J., Martin, M. P., De Ridder, F., Lacarce, E., Wetterlind, J., De Baets, S., Bas, C. Le, Louis, B. P.,
585 Orton, T. G., Bispo, A., and Arrouays, D.: A novel soil organic C model using climate, soil type and
586 management data at the national scale in France, Agron. Sustain. Dev., 32, 873–888,
587 <https://doi.org/10.1007/s13593-012-0085-x>, 2012b.
- 588 Minasny, B., McBratney, A. B., Malone, B. P., and Wheeler, I.: Digital Mapping of Soil Carbon, Elsevier, 1–47
589 pp., <https://doi.org/10.1016/B978-0-12-405942-9.00001-3>, 2013.
- 590 Mulder, V. L., Lacoste, M., Richer-de-Forges, A. C., Martin, M. P., and Arrouays, D.: National versus global
591 modelling the 3D distribution of soil organic carbon in mainland France, Geoderma, 263, 16–34,
592 <https://doi.org/10.1016/J.GEODERMA.2015.08.035>, 2016.
- 593 Padarian, J., Minasny, B., and McBratney, A. B.: Machine learning and soil sciences: A review aided by
594 machine learning tools, 6, 35–52, <https://doi.org/10.5194/soil-6-35-2020>, 2020.
- 595 Pei, T., Qin, C. Z., Zhu, A. X., Yang, L., Luo, M., Li, B., and Zhou, C.: Mapping soil organic matter using the
596 topographic wetness index: A comparative study based on different flow-direction algorithms and kriging
597 methods, Ecol. Indic., 10, 610–619, <https://doi.org/10.1016/j.ecolind.2009.10.005>, 2010.



- 598 Peterson, B., Ulrich, J., and Boudt, K.: Package ‘DEoptim,’ <https://doi.org/10.18637/jss.v040.i06>, 2021.
- 599 Poeplau, C., Bolinder, M. A., Eriksson, J., Lundblad, M., and Kätterer, T.: Positive trends in organic carbon
600 storage in Swedish agricultural soils due to unexpected socio-economic drivers, 12, 3241–3251,
601 <https://doi.org/10.5194/bg-12-3241-2015>, 2015.
- 602 Poeplau, C., Jacobs, A., Don, A., Vos, C., Schneider, F., Wittnebel, M., Tiemeyer, B., Heidkamp, A., Prietz, R.,
603 and Flessa, H.: Stocks of organic carbon in German agricultural soils—Key results of the first comprehensive
604 inventory, *J. Plant Nutr. Soil Sci.*, 183, 665–681, <https://doi.org/10.1002/jpln.202000113>, 2020.
- 605 Prechtel, A., Von Lützow, M., Schneider, B. U., Bens, O., Bannick, C. G., Kögel-Knabner, I., and Hüttl, R. F.:
606 Organic Carbon in soils of Germany: Status quo and the need for new data to evaluate potentials and trends of
607 soil carbon sequestration, *J. Plant Nutr. Soil Sci.*, 172, 601–614, <https://doi.org/10.1002/jpln.200900034>, 2009.
- 608 Probst, P., Wright, M. N., and Boulesteix, A. L.: Hyperparameters and tuning strategies for random forest, *Wiley*
609 *Interdiscip. Rev. Data Min. Knowl. Discov.*, 9, 1–15, <https://doi.org/10.1002/widm.1301>, 2019.
- 610 Qin, A. K., Huang, V. L., and Suganthan, P. N.: Differential evolution algorithm with strategy adaptation for
611 global numerical optimization, *IEEE Trans. Evol. Comput.*, 13, 398–417,
612 <https://doi.org/10.1109/TEVC.2008.927706>, 2009.
- 613 Ramifehiarivo, N., Brossard, M., Grinand, C., Andriamananjara, A., Razafimbelo, T., Rasolohery, A.,
614 Razafimahatratra, H., Seyler, F., Ranaivoson, N., Rabenarivo, M., Albrecht, A., Razafindrabe, F., and
615 Razakamanarivo, H.: Mapping soil organic carbon on a national scale: Towards an improved and updated map
616 of Madagascar, *Geoderma Reg.*, 9, 29–38, <https://doi.org/10.1016/j.geodrs.2016.12.002>, 2017.
- 617 Rawlins, B. G., Marchant, B. P., Smyth, D., Scheib, C., Lark, R. M., and Jordan, C.: Airborne radiometric survey
618 data and a DTM as covariates for regional scale mapping of soil organic carbon across Northern Ireland, *Eur. J.*
619 *Soil Sci.*, 60, 44–54, <https://doi.org/10.1111/j.1365-2389.2008.01092.x>, 2009.
- 620 Reeves, D. W.: The role of soil organic matter in maintaining soil quality in continuous cropping systems, *Soil*
621 *Tillage Res.*, 43, 131–167, [https://doi.org/10.1016/S0167-1987\(97\)00038-X](https://doi.org/10.1016/S0167-1987(97)00038-X), 1997.
- 622 Richter, A., Adler, G. H., Fahrak, M., and Eckelmann, W.: Erläuterungen zur nutzungsdifferenzierten
623 Bodenübersichtskarte der Bundesrepublik Deutschland im Maßstab 1:1.000.000, 000, 1–54, 2007.
- 624 Ritchie, J. C., McCarty, G. W., Venteris, E. R., and Kaspar, T. C.: Soil and soil organic carbon redistribution on
625 the landscape, 89, 163–171, <https://doi.org/10.1016/j.geomorph.2006.07.021>, 2007.
- 626 Roßberg, D., Michel, V., Graf, R., Neukampf, R.: Definition von Boden-Klima-Räumen für die Bundesrepublik
627 Deutschland, *Nachrichtenblatt des Dtsch. Pflanzenschutzdienstes*, 59, 155–161, 2007.
- 628 Roßkopf, N., Fell, H., and Zeitz, J.: Organic soils in Germany, their distribution and carbon stocks, 133, 157–
629 170, <https://doi.org/10.1016/j.catena.2015.05.004>, 2015.
- 630 Santos, C. E. da S., Sampaio, R. C., Coelho, L. dos S., Bestarsd, G. A., and Llanos, C. H.: Multi-objective
631 adaptive differential evolution for SVM/SVR hyperparameters selection, *Pattern Recognit.*, 110, 107649,
632 <https://doi.org/10.1016/j.patcog.2020.107649>, 2021.
- 633 Schapire, R. E.: The Boosting Approach to Machine Learning: An Overview, 149–171,
634 https://doi.org/10.1007/978-0-387-21579-2_9, 2003.
- 635 Schneider, F., Amelung, W., and Don, A.: Origin of carbon in agricultural soil profiles deduced from depth
636 gradients of C:N ratios, carbon fractions, $\delta^{13}\text{C}$ and $\delta^{15}\text{N}$ values, *Plant Soil*, 460, 123–148,
637 <https://doi.org/10.1007/s11104-020-04769-w>, 2021.
- 638 Smola, A. J. and Schölkopf, B.: A tutorial on support vector regression, *Stat. Comput.*, 14, 199–222,
639 <https://doi.org/10.1023/B:STCO.0000035301.49549.88>, 2004.
- 640 Storn, R. and Price, K.: Differential Evolution - A Simple and Efficient Heuristic for Global Optimization over
641 Continuous Spaces, *J. Glob. Optim.*, 11, 341–359, <https://doi.org/10.1023/A:1008202821328>, 1997.
- 642 Taghizadeh-Toosi, A., Olesen, J. E., Kristensen, K., Elsgaard, L., Østergaard, H. S., Lægdsmand, M., Greve, M.
643 H., and Christensen, B. T.: Changes in carbon stocks of Danish agricultural mineral soils between 1986 and
644 2009, *Eur. J. Soil Sci.*, 65, 730–740, <https://doi.org/10.1111/ejss.12169>, 2014.
- 645 Tóth, G., Jones, A., and Montanarella, L.: LUCAS Topsoil Survey: Methodology, Data, and Results,



- 646 <https://doi.org/10.2788/97922>, 2013.
- 647 Tziachris, P., Aschonitis, V., Chatzistathis, T., Papadopoulou, M., and Doukas, I. J. D.: Comparing machine
648 learning models and hybrid geostatistical methods using environmental and soil covariates for soil pH prediction,
649 *ISPRS Int. J. Geo-Information*, 9, <https://doi.org/10.3390/ijgi9040276>, 2020.
- 650 Varma, S. and Simon, R.: Bias in error estimation when using cross-validation for model selection, *BMC*
651 *Bioinformatics*, 7, 1–8, <https://doi.org/10.1186/1471-2105-7-91>, 2006.
- 652 Wadoux, A. M. J. C., Minasny, B., and McBratney, A. B.: Machine learning for digital soil mapping:
653 Applications, challenges and suggested solutions, *Earth-Science Rev.*, 210,
654 <https://doi.org/10.1016/j.earscirev.2020.103359>, 2020.
- 655 Wang, S., Xu, L., Zhuang, Q., and He, N.: Investigating the spatio-temporal variability of soil organic carbon
656 stocks in different ecosystems of China, *Sci. Total Environ.*, 758,
657 <https://doi.org/10.1016/j.scitotenv.2020.143644>, 2021.
- 658 Wang, X., Zhang, Y., Atkinson, P. M., and Yao, H.: Predicting soil organic carbon content in Spain by
659 combining Landsat TM and ALOS PALSAR images, *Int. J. Appl. Earth Obs. Geoinf.*, 92, 102182,
660 <https://doi.org/10.1016/j.jag.2020.102182>, 2020.
- 661 Ward, K. J., Chabrillat, S., Neumann, C., and Foerster, S.: A remote sensing adapted approach for soil organic
662 carbon prediction based on the spectrally clustered LUCAS soil database, *Geoderma*, 353, 297–307,
663 <https://doi.org/10.1016/j.geoderma.2019.07.010>, 2019.
- 664 Were, K., Bui, D. T., Dick, Ø. B., and Singh, B. R.: A comparative assessment of support vector regression,
665 artificial neural networks, and random forests for predicting and mapping soil organic carbon stocks across an
666 Afromontane landscape, *Ecol. Indic.*, 52, 394–403, <https://doi.org/10.1016/j.ecolind.2014.12.028>, 2015.
- 667 Wiesmeier, M., Spörlein, P., Geuß, U., Hangen, E., Haug, S., Reischl, A., Schilling, B., von Lützwow, M., and
668 Kögel-Knabner, I.: Soil organic carbon stocks in southeast Germany (Bavaria) as affected by land use, soil type
669 and sampling depth, *Glob. Chang. Biol.*, 18, 2233–2245, <https://doi.org/10.1111/j.1365-2486.2012.02699.x>,
670 2012.
- 671 Wright, M. N. and Ziegler, A.: Ranger: A fast implementation of random forests for high dimensional data in
672 C++ and R, *J. Stat. Softw.*, 77, 1–17, <https://doi.org/10.18637/jss.v077.i01>, 2017.
- 673 Yu, D., Hu, F., Zhang, K., Liu, L., and Li, D.: Available water capacity and organic carbon storage profiles in
674 soils developed from dark brown soil to boggy soil in Changbai Mountains, China, *Soil Water Res.*, 16, 11–21,
675 <https://doi.org/10.17221/150/2019-SWR>, 2021.
- 676 Zhang, J., Niu, Q., Li, K., and Irwin, G. W.: Model selection in SVMs using Differential Evolution, *IFAC*,
677 14717–14722 pp., <https://doi.org/10.3182/20110828-6-IT-1002.00584>, 2011.
- 678 Zhong, Z., Chen, Z., Xu, Y., Ren, C., Yang, G., Han, X., Ren, G., and Feng, Y.: Relationship between soil
679 organic carbon stocks and clay content under different climatic conditions in Central China, 9, 1–14,
680 <https://doi.org/10.3390/f9100598>, 2018.
- 681 Zhou, T., Geng, Y., Ji, C., Xu, X., Wang, H., Pan, J., Bumberger, J., Haase, D., and Lausch, A.: Prediction of
682 soil organic carbon and the C:N ratio on a national scale using machine learning and satellite data: A comparison
683 between Sentinel-2, Sentinel-3 and Landsat-8 images, *Sci. Total Environ.*, 755,
684 <https://doi.org/10.1016/j.scitotenv.2020.142661>, 2021.
- 685