

~~Performance of three machine learning algorithms for predicting soil organic carbon in German agricultural soil~~ Spatial prediction of organic carbon in German agricultural topsoil using machine learning algorithms

Ali Sakhaee¹, Anika Gebauer², Mareike Ließ², Axel Don¹

¹Thünen Institute of Climate Smart Agriculture, Braunschweig, Germany

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

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

Abstract

~~As the largest terrestrial carbon pool, S~~soil organic carbon (SOC), ~~as the largest terrestrial carbon pool,~~ has the potential to influence and mitigate climate change ~~and mitigation, and consequently~~ hence the importance of SOC monitoring ~~is important~~ in the frameworks of different various international treaties. ~~There is therefore a need for~~ hHigh resolution SOC maps are therefore required. Machine learning (ML) offers new opportunities to ~~do this~~ develop these due to its capability for data mining of large datasets. The aim of this study, ~~therefore,~~ was to ~~test~~ apply three ~~commonly used~~ algorithms commonly used in digital soil mapping – random forest (RF), boosted regression trees (BRT) and support vector machine for regression (SVR) – on the first German Agricultural Soil Inventory to model agricultural topsoil (0-30 cm) SOC content; ~~and develop a two-model approach to address the~~ high variability of SOC in German agricultural soils. Model performance is often limited by the size and quality of the soil dataset available for calibration and validation. Therefore, the impact of enlarging the training data was tested by including data from the European Land Use/Land Cover Area Frame Survey for agricultural sites in Germany. Nested cross-validation was implemented for model evaluation and parameter tuning. ~~Moreover, G~~grid search and the differential evolution algorithm were also applied to ensure that each algorithm was appropriately tuned ~~and optimised suitably~~. The SOC content of the German Agricultural Soil Inventory was highly variable, 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 considered organic or degraded organic soils. The results showed ed that SVR ~~provided~~ produced the best performance with an RMSE of 32 g kg⁻¹ when the algorithms were trained on the full dataset. However, the average RMSE of all algorithms decreased by 34% when mineral and organic soils were modelled separately, with the best result from SVR with a RMSE of 21 g kg⁻¹. ~~Model performance is often limited by the size and quality of the~~ available soil dataset for calibration and validation. Therefore, the impact of enlarging the training data was tested by including 1223 data points from the European Land Use/Land Cover Area Frame Survey for agricultural sites in Germany. The model performance was enhanced for maximum by up to 1% for mineral soils and by 2% for organic soils. Despite the capability of machine learning algorithms in general, and particularly SVR in particular, ~~in to~~ modelling SOC on a national scale, the study showed that the most important aspect to for improving the model performance was to separate the modelling of mineral and organic soils.

36 1 Introduction

37 Soil organic carbon (SOC) is the largest terrestrial carbon pool (Wang et al., 2020) and plays an essential role in
38 agriculture. Since SOC influences various physical, chemical and biological properties of soil (Reeves, 1997),
39 numerous studies recognise it as a crucial indicator of soil quality (Castaldi et al., 2019; Meersmans et al., 2012a;
40 Reeves, 1997) ~~(Castaldi et al., 2019; Meersmans et al., 2012; Reeves, 1997) and therefore. Thus,~~ its decline is
41 identified as a threat that leads to soil degradation (Castaldi et al., 2019; Poeplau et al., 2020). Moreover, when
42 considering carbon sequestration, the SOC pool provides the option for climate change mitigation (Meersmans et
43 al., 2012a; Ward et al., 2019) ~~(Meersmans et al., 2012; Ward et al., 2019). Consequently,~~ SOC monitoring is
44 therefore important in the frameworks of various international treaties such as the European Union Soil Thematic
45 Strategy and the United Nations Framework Convention on Climate Change (Meersmans et al., 2012b; Poeplau et
46 al., 2020) ~~(Meersmans et al., 2012; Poeplau et al., 2020) and T~~ there is therefore growing interest in understanding
47 the spatial distribution of SOC at different scales in response to an increasing demand for a better assessment of
48 SOC (Minasny et al., 2013). This is particularly important for agricultural land due to its potential for carbon
49 sequestration (Lal, 2004).

50 In digital soil mapping (DSM), a soil attribute is formulated as described by an empirical quantitative function of
51 seven factors: soil properties, climate, organisms, topography, parent material, time, and spatial position
52 (McBratney et al., 2003). ~~Therefore, T~~ This function, known as the SCORPAN model, can be applied to spatially
53 predict the soil attribute-property of interest (Minasny et al., 2013). Within this framework, machine learning
54 algorithms aim to automatically extract ~~the~~ information from the data for predictive purposes (Behrens et al.,
55 2005). This is of particularly intriguing interest in view of the recent expansion of databases at a different scale in
56 soil sciences soil databases and the the complexity of the covariates vast amount of data to approximate the soil
57 forming factors in recent years (McBratney et al., 2003; Wadoux et al., 2020), thus making DSM cost-effective,
58 time-efficient and applicable over large areas with good results (Behrens and Scholten, 2006; Camera et al., 2017).

59 Despite the advantages of DSM, it is crucial to consider-note that its application requires soil databases of an
60 adequate sample size for training and testing. Furthermore, consistent and quality-checked datasets are a
61 prerequisite for DSM. Several soil inventories and monitoring networks for SOC have been formed-established on
62 a national scale in countries such as Sweden (Poeplau et al., 2015), France (Belon et al., 2012; Arrouays et al.,
63 2002) ~~(Belon et al., 2012; Meersmans et al., 2012)~~ Denmark (Taghizadeh-Toosi et al., 2014) and Scotland
64 (Chapman et al., 2013). ~~Nonetheless~~ However, in Germany the most critical shortcomings of soil inventories ~~in~~
65 Germany concern are the lack of a large-scale, high-quality SOC inventory-monitoring (Wiesmeier et al., 2012)
66 with a periodic and standardised sampling focused on agricultural soils (Prechtel et al., 2009). These issues have
67 now been solved-addressed in the first German Agricultural Soil Inventory (Poeplau et al., 2020). This inventory
68 was conducted-carried out on a national scale with a sampling depth down to 100 cm considering a sampling depth
69 of 1 m at 3104 sampling sites covering agricultural land. Furthermore, on a European scale, the Land Use/Land
70 Cover Area Frame Survey (LUCAS) undertaken in 2009 is the first harmonised topsoil survey with physico-
71 chemical analyses of georeferenced topsoil samples in-from 23 European states (Tóth et al., 2013). Therefore, by
72 taking advantage of DSM and of both the German Agricultural Soil Inventory and ~~the~~ LUCAS survey, it is possible
73 to regionalise from single-point measurements to obtain complete high-resolution cover soil data nationwide and
74 thus provide a baseline for both SOC monitoring as well as for and environmental and climatic modelling for
75 Germany.

76 Boosted regression trees (BRT), random forest (RF) and support vector machine for regression (SVR) are among
77 the most widely used algorithms in DSM (Padarian et al., 2020). For example, Martin et al. (2014) predicted topsoil
78 SOC on a national scale for France using the BRT algorithm ~~and comparing~~ its results when the same algorithm
79 was coupled with a geostatistical approach. They concluded that due to the large distances between sampling sites,
80 spatial autocorrelation is unlikely since spatial autocorrelation is not feasible in most national inventories, and the
81 BRT algorithm alone is sufficient for this purpose. This algorithm has was also been used on a national scale in
82 China for data from the 1980s and 2010s in order to predict topsoil SOC and its spatial-temporal change, as well
83 as the main drivers of its variability (Wang et al., 2021). ~~Moreover,~~ RF has also become more popular in DSM due
84 to its relative simplicity and performance. For example, this algorithm was implemented to map topsoil SOC on a
85 national scale in Madagascar and ~~obtain identify~~ its main drivers (Ramifehiarivo et al., 2017). Ramifehiarivo et al.
86 (2017) ~~concluded that~~ the uncertainty of the map generated by RF model training ~~uncertainty of the algorithm~~
87 was lower when compared with the maps that were formerly generated for the country. Moreover, this algorithm
88 was compared with the Cubist ~~model algorithm~~ for mapping SOC at different resolutions on a regional scale in
89 China and ~~could was found to~~ outperformed it (Li et al., 2021). Fewer studies have used SVR than RF to predict
90 SOC ~~than RF~~. Studies have mainly implemented SVR on a regional scale with a limited number of samples
91 (Forkuor et al., 2017; Were et al., 2015) or on a national scale (Switzerland) with very few samples (150 samples
92 from the European LUCAS survey) (Zhou et al., 2021). However, in a study comparing different algorithms,
93 including SVR and RF, on a continental scale and within each country in Latin America, the results indicated that
94 the best-performing algorithm varied in different countries from country to country (Guevara et al., 2018). ~~Theis~~
95 difference mainly depended on sample data density, quality, dispersity and representativeness and also country
96 size, which affects the heterogeneity of land use and environmental conditions.

97 Another important consideration when applying machine learning is the impact of the parameter-tuning strategy
98 in algorithm performance. This is particularly crucial when the objective of the study is ~~the comparisons of to~~
99 compare different machine learning algorithms. Although some algorithms are less sensitive to tuning, this step is
100 more important for others, particularly those with a higher number of parameters (Tziachris et al., 2020; Wadoux
101 et al., 2020). Furthermore, as algorithms differ by ~~the~~ type of ~~their~~ parameters, continuous or discrete, the chosen
102 strategy should be aligned in accordance with this difference. (Ließ et al., 2021). ~~This is particularly more important~~
103 ~~for algorithms with continuous parameters.~~ For example, ~~it has been shown that that~~ the performance of SVR and
104 BRT ~~is has been shown to be~~ better and more stable when optimised by a differential evolution (DE) algorithm
105 than tuned by grid search (Zhang et al., 2011; Gebauer et al., 2020). Despite this importance, in a review of studies
106 that have applied DSM, Wadoux et al. (2020) state that almost half of them implemented parameter tuning, with
107 grid search the most common strategy applied for this purpose. This finding indicates that the role of parameter
108 tuning and optimisation is unfortunately undermined in DSM. This is particularly evident when the application of
109 machine learning in this field is compared with other fields, where various studies have shown the impact of
110 parameter-tuning strategies on the performance of algorithms such as SVR and BRT (Liang et al., 2011; Santos et
111 al., 2021; Bhadra et al., 2012; Deng et al., 2019).

112 The aims of the present study ~~were was~~ therefore: i) to address the above-mentioned parameter-tuning issue and
113 consequently provide a true comparison of the performance of BRT, RF, and SVR in modelling the SOC contents
114 of German agricultural topsoils (0-30 cm), ii) to assess the impact of training data size by extending the data of the
115 German Agricultural Soil Inventory with LUCAS data for model calibration, and iii) to develop a two-model

116 approach to address the high variability of SOC in German agricultural soils and compare it with a single-model
117 approach.

118 2 Materials and methods

119 2.1 Soil data

120 The models were built using SOC content data from two soil inventories. The first dataset was from the German
121 Agricultural Soil Inventory, which ~~consists-comprise~~ of 3104 sites ~~with a fixed~~ collected along a grid of 8x8 km
122 throughout Germany (Poeplau et al., 2020). The sites were sampled and analysed for different soil properties,
123 including SOC content measured via dry combustion, for the upper 30 cm of the soil between 2012 and 2018. The
124 second dataset was the European LUCAS survey that provides SOC content, similarly also measured via dry
125 combustion, ~~for all EU countries,~~ with the sampling depth limited to 0-20 cm (Tóth et al., 2013). For Germany,
126 data collected on agricultural soils cover 1223 sites. Therefore, in order to harmonise the depths of both datasets,
127 they ~~se~~ were subdivided into two classes: mineral and organic soils ~~classes~~ according to a SOC threshold value of
128 87.0 g kg⁻¹. Accordingly, ~~considering~~ all soils above this threshold were considered as organic soils comprising
129 peat soils and disturbed and degraded peat soils (Poeplau et al., 2020). Linear regression functions were derived
130 for both mineral, Eq. 1, and organic, Eq. 2, soil classes on behalf of the data of the German Agricultural Soil
131 Inventory to relate the SOC content of 0-30 cm to that of 0-20 cm. ~~Linear correlation functions between 0-30 cm~~
132 ~~and 0-20 cm were derived for each soil class of the German Agricultural Soil Inventory separately.~~ These functions
133 were then applied to the corresponding soil class from the LUCAS data in order to estimate 0-30 cm topsoil SOC.
134 ~~With a slope of 0.881 for mineral soils and 1.02 for organic soils, they changed the mean of LUCAS data by less~~
135 ~~than 6%. The depth extrapolated values of mineral and organic soils were then combined to form the complete~~
136 ~~dataset.~~ The 0-30 cm LUCAS data generated and the original 0-20 cm LUCAS data were then used by each
137 algorithm to check the effect of depth extrapolation.

$$138 \quad y = 1.01 + 0.881x \quad (1)$$

$$139 \quad y = 1.6 + 1.02x \quad (2)$$

140 2.2 Covariates

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

148 Following the SCORPAN framework, 24 covariates including x and y coordinates for spatial positions were
149 compiled. In order to ~~capture-represent the~~ climate factors (C factor), precipitation (DWD, 2018c), sunshine
150 duration (DWD, 2017), summer days (DWD, 2018b), and minimum temperature (DWD, 2018a) were used-applied
151 according to the study of Schneider et al. (2021). Using principal component analysis, these four covariates were
152 indicated-identified to be the most important ~~among-out of~~ 34 available climate factors for SOC in the German
153 Agricultural Soil Inventory dataset. Moreover, type of agricultural land use is one of the main drivers of SOC

154 variability ~~at-on~~ a national scale (Poeplau et al., 2020). ~~Thus, therefore~~ the land-use map from the official
155 topographic-~~cartographic~~ information system (BKG, 2019) with its corresponding classes according to the German
156 Agricultural Soil Inventory was rasterised and included. This is a categorical covariate, representing the organism
157 factor of SCORPAN (O factor), ~~that-which~~ distinguishes croplands from grasslands and captures their spatial
158 distribution throughout Germany.

159 The European Digital Elevation Model (EUEM) (European Union Copernicus Land Monitoring Service, 2016)
160 ~~with original resolution of 25 m was resampled to 100 m. and s~~ Six covariates derived from the ~~is-layer~~ ~~resampled~~
161 ~~layer~~ were also added to integrate the ~~topography and~~ relief parameters (R factor). Slope, plan curvature and profile
162 curvature, generated ~~on-with~~ SAGA (Conrad et al., 2015), were included to capture the slope's gradient, convexity-
163 concavity and convergence-divergence. These factors influence the soil distribution throughout the landscape, e.g.
164 affecting flow over the surface, thus impacting SOC and its dynamic (Ritchie et al., 2007). Moreover, ~~north-south~~
165 ~~and east-west aspects were~~ slope exposition (aspect) was ~~obtained-calculated from the from~~ EUEM as ~~these-it~~
166 influences soil development and subsequently affects SOC (Carter and Ciolkosz, 1991). ~~The circular variable was~~
167 ~~then decomposed into northness and eastness.~~ The Topographic Wetness Index (TWI), generated on SAGA
168 ~~(Conrad et al., 2015)~~, was also added since it captures the soil moisture distribution of the landscape and ~~some~~
169 ~~studies have shown its direct correlation has a direct correlation~~ with SOC (Pei et al., 2010). A geomorphographic
170 map of Germany ~~(Federal Institute for Geosciences and Natural Resources (BGR), 2007)~~ ~~-containing-featuring~~ 25
171 geomorphic categories was also used to distinguish between four different landscape areas of the country: North
172 German lowlands, highlands, Alpine foothills and the Alps.

173 Continuing with the framework, a large-scale soil landscape unit map ("~~Soil Scapes in~~
174 ~~Germany~~~~Bodengrosslandschaft~~") (~~Federal Institute for Geosciences and Natural Resources (BGR), 2008~~)
175 comprising 38 classes was used. This covariate divides Germany by various geo-factors that can be compiled into
176 a map with 12 soil regions ~~representing mainly the parent materials~~. Similarly, ~~large-scale~~ the soil-climate region
177 map ("~~Bodenklima~~") (Roßberg et al., 2007) with 50 classes was added. Moreover, ~~the Hydrogeological unit~~
178 ~~according Hydrogeological map of Germany~~ ~~Germany's hydrogeological unit map~~ (BGR and SGD, 2019). ~~The~~
179 ~~hydrogeological map~~ provides information about ~~-hydrogeologically relevant attributes including consolidation,~~
180 ~~type of porosity, permeability, type of rock and geochemical classification, lithology and its hydrological~~
181 ~~characteristics.~~ These categorical maps were rasterised and applied to the model as the P factor of SCORPAN.
182 Moreover, the soil factor of the framework (S factor) was captured by eight covariates that represent different
183 aspects of its properties: the map of organic soils (Roßkopf et al., 2015) that distinguishes mineral soils from
184 organic ones and explains their spatial distribution throughout the country, as well as the maps of nitrogen
185 (Ballabio et al., 2019) and clay content (Ballabio et al., 2016) since they directly correlate with SOC. As nitrogen
186 is a crucial component of soil organic matter, regions with higher total nitrogen have higher SOC (Ballabio et al.,
187 2019). Also for clay content, different studies have shown that coarser soil textures tend to have a lower
188 accumulation of SOC (Zhong et al., 2018; Hoyle et al., 2011). ~~The M~~ map of pH ~~from~~ (Ballabio et al., (2019) ~~was~~
189 ~~included~~ since soil pH directly impacts microbial activities that influence the turnover of soil organic matter, and
190 consequently negatively correlates with SOC (Malik et al., 2018). Furthermore, ~~the~~ map of available water capacity
191 (Ballabio et al., 2016) ~~was used~~ as this soil ~~properties-property~~ is another interactive factor with SOC through plant
192 productivity and soil texture (Burke et al., 1989; Yu et al., 2021). Soil erosion is also a key factor in the SOC cycle
193 (Li et al., 2019), which was added through the ~~soil-map of Europe's net soil erosion and deposition rates~~ ~~erosion~~

194 ~~map of Europe~~ (Borrelli et al., 2018). Based on the WaTEM/SEDEM model, this map illustrates the potential
195 spatial displacement and transport of soil sediments due to water erosion (Borrelli et al., 2018). Figure S1 provides
196 a more detailed view for better visualisation of the covariates that were used in this study.

197 2.3 Boosted Regression Trees

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

209 2.4 Random Forest

210 Similar to BRT, RF is another tree-based algorithm. RF uses bootstrap sampling of the dataset for growing a
211 decision tree. Subsequently, by aggregating the results of a large number of decision trees, the bias and variance
212 of the final model can be reduced (Breiman, 1999). The method of bootstrapping in conjunction with aggregating,
213 known as bagging, increases the robustness and stability of RF. However, the trees from different bootstraps may
214 form a similar structure if all covariates participate in a split of each node. Thus, the variance cannot be reduced
215 optimally through the bagging process (Kuhn and Johnson, 2013). In order to avoid this tree correlation, a random
216 subset of covariates, i.e. predictors, is selected at each split. The parameter m_{try} defines the number of predictors
217 included in this subset and should be tuned (Kuhn and Johnson, 2013). The RF algorithm was implemented by
218 setting the number of trees to 1000 and using the “Ranger” package (Wright and Ziegler, 2017) in R.

219 2.5 Support Vector Regression

220 SVR is a form of support vector machine adopted for regression. From all possible solutions, i.e. estimation
221 function, for the problem, SVR tries to obtain an estimation function ~~with- that has at most the maximum- ϵ error~~
222 deviation from the response values of the training data while minimising model complexity (Smola and Schölkopf,
223 2004). Thus, a symmetrical tolerance threshold, ϵ -insensitivity zone, is created around the estimation function
224 ~~within which the vectors are not penalised~~ (Awad and Khanna, 2015). ~~However,~~ the data vectors of the samples
225 that lie on the boundary of the ϵ -insensitivity zone are called support vectors. The vector lying within the
226 insensitivity zone are not penalized. ~~Therefore,~~ ϵ is an optimisable parameter that controls the width of ϵ -
227 insensitivity, alters the model complexity and inversely impacts the number of support vectors ~~inversely~~
228 (Cherkassky and Ma, 2004). Moreover, the trade-off between model complexity and tolerance of ϵ deviation is
229 controlled by a parameter named C (Smola and Schölkopf, 2004; Cherkassky and Ma, 2004). Optimising the C
230 parameter has a crucial impact on SVR performance since a high C can lead to overfitting, while a low C can cause
231 under fitting (Kuhn and Johnson, 2013). The use of kernel functions makes SVR a powerful tool for nonlinear
232 problems. By implementing these functions, SVR can map the data space from its original dimension to a higher

233 dimensional space where a nonlinear problem can be solved linearly. In this study, the Radial Basis Function
 234 (RBF) kernel was used with gamma as its tuneable parameter. This parameter affects the generalisation
 235 performance of SVR by inversely controlling the influence of support vectors ~~inversely~~ (Battinini et al., 2019).
 236 SVR was implemented from the package e1071 in R (Hornik et al., 2021).

237 2.6 Performance evaluation

238 When training a predictive model, it is important to evaluate its generalisation performance on unseen data of the
 239 same type (Hawkins et al., 2003). ~~However, as the number of available samples is usually a limiting factor, the~~
 240 ~~evaluation process is often done by k-fold cross validation (CV). Therefore, the dataset is divided into k folds and~~
 241 ~~k – 1 folds are used for training the model and one fold for testing. This process is repeated k times so each fold~~
 242 ~~participate in train and test. However, as the number of available samples is usually a limiting factor, the evaluation~~
 243 ~~process is often done by randomly splitting the available dataset into training and testing sets multiple times, i.e.~~
 244 ~~cross-validation (CV). Although this process is effective, it is not entirely immune from biased estimation of error.~~
 245 ~~However, to ensure that the estimated error in model evaluation is as unbiased as possible, However, to ensure the~~
 246 ~~robustness of the model, every each~~ model training step should be performed within the CV. This includes finding
 247 the best parameter sets for the chosen algorithm (Varma and Simon, 2006). Thus, the algorithms in this study were
 248 applied on a stratified nested CV.

249 First, to ensure that the SOC distribution was represented in the CV scheme, Germany was divided into 50 strata
 250 using a 100x100 km INSPIRE grid ~~into 50 strata~~. Random samples from each stratum were then taken and
 251 compiled into a fold. This procedure was continued to create five folds and was repeated five times, forming the
 252 outer loop of CV used for model evaluation. ~~Large A long~~ distance between neighboring samples, 8120 m on
 253 average, prevents train and test data from being spatially autocorrelated. Since the aim was to tune the algorithms'
 254 ~~parameters of the algorithms~~, the training set of the outer loop of CV was nested, creating five folds as the inner
 255 loop on which the parameter tuning was performed. To evaluate the performance of algorithms, root-mean-squared
 256 error (RMSE), Eq. 34, mean absolute error (MAE), Eq. 42, and mean absolute percentage error (MAPE), Eq. 53,
 257 were used. ~~Furthermore, AIC, Eq. 6, BIC, Eq. 7, and %Bias, Eq. 8, are also included in Table S2 for more detailed~~
 258 comparison.

$$259 \quad RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (P_i - O_i)^2} \quad (34)$$

$$260 \quad MAE = \frac{1}{n} \sum_{i=1}^n (P_i - O_i) \quad (42)$$

$$261 \quad MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{P_i - O_i}{O_i} \right| \times 100 \quad (53)$$

$$262 \quad AIC = -2\ln(L) + 2k \quad (6)$$

$$263 \quad BIC = -2\ln(L) + \log(n) \log k \quad (7)$$

$$264 \quad \%BIAS = \frac{1}{n} \sum_{i=1}^n \frac{(P_i - O_i)}{O_i} \times 100 \quad (8)$$

265 ~~w~~where n is the number of samples, L is likelihood, k is the number of parameters, and P_i and O_i are the predicted
 266 and observed values, respectively. ~~were~~

267 2.6.1 Parameter tuning

268 As mentioned in [Sect. 4](#) previously, choosing a suitable strategy for parameter tuning is a crucial step in machine
269 learning particularly ~~for when~~ comparing [the performance of the](#) algorithms. Therefore, two strategies were applied
270 depending on the algorithm: 1) a grid search for RF and 2) optimisation with the DE algorithm for BRT and SVR.
271 ~~The first strategy was an exhaustive search over a defined space consisting of lower bound, upper bound and n
272 steps in between for the target parameter.~~ One major problem with applying the grid search strategy for algorithms
273 ~~that comprise continuous parameters such as BRT and SVR is that it is impossible to consider the whole continuous~~
274 ~~parameter space. Thus, the parameter combination for testing should be determined. However, this is not~~
275 ~~problematic for tuning RF in the present case since m_{try} is a parameter with discrete values. The DE algorithm~~
276 ~~however, is a stochastic approach to solve an optimisation problem that can be applied to both continuous and~~
277 ~~discrete parameters. Therefore, the target parameters in this strategy should be discrete or discretised beforehand~~
278 ~~if they are continuous (Probst et al., 2019). This strategy was applied to RF since the tuning parameter is discrete.~~
279 ~~However, the second strategy is a stochastic approach of searching over a continuous space in order to solve an~~
280 ~~optimisation problem (Qin et al., 2009) and. This method~~ is described in more detail by Storn & Price (1997).
281 Therefore, SVR and BRT ~~were~~ are optimised by this strategy as [the former algorithm has continuous parameters](#)
282 ~~and the latter one has both continuous and discrete parameters~~ they have continuous parameters. For the
283 optimisation task in the present study, the R package “DEoptim” was applied (Peterson et al., 2021). Table S1
284 shows the parameters and their tuning range for each algorithm.

285 2.6.2 Variable importance

286 Variable importance was assessed by permutation (Ließ et al., 2021). ~~Therefore, the values of each a particular~~
287 ~~covariate in the test set was were shuffled 10 times and on each occasion the prior to applying the respective model~~
288 ~~to eliminate any predictor-response relationship present with regards to that predictor trained model corresponding~~
289 ~~to that test set was applied. The variable importance corresponds to the relative increase in the test set RMSE. This~~
290 ~~process procedure was repeated 10 times for each covariate. The resulting values were averaged and the~~
291 ~~population of RMSE was averaged and its to be used for calculation of relative change to the RMSE of the original~~
292 ~~test set was calculated.~~ Thus, the variable importance of each covariate in terms of ~~percentage~~ relative change in
293 RMSE was obtained.

294 2.7 Modelling approaches

295 ~~Three approaches were designed~~ We followed a two-by-two strategy resulting in four modelling approaches to test
296 the performance of the algorithms (Table 1). ~~The models were built based on nested CV, while the train and test~~
297 ~~sets remained identical for the three algorithms to make the results comparable.~~

298 [Table 1: Modelling approaches](#)

	Dataset 1: German Agricultural Soil Inventory	Dataset 2: German Agricultural Soil Inventory + LUCAS
One-Model-Approach	AP1	AP1L
Two-Model-Approach	AP2	AP2L

299

300 ~~On the one hand, we~~The first approach (AP1) only used the SOC ~~content data~~ from the German Agricultural Soil
301 Inventory and corresponding values from the covariates ~~were used to build/train~~ the models (AP1). ~~Thus, the~~
302 ~~dataset was cross validated and used by BRT, RF and SVR to predict the SOC content of German agricultural~~
303 ~~soils. The results of this approach served as a baseline on which the model improvement for each algorithm in the~~
304 ~~other two approaches was assessed.~~

305 Due to the high variability of SOC in ~~the~~ agricultural soils of Germany, ~~we then trained~~ two separate models for
306 organic and mineral soils (AP2) ~~was developed and tested~~ to identify whether ~~it this~~ could improve model
307 performance. Accordingly, the German Agricultural Soil Inventory was subdivided by the threshold 87 g kg⁻¹ into
308 mineral and organic soils, ~~and then two were used to train separate models were trained.~~ This approach was named
309 AP2. ~~The same nested CV procedure was applied for both data subsets. The results of BRT, RF and SVR were~~
310 ~~compared to identify which one had better performance under mineral and organic soils separately. Finally, each~~
311 ~~algorithm's predicted SOC values from the two separate models was were combined, and the error metrics were~~
312 ~~calculated for the full data set to identify the impact of AP2 on model performance. The CV folds for this procedure~~
313 ~~match the one from the AP1 models.~~

314 The impact of enlarging the training set on model performance was ~~then~~ examined for both AP1 and AP2
315 ~~approaches~~. Thus, 1223 depth-extrapolated samples of the LUCAS data were added to the training sets of AP1.
316 ~~The corresponding modelling approach was and~~ named AP1L. Moreover, the same threshold (87 g kg⁻¹) was used
317 to subdivide this dataset and each soil class was included to the training set of the corresponding soil class of AP2.
318 ~~This modelling approach was then and~~ named AP2L.

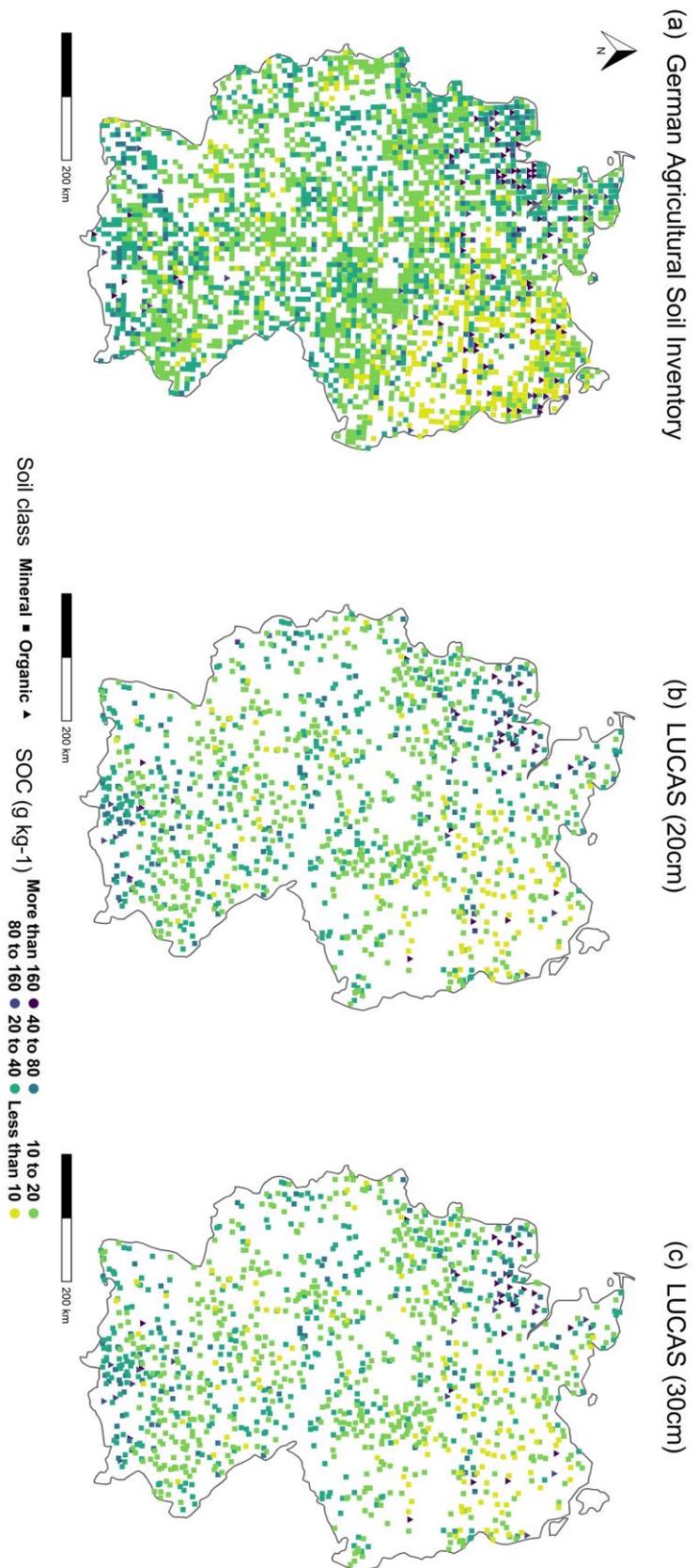
319 The test sets ~~for the model performance evaluation of the CV procedure~~ remained the same ~~for all four approaches.~~
320 ~~The models were built based on nested CV, while the train and test sets remained identical for the three algorithms~~
321 ~~to make the results comparable. Thus, the dataset was cross validated and used by BRT, RF and SVR to predict~~
322 ~~the SOC content of German agricultural soils. The results of this the AP1 approach served as a baseline on which~~
323 ~~the model improvement for each algorithm in the other two approaches were was assessed.~~

324 3 Results and Discussion

325 3.1 Comparison of algorithms on the data from the German Agricultural Soil Inventory (AP1)

326 The range of ~~the topsoil~~ SOC content ~~of topsoil~~ for the German Agricultural Soil Inventory dataset was 4 g kg⁻¹ to
327 480 g kg⁻¹, with a mean of 27 g kg⁻¹ and ~~a~~ median of 16 g kg⁻¹. Figure 1 shows the spatial distribution of the
328 ~~implemented~~ data. ~~For the first approach (AP1), BRT, RF, and SVR were applied to model SOC using data from~~
329 ~~German Agricultural Soil Inventory.~~ The RMSE and MAPE indicated ~~d~~ that SVR had a better general performance
330 than the ~~two~~ other ~~two~~ algorithms (Fig. 2). In this respect, the RMSE of SVR was 5% lower than that from RF
331 and 4% lower than that from BRT. Furthermore, its MAPE was 3% and 7% lower than that from RF and BRT
332 respectively. However, despite the difference in overall performance, the spatial distribution of relative residuals
333 indicated that all three algorithms were less accurate in ~~the northern~~ of Germany compared with the centre and
334 south of the country (Fig. 3A). This can be explained by the characteristics of this region and its higher SOC
335 variability. The northern part of Germany is ~~a~~ lowland dominated by ~~a~~ sandy soil texture from pleistocene
336 sedimentation with geomorphological structures such as ground moraines, terminal moraines and aprons (Roßkopf
337 et al., 2015). Despite general geomorphological and pedological similarities throughout the region, 1) organic soils
338 ~~in under agricultural use Germany~~ are mainly located in the north and 2) mineral soils with the lowest and ~~the~~

339 highest SOC contents are also located in the northeast and northwest respectively. Therefore, this region has the
340 ~~highest widest~~ SOC range ~~on agricultural soils~~.



341

342 Figure 1: Soil organic carbon content in the topsoil of two soil inventories: A) German Agricultural Soil Inventory (0-
 343 30 cm), B) LUCAS at its original sampling depth (0-20 cm) and, C) LUCAS after depth extrapolation (0-30 cm)





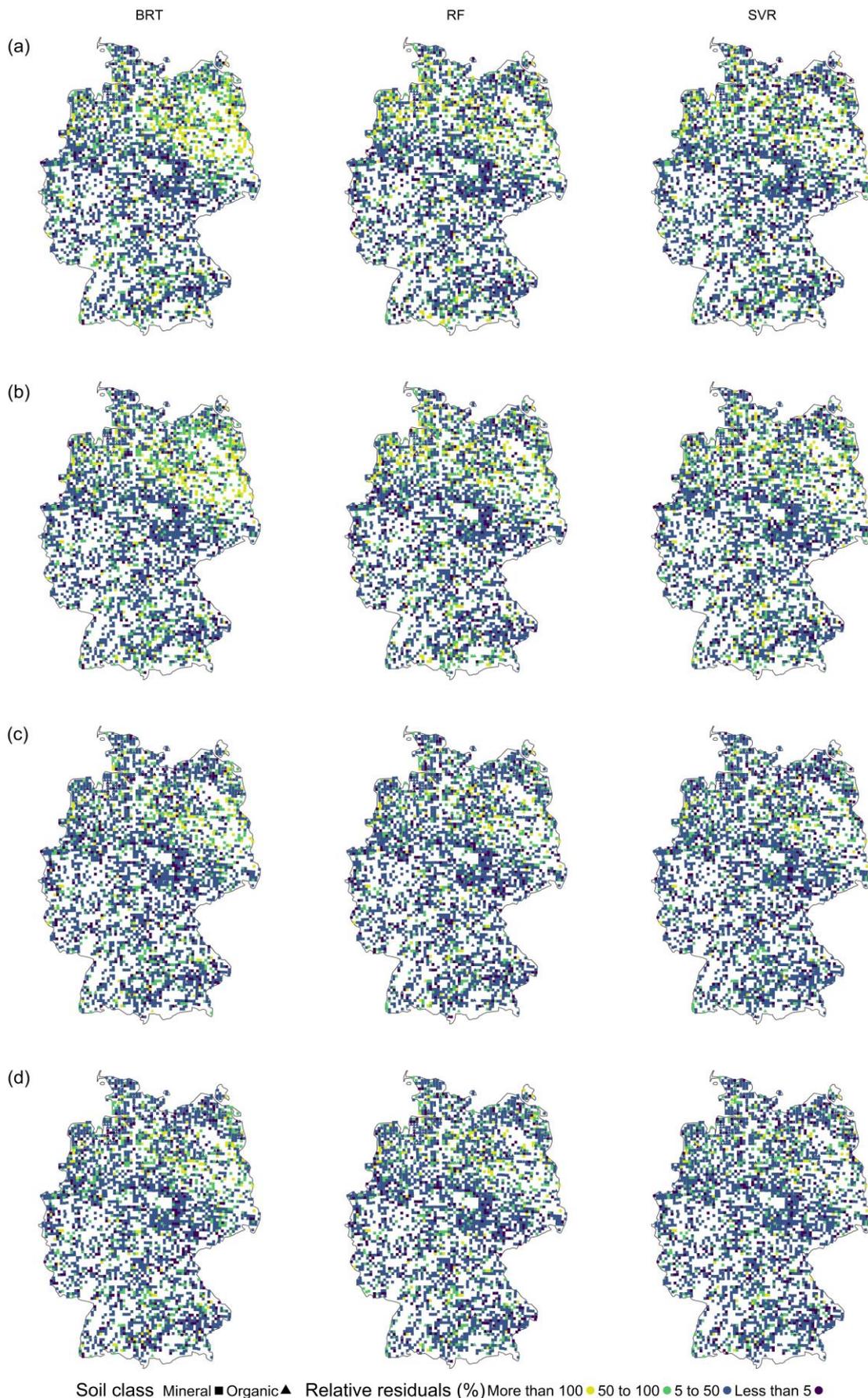
345

346 **Figure 2: Performance indicators of the three algorithms. One-model approach (without LUCAS data AP1 and with**
 347 **LUCAS data AP1L) versus the two-model approach (AP2 and AP2L) for A) RMSE (g kg^{-1}), B) MAE (g kg^{-1}) and C)**
 348 **MAPE (%). The whiskers of boxplots show 1.5 times the interquartile range. Please note that the y-axis is shortened for**
 349 **better visibility and does not display a zero. BRT = boosted regression trees, RF = random forest, and SVR = support**
 350 **vector regression.**

351 Consequently, the variable importance (Fig. 4A) indicated that the map of organic soils ~~contains the highest~~
 352 ~~available information among all~~ was the most important covariates for the algorithms. The value ~~for of the~~ variable
 353 importance for this covariate was 65% in SVR, 72% in RF and 84% in BRT. These values firstly show the crucial
 354 role of the map of organic soils for the algorithms in explaining the variability of SOC and, secondly, ~~the~~
 355 ~~comparatively greater importance of this predictor and the lower variable importance of other predictors in the~~
 356 ~~BRT model compared with the SVR model.~~ ~~how BRT mainly relies on the map of organic soils to predict SOC~~
 357 ~~compared with SVR.~~ Despite the importance of the organic soil map, the scatterplots (Fig. 5A) show that all three
 358 algorithms underpredicted the SOC of ~~the~~ organic soils and had similar heteroscedasticity patterns in their
 359 residuals. Thus, while most residuals from mineral soils followed the 1:1 line, they became more scattered in soils
 360 with a higher SOC content. The underprediction of SOC in organic soils can be explained by their ~~low-small~~
 361 sample size, resulting in a dataset with a ~~high-wide~~ SOC range and a unimodal distribution that leaves these soils
 362 in the tail. Consequently, the organic soils were underrepresented and the results were systematically pulled
 363 towards mineral soils, ~~regardless-irrespective~~ of the choice of algorithm. Different studies have shown that

364 predicting soil properties with mineral and organic soils combined can lead to underprediction or overprediction
365 of one soil class, depending on the distribution of the dataset (Brogniez et al., 2015; Guio Blanco et al., 2018;
366 Mulder et al., 2016).

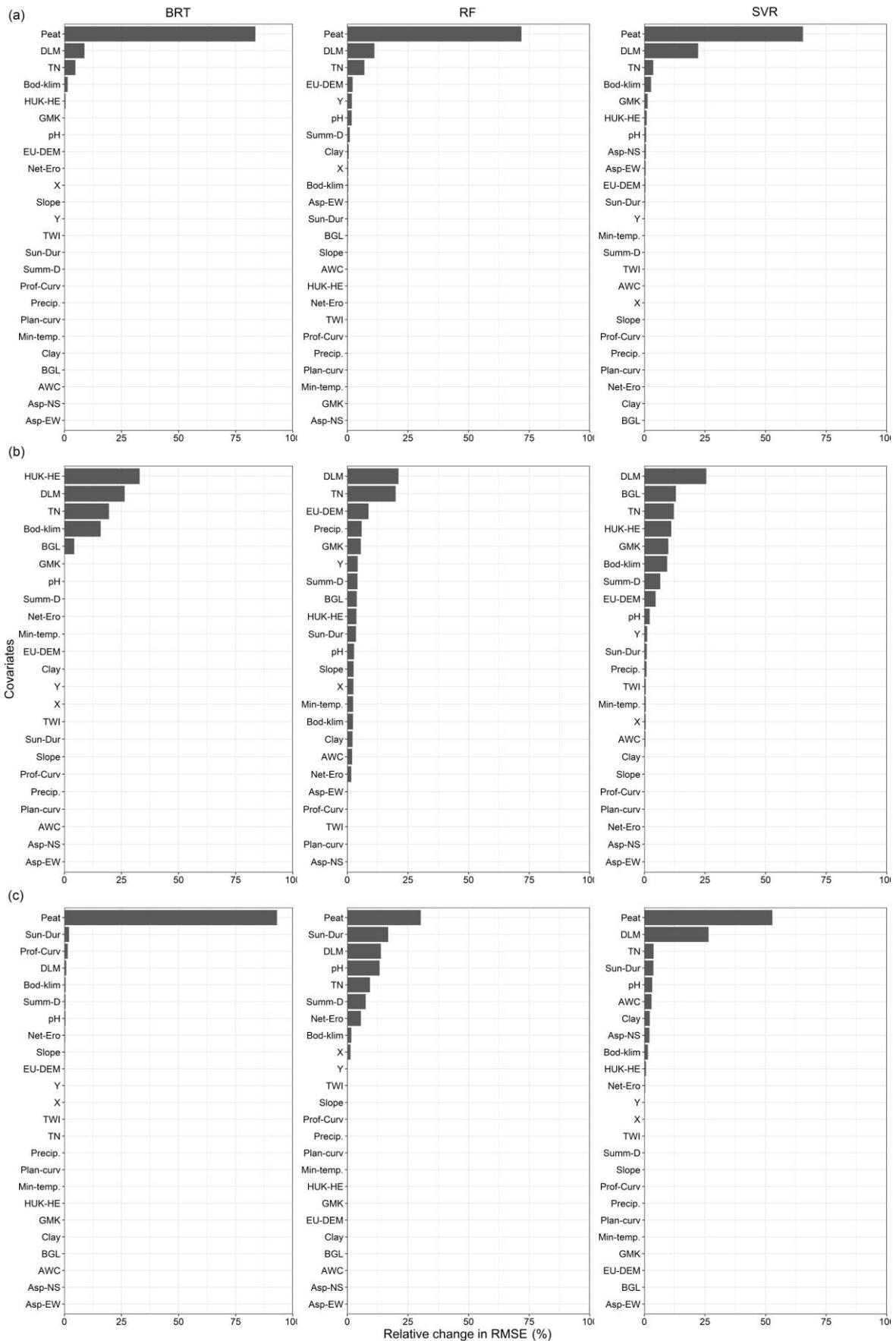
367 Although the map of organic soils was able to distinguish between the two soil classes, i.e. between mineral and
368 organic soil, it could not separate the mineral soils with a low SOC content in the northeast from those with a high
369 SOC content in the northwest. The spatial distribution of the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



381

382
383

Figure 3: Spatial distribution of relative residuals. A) AP1 approach, B) AP1L approach, C) AP2 approach and D) AP2L approach. **BRT = boosted regression trees, RF = random forest, and SVR = support vector regression.**



384

385
386
387

Figure 4: Variable importance in terms of average relative change (%) in RMSE. A) AP1, B) mineral soil subset of AP2 and C) organic soil subset of AP2. The full name for each abbreviation is presented in Table S43. BRT = boosted regression trees, RF = random forest, and SVR = support vector regression.

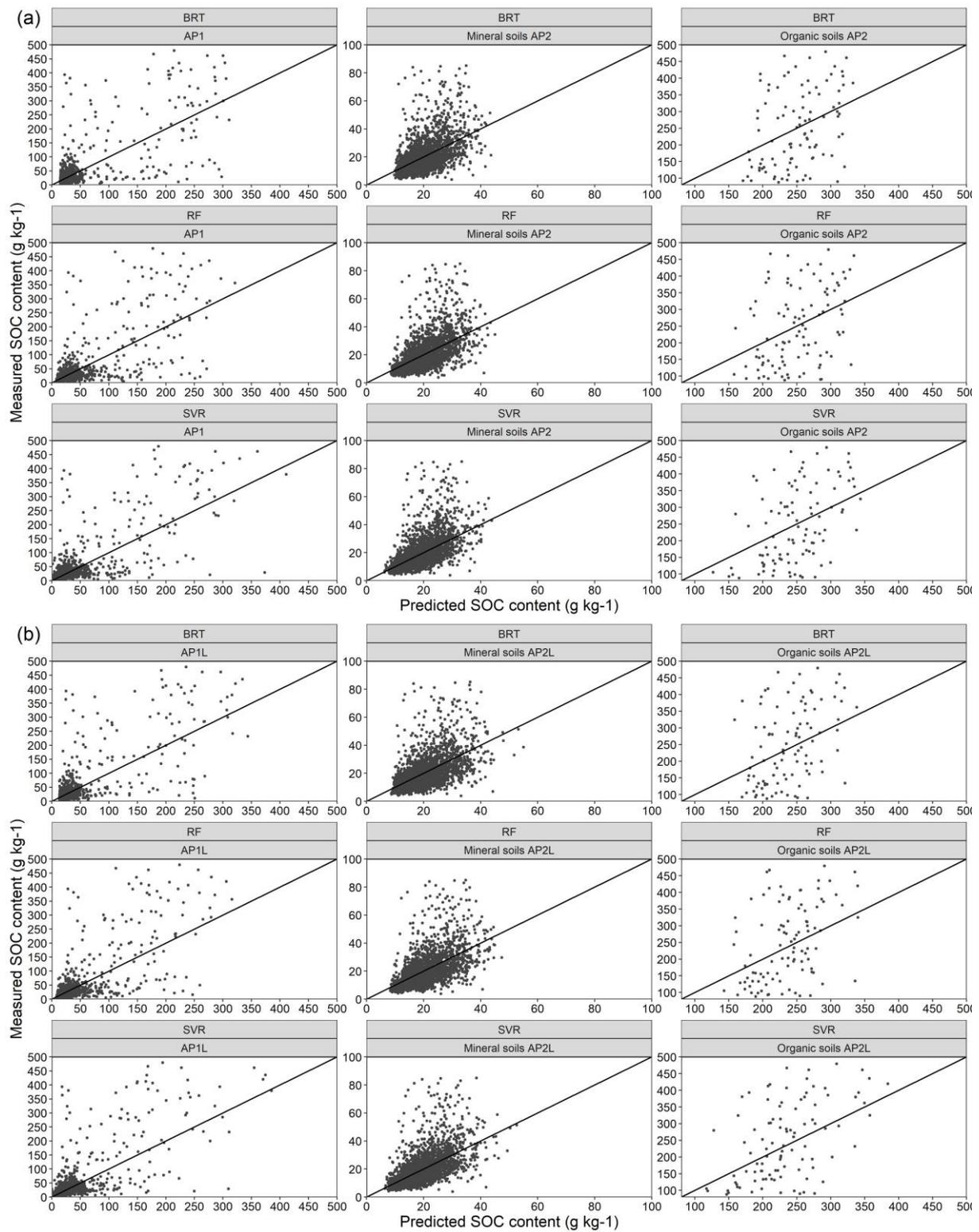
3.2 Enlarging the dataset with additional soil inventories (APIL)

A larger soil dataset may provide additional information and consequently improve model performance. This possibility was explored in the APIL approach ~~with~~ by adding ~~the~~ LUCAS data. The SOC content of LUCAS data at 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 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 median of 17 g kg⁻¹. The spatial distribution of LUCAS data at their original and extrapolated depth is shown in Figure 1.

A statistical test was performed on the residuals of models built on LUCAS data with the original and extrapolated depths. That was done to identify whether extrapolating the depth of LUCAS data to that of the German Agricultural Soil Inventory would significantly affect model performance after their inclusion in the training set. With the Shapiro-Wilk test rejecting the normality assumption of residuals of all corresponding algorithms at 20 cm and 30 cm, the non-parametric Kruskal-Wallis test showed no significant difference between the residuals at ~~both~~ either depths. Thus, the extrapolation of ~~the~~ soil depth had no significant impact on ~~the~~ data quality to regionalisze SOC. As a result, any further change in the performance of the algorithms after adding LUCAS data was due to enlargement of the training set ~~being enlarged~~. The result of the algorithms at both depths can be found in the supplementary information (Fig. S34).

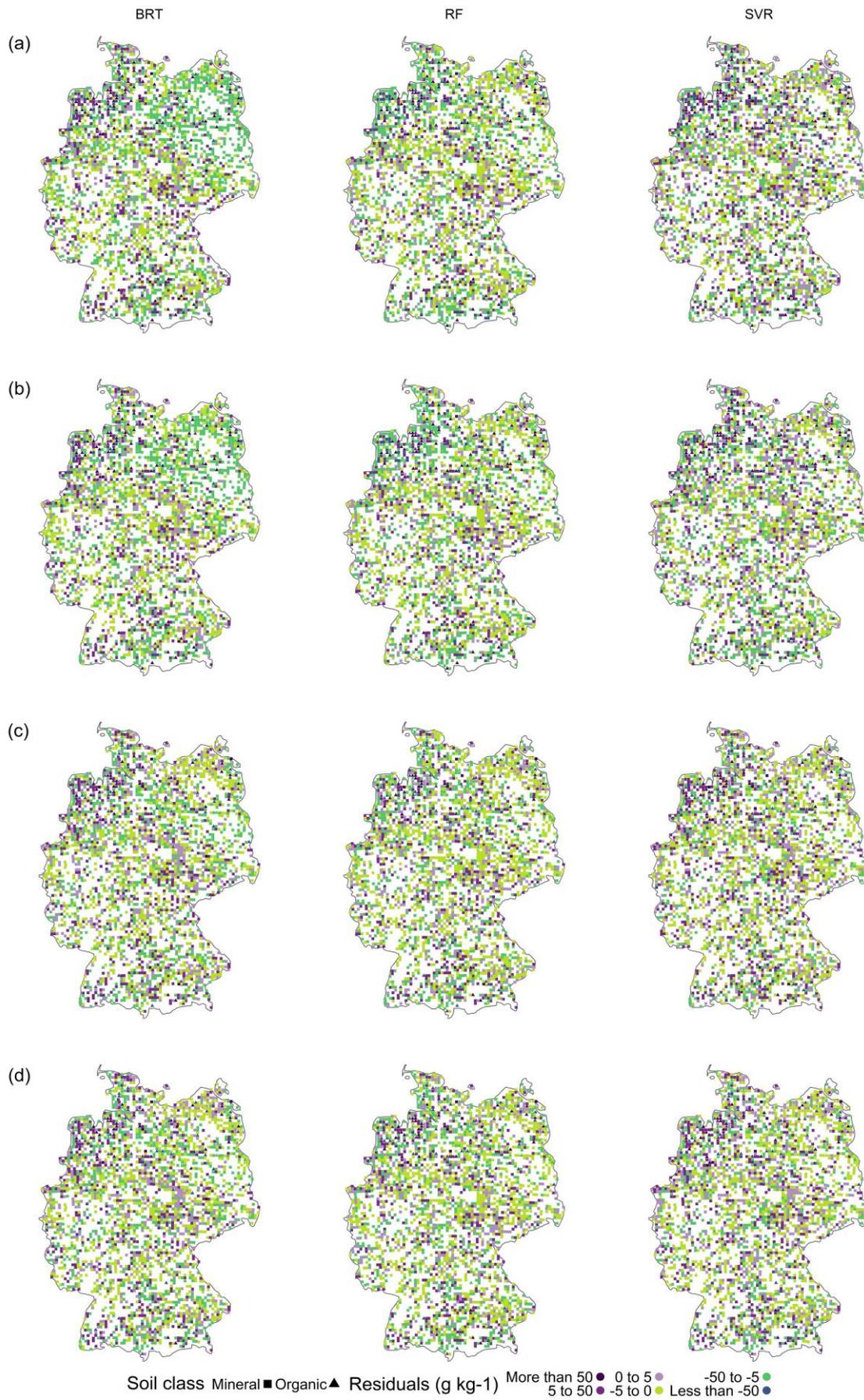
After enlarging the training set from 2278 to 3501 sampling points, BRT obtained the lowest RMSE (Fig. 2A1) and MAE among the algorithms (Fig. 2B1). A comparison of the error metrics of corresponding algorithms from the AP1 approach with those from the APIL approach showed that BRT had the highest error reduction at 7% in the MAPE and 5% in the RMSE and MAE. Furthermore, although the error metrics of RF did not improve as much as those of BRT, additional training points were still beneficial for this algorithm. However, SVR did not follow any systematic change under the APIL. Despite a 2% decrease in MAPE, the RMSE increased by 3% and MAE remained unchanged. To explore the potential explanation for this behaviour by SVR, the residuals of mineral soils were separated from those of organic soils. Additional samples reduced the RMSE in mineral soils for all algorithms by between 9% and 13%. However, this error increased by 9% in the organic subset for SVR, while it increased by just 1% for RF and even decreased by 1% for BRT. This indicated that enlarging the training set by data with similar characteristics had a greater influence on systematic error of the underrepresented soil class in SVR. This influence is understandable when considering the higher optimised ϵ in the APIL approach compared with that of ~~the~~ AP1 approach. The higher value of ϵ means that the hyperplane for the training set is less complex (Cherkassky and Ma, 2004) and more suitable for predicting most soil samples, i.e. mineral soils. Thus, when this hyperplane was fitted to the test set identical to the AP1, the generalisation performance was hindered because it could not capture the variability of samples with higher SOC values, i.e. organic soils.

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



426

427 **Figure 5: Scatterplot of residuals. A) AP1 approach and mineral and organic soils of AP2 and B) AP1L approach and**
 428 **mineral and organic soils of AP2L. BRT = boosted regression trees, RF = random forest, and SVR = support vector**
 429 **regression.**



430

431

432

Figure 6: Spatial distribution of residuals. A) AP1 approach, B) AP1L approach, C) AP2 approach and D) AP2L approach. BRT = boosted regression trees, RF = random forest, and SVR = support vector regression.

3.3 Subdividing soil inventories into mineral and organic subsets (AP2 and AP2L)

As ~~presented-outlined~~ in the sections above, the modelling of SOC content when mineral and organic soils were combined led to a systematic underprediction of soils with higher SOC values by all three algorithms, ~~regardless irrespective~~ of the number of training samples. Therefore, by implementing the AP2 approach with two models one for mineral soils and one for organic soils, a noticeable improvement in the performance of all algorithms was ~~observed,observed~~ (Table S3B), with SVR showing the best error metrics (Fig. 2A6, Fig. 2B6, Fig. 2C6). This meant 34% lower RMSE, 30% lower MAE, and 32% lower MAPE than when this algorithm was trained under the AP1 approach with one model for all soils. As the high variability of SOC was initially hard to capture, the subdivision of the dataset provided a range that better represented each soil class. This was particularly beneficial for mineral soils (ranging from 4 g kg⁻¹ to 85 g kg⁻¹) since the number of samples did not reduce drastically (only by 99 samples). Thus, the algorithms could better capture the relationship between SOC and covariates. Consequently, the overall performance improved when the underrepresented soil class was modelled separately. This is in line with the study of Rawlins et al. (2009) ~~which-that~~ recommends ~~the~~ separate modelling of mineral and organic soils.

Nonetheless, following the AP2L approach with additional data, the RMSE and MAPE of the algorithms improved by less than 2% compared with AP2 (Table S3E). However, the greatest change was observed in the MAE of SVR with a 2% improvement. Therefore, additional training samples did not ~~considerably-greatly~~ influence the performance since the majority of these samples were in mineral soils, while the limiting factor was the high variability of organic soils combined with its low number of samples. ~~NeverthelessHowever~~, an improvement was noted in relation to ~~the~~ all error metrics of SVR in the AP2L approach. This ~~was-in-contrasted~~ ~~to-with~~ when the training set was enlarged without subdividing the data, i.e. AP1L. Therefore, it further confirmed that it is more important for SVR than ~~for~~ BRT and RF to model the soil classes separately when ~~its-the~~ training set is enlarged by datasets with similar characteristics.

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

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

473 subdividing the data, the algorithms differed ~~differentiated~~ from each other by the order of covariates in their variable
474 importance (Figure 4).

475 A comparison of the error metrics of each soil class in AP2 with its counterpart in AP2L revealed that the additional
476 1177 samples had a minor influence on the performance (from zero to a maximum of 2%) of the algorithms in
477 mineral soils (Table S2). These results indicated that the German Agricultural Soil Inventory offers a good
478 representation of the spatial variability of SOC in mineral soil under agricultural use throughout the country and
479 ~~that the inclusion of including~~ more sample points ~~do did~~ not provide additional information about SOC variability
480 in this soil class.

481 However, 46 additional organic soil samples from the LUCAS dataset improved ~~the~~ MAPE and MAE by 12% and
482 6% for SVR, by 10%, and 4% for RF, and by 7% and 2% for BRT, respectively, but the RMSE of the three
483 algorithms was improved by less than 2%. Thus, additional organic samples mainly influenced the average
484 magnitude of the error. This could be explained by organic soils having a wide range of SOC and the number of
485 samples ~~was-being~~ limited. Thus, the addition of LUCAS data to the training set ~~offered-gave~~ the algorithms more
486 information about spatial variability of SOC in this soil class. Despite this limitation, SVR had the best overall
487 performance among the algorithms in AP2 and AP2L. It should be noted that training samples must span the
488 complexity of the parameter space in order for the model to be able to ~~effectively~~ match the training data ~~effectively~~
489 and ~~to-generalize~~ unseen data. ~~A s~~Small sample size can therefore negatively influence the predictive power of
490 the algorithms. This complexity can be addressed by structural risk minimisation (SRM) (Al-Anazi and Gates,
491 2012). Implementation of SRM makes SVR capable of performing well in such datasets. Other studies have
492 compared the performance of algorithms on different sample sizes ~~for-in~~ predicting soil properties and shown that
493 SVR is one of the best choices, if not the best, when the number of samples is a limiting factor (Al-Anazi and
494 Gates, 2012; Khaledian and Miller, 2020). In contrast, in a study by Zhou et al. (2021), 150 samples with different
495 sets of covariates at different resolutions were used to compare RF, BRT and SVR to predict SOC content in
496 Switzerland. Their results showed that the best-performing algorithm varied depending on the resolution and
497 covariates. However, the best performance throughout all scenarios was obtained by BRT. The discrepancy
498 between their results and the results of the present study may be due to the parameter-tuning method of the
499 algorithms, as they only used grid search, or other factors, including the spatial distribution of samples or the
500 chosen set of covariates.

501 **Table 2: Mean of error metrics of the three models for each approach.**

<u>Approach</u>	<u>Mean RMSE</u> <u>(g kg⁻¹)</u>	<u>Mean MAE</u> <u>(g kg⁻¹)</u>	<u>Mean MAPE</u> <u>(%)</u>
<u>AP1</u>	<u>32.6</u>	<u>12.3</u>	<u>49.0</u>
<u>AP1L</u>	<u>32.1</u>	<u>12.1</u>	<u>46.9</u>
<u>AP2</u>	<u>21.6</u>	<u>8.8</u>	<u>34.4</u>
<u>AP2L</u>	<u>21.3</u>	<u>8.7</u>	<u>34.3</u>

502 Overall, the change in performance across different sample sizes, different algorithms and different approaches
503 (Table S3) indicated that the most important aspect of modeling SOC content of German agricultural topsoil is a
504 two-model approach. Although combining soil inventories for more training samples can possibly improve model
505 performance, the effect was not noticeable compared to when each soil class was predicted by its dedicated model

(Table S3B and Table S3D). The advantage of two-model approach can also be seen in the average error metrics of the three models (Table 2). While the average RMSE of the models reduces by less than 1 g kg⁻¹ after enlarging the training set, the same error metrics reduces by more than 10 g kg⁻¹ in AP2 and AP2L (Table 2). Therefore, it is also recommended to consider the two-model approach in soil-landscape settings similar to Germany or situations where one-model approach cannot have good predictive performance.

The map of organic soil was used to spatially distinguish each soil class to map the SOC content of the class by its corresponding model. Figure S5 shows the spatial distribution of SOC content using the AP2L approach for the three algorithms. Although SVR captured a wider range of SOC, 2 g kg⁻¹ to 371.5 g kg⁻¹, than BRT, 8 g kg⁻¹ to 341.1 g kg⁻¹, and RF, 7.7 g kg⁻¹ to 354.6 g kg⁻¹, all three algorithms showed a relatively similar distribution of SOC content across the country particularly in mineral soils. As shown in Figure S5, organic soils are mainly distributed in the north. These soils are mostly bogs in the northwest and fens in the northeast (Roßkopf et al., 2015). There is also a small distribution of organic soil in the foothills of Alps in the south. In mineral soils, a higher SOC content is mainly found in northwest and south of the country. As explained in the previous sections, one of the main reasons for this distribution is land use since these regions are mainly under grassland while low SOC content regions are found under cropland.

4 Conclusions

The three ~~most commonly used~~ algorithms most commonly used in DSM were ~~implemented~~ applied to predict the SOC content of German agricultural soils under different approaches. Suitable tuning strategies for each algorithm ensured optimum parameter tuning and made their performance truly comparable. Machine learning algorithms was shown to be powerful ~~in~~ at modelling SOC on a national scale. However, the study showed that separate modelling of mineral and organic soils was a better approach for modelling SOC compared ~~to~~ with using just one model. Thus, this approach ~~has~~ takes priority ~~to~~ over the choice of algorithm and number of training samples. ~~We recommend~~ Further testing of this approach ~~to be further tested~~ is recommended in countries and regions that cover both of these soil classes. Nonetheless, SVR had a better performance than RF and BRT, except when the number of samples in training was increased by additional dataset. This was disadvantageous for SVR and advantageous for BRT unless mineral and organic soils were modelled separately. In general, increasing the number of training samples led to limited improvement of performance. Therefore, this approach should be ~~done~~ adopted with giving consideration of the algorithm and the characteristics of the data. Furthermore, the better performance of SVR ~~over~~ compared with that of RF and BRT was particularly highlighted when predicting SOC in organic soils. ~~Thus, this~~ The good performance of algorithm ~~SVR suggests that this algorithm~~ should therefore be taken into greater account in DSM ~~when the number of samples is limited~~.

Data availability

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

Author contribution

AS and AD conceptualised and developed the methodology of the presented work, with input from ML. AS gathered the predictors with contributions from AD. AS executed the programming, testing of existing code

543 components, formal analysis and ~~visualization~~visualisation. AG contributed to the programming. The preparation
544 of the paper was done by all authors.

545 **Competing interests**

546 The authors declare that they have no conflict of interest except the author AD is a member of the journal's
547 editorial board ~~of the journal~~.

548 **Acknowledgements**

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

552

553

554

References

- 555 Al-Anazi, A. F. and Gates, I. D.: Support vector regression to predict porosity and permeability: Effect of sample
556 size, *Comput. Geosci.*, 39, 64–76, <https://doi.org/10.1016/j.cageo.2011.06.011>, 2012.
- 557 Arrouays, D., Jolivet, C., Boulonne, L., Bodineau, G., Saby, N., & Grolleau, E.: A new projection in France: a
558 multi-institutional soil quality monitoring network, *Comptes Rendus l'Académie d'Agriculture Fr.*, 88, 93–103,
559 2002.
- 560 Awad, M. and Khanna, R.: Support Vector Regression, in: *Efficient Learning Machines*, Apress, Berkeley, CA,
561 67–80, https://doi.org/10.1007/978-1-4302-5990-9_4, 2015.
- 562 Ballabio, C., Panagos, P., and Monatanarella, L.: Mapping topsoil physical properties at European scale using the
563 LUCAS database, *Geoderma*, 261, 110–123, <https://doi.org/10.1016/j.geoderma.2015.07.006>, 2016.
- 564 Ballabio, C., Lugato, E., Fernández-Ugalde, O., Orgiazzi, A., Jones, A., Borrelli, P., Montanarella, L., and
565 Panagos, P.: Mapping LUCAS topsoil chemical properties at European scale using Gaussian process regression,
566 *Geoderma*, 355, 113912, <https://doi.org/10.1016/j.geoderma.2019.113912>, 2019.
- 567 Battineni, G., Chintalapudi, N., and Amenta, F.: Machine learning in medicine: Performance calculation of
568 dementia prediction by support vector machines (SVM), *Informatics Med. Unlocked*, 16, 100200,
569 <https://doi.org/10.1016/j.imu.2019.100200>, 2019.
- 570 Behrens, T. and Scholten, T.: Digital soil mapping in Germany - A review, *J. Plant Nutr. Soil Sci.*, 169, 434–443,
571 <https://doi.org/10.1002/jpln.200521962>, 2006.
- 572 Behrens, T., Förster, H., Scholten, T., Steinrücken, U., Spies, E. D., and Goldschmitt, M.: Digital soil mapping
573 using artificial neural networks, *J. Plant Nutr. Soil Sci.*, 168, 21–33, <https://doi.org/10.1002/jpln.200421414>, 2005.
- 574 Belon, E., Boisson, M., Deportes, I. Z., Eglin, T. K., Feix, I., Bispo, A. O., Galsomies, L., Leblond, S., and Guellier,
575 C. R.: An inventory of trace elements inputs to French agricultural soils, *Sci. Total Environ.*, 439, 87–95,
576 <https://doi.org/10.1016/j.scitotenv.2012.09.011>, 2012.
- 577 Bhadra, T., Bandyopadhyay, S., and Maulik, U.: Differential Evolution Based Optimization of SVM Parameters
578 for Meta Classifier Design, *Procedia Technol.*, 4, 50–57, <https://doi.org/10.1016/j.protcy.2012.05.006>, 2012.
- 579 Borrelli, P., Van Oost, K., Meusburger, K., Alewell, C., Lugato, E., and Panagos, P.: A step towards a holistic
580 assessment of soil degradation in Europe: Coupling on-site erosion with sediment transfer and carbon fluxes,
581 *Environ. Res.*, 161, 291–298, <https://doi.org/10.1016/j.envres.2017.11.009>, 2018.
- 582 Breiman, L.: *Randon Forests*, 1–35, 1999.
- 583 de Brogniez, D., Ballabio, C., Stevens, A., Jones, R. J. A., Montanarella, L., and van Wesemael, B.: A map of the
584 topsoil organic carbon content of Europe generated by a generalized additive model, *Eur. J. Soil Sci.*, 66, 121–
585 134, <https://doi.org/10.1111/ejss.12193>, 2015.
- 586 Burke, I. C., Yonker, C. M., Parton, W. J., Cole, C. V., Flach, K., and Schimel, D. S.: Texture, Climate, and
587 Cultivation Effects on Soil Organic Matter Content in U.S. Grassland Soils, *Soil Sci. Soc. Am. J.*, 53, 800–805,
588 <https://doi.org/10.2136/sssaj1989.03615995005300030029x>, 1989.
- 589 Camera, C., Zomeni, Z., Noller, J. S., Zissimos, A. M., Christoforou, I. C., and Bruggeman, A.: A high resolution
590 map of soil types and physical properties for Cyprus: A digital soil mapping optimization, *Geoderma*, 285, 35–49,
591 <https://doi.org/10.1016/j.geoderma.2016.09.019>, 2017.
- 592 Carter, B. J. and Ciolkosz, E. J.: Slope gradient and aspect effects on soils developed from sandstone in
593 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.
- 594 Castaldi, F., Hueni, A., Chabrilat, S., Ward, K., Buttafuoco, G., Bomans, B., Vreys, K., Brell, M., and van
595 Wesemael, B.: Evaluating the capability of the Sentinel 2 data for soil organic carbon prediction in croplands,
596 *ISPRS J. Photogramm. Remote Sens.*, 147, 267–282, <https://doi.org/10.1016/j.isprsjprs.2018.11.026>, 2019.
- 597 Chapman, S. J., Bell, J. S., Campbell, C. D., Hudson, G., Lilly, A., Nolan, A. J., Robertson, A. H. J., Potts, J. M.,
598 and Towers, W.: Comparison of soil carbon stocks in Scottish soils between 1978 and 2009, *Eur. J. Soil Sci.*, 64,
599 455–465, <https://doi.org/10.1111/ejss.12041>, 2013.
- 600 Cherkassky, V. and Ma, Y.: *FL Methods New Genetic Technology.pdf*, 17, 113–126, 2004.
- 601 Conrad, O., Bechtel, B., Bock, M., Dietrich, H., Fischer, E., Gerlitz, L., Wehberg, J., Wichmann, V., and Böhner,

602 J.: System for Automated Geoscientific Analyses (SAGA) v. 2.1.4, *Geosci. Model Dev.*, 8, 1991–2007,
603 <https://doi.org/10.5194/gmd-8-1991-2015>, 2015.

604 Deng, S., Wang, C., Wang, M., and Sun, Z.: A gradient boosting decision tree approach for insider trading
605 identification: An empirical model evaluation of China stock market, *Appl. Soft Comput. J.*, 83, 105652,
606 <https://doi.org/10.1016/j.asoc.2019.105652>, 2019.

607 DWD Climate Data Center (CDC): Multi-annual grids of annual sunshine duration over Germany 1981-
608 2010, version v1.0, 2017.

609 DWD Climate Data Center (CDC): Multi-annual grids of monthly averaged daily minimum air temperature (2m)
610 over Germany, version v1.0, 2018a.

611 DWD Climate Data Center (CDC): Multi-annual grids of number of summer days over Germany, version v1.0,
612 2018b.

613 DWD Climate Data Center (CDC): Multi-annual grids of precipitation height over Germany 1981-2010, version
614 v1.0, 2018c.

615 Environmental Systems Research Institute (ESRI): ArcGIS 10.2 for Desktop, 2013.

616 European Union Copernicus Land Monitoring Service, and E. E. A. (EEA): European Digital Elevation Model
617 (EU-DEM), Version 1.1, 2016.

618 Eurostat grid generation tool for ArcGIS: [https://www.efgs.info/information-base/best-practices/tools/eurostat-
619 grid-generation-tool-arcgis/](https://www.efgs.info/information-base/best-practices/tools/eurostat-grid-generation-tool-arcgis/).

620 Federal Agency for Cartography and Geodesy (BKG): Digitales Basis-Landschaftsmodell (Basis-DLM), Leipzig,
621 2019.

622 Federal Institute for Geosciences and Natural Resources (BGR): Geomorphographic Map of Germany
623 (GMK1000), Hanover, 2007.

624 Federal Institute for Geosciences and Natural Resources (BGR): Soil scapes in Germany 1:5,000,000 (BGL5000),
625 Hanover, 2008.

626 Federal Institute for Geosciences and Natural Resources (BGR) and German State Geological Surveys (SDG):
627 Hydrogeological Map of Germany 1:250,000 (HÜK250), Hanover, 2019.

628 Forkuor, G., Hounkpatin, O. K. L., Welp, G., and Thiel, M.: High resolution mapping of soil properties using
629 Remote Sensing variables in south-western Burkina Faso: A comparison of machine learning and multiple linear
630 regression models, *PLoS One*, 12, 1–21, <https://doi.org/10.1371/journal.pone.0170478>, 2017.

631 Friedman, J., Tibshirani, R., and Hastie, T.: Additive logistic regression: a statistical view of boosting (With
632 discussion and a rejoinder by the authors), *Ann. Stat.*, 28, 337–407, <https://doi.org/10.1214/aos/1016120463>, 2000.

633 Friedman, J., Hastie, T., and Tibshirani, R.: *The Elements of Statistical Learning*, second., Springer New York,
634 New York, NY, 158-61 pp., <https://doi.org/10.1007/b94608>, 2009.

635 Friedman, J. H.: Stochastic gradient boosting, *Comput. Stat. Data Anal.*, 38, 367–378,
636 [https://doi.org/10.1016/S0167-9473\(01\)00065-2](https://doi.org/10.1016/S0167-9473(01)00065-2), 2002.

637 Gebauer, A., Ellinger, M., Brito Gomez, V., and Ließ, M.: Development of pedotransfer functions for water
638 retention in tropical mountain soil landscapes: Spotlight on parameter tuning in machine learning, 6, 215–229,
639 <https://doi.org/10.5194/soil-6-215-2020>, 2020.

640 Greenwell, B., Boehmke, B., and Cunningham, J.: Package “gbm” - Generalized Boosted Regression Models,
641 CRAN Repos., 39, 2019.

642 Guevara, M., Olmedo, G. F., Stell, E., Yigini, Y., Aguilar Duarte, Y., Arellano Hernández, C., Arévalo, G., Arroyo-
643 Cruz, C. E., Bolivar, A., Bunning, S., Bustamante Cañas, N., Cruz-Gaistardo, C. O., Davila, F., Dell Acqua, M.,
644 Encina, A., Figueredo Tacona, H., Fontes, F., Hernández Herrera, J. A., Ibelle Navarro, A. R., Loayza, V.,
645 Manueles, A., Mendoza Jara, F., Olivera, C., Osorio Hermosilla, R., Pereira, G., Prieto, P., Alexis Ramos, I., Rey
646 Brina, J. C., Rivera, R., Rodríguez-Rodríguez, J., Roopnarine, R., Rosales Ibarra, A., Rosales Riveiro, K. A.,
647 Schulz, G. A., Spence, A., Vasques, G., Vargas, R., and Vargas, R.: No Silver Bullet for Digital Soil Mapping:
648 Country-specific Soil Organic Carbon Estimates across Latin America, *SOIL Discuss.*, 1–20,
649 <https://doi.org/10.5194/soil-2017-40>, 2018.

- 650 Guio Blanco, C. M., Brito Gomez, V. M., Crespo, P., and Ließ, M.: Spatial prediction of soil water retention in a
651 Páramo landscape: Methodological insight into machine learning using random forest, *Geoderma*, 316, 100–114,
652 <https://doi.org/10.1016/j.geoderma.2017.12.002>, 2018.
- 653 Hawkins, D. M., Basak, S. C., and Mills, D.: Assessing model fit by cross-validation, *J. Chem. Inf. Comput. Sci.*,
654 43, 579–586, <https://doi.org/10.1021/ci025626i>, 2003.
- 655 Hornik, K., Weingessel, A., Leisch, F., and Davidmeyer-projectorg, M. D. M.: Package ‘e1071,’ 2021.
- 656 Hoyle, F. C., Baldock, J. A., and Murphy, D. V.: Soil Organic Carbon – Role in Rainfed Farming Systems, *Rainfed
657 Farming Syst.*, <https://doi.org/10.1007/978-1-4020-9132-2>, 2011.
- 658 Khaledian, Y. and Miller, B. A.: Selecting appropriate machine learning methods for digital soil mapping, *Appl.
659 Math. Model.*, 81, 401–418, <https://doi.org/10.1016/j.apm.2019.12.016>, 2020.
- 660 Kuhn, M. and Johnson, K.: *Applied predictive modeling*, 1st ed., Springer-Verlag, New York, 1–600 pp.,
661 <https://doi.org/10.1007/978-1-4614-6849-3>, 2013.
- 662 Lal, R.: Soil carbon sequestration impacts on global climate change and food security, *Science (80-.)*, 304, 1623–
663 1627, <https://doi.org/10.1126/science.1097396>, 2004.
- 664 Li, T., Zhang, H., Wang, X., Cheng, S., Fang, H., Liu, G., and Yuan, W.: Soil erosion affects variations of soil
665 organic carbon and soil respiration along a slope in Northeast China, *Ecol. Process.*, 8,
666 <https://doi.org/10.1186/s13717-019-0184-6>, 2019.
- 667 Li, X., Ding, J., Liu, J., Ge, X., and Zhang, J.: Digital mapping of soil organic carbon using sentinel series data: A
668 case study of the ebinur lake watershed in xinjiang, *Remote Sens.*, 13, 1–19, <https://doi.org/10.3390/rs13040769>,
669 2021.
- 670 Liang, W., Zhang, L., and Wang, M.: The chaos differential evolution optimization algorithm and its application
671 to support vector regression machine, *J. Softw.*, 6, 1297–1304, <https://doi.org/10.4304/jsw.6.7.1297-1304>, 2011.
- 672 Ließ, M., Gebauer, A., and Don, A.: Machine Learning With GA Optimization to Model the Agricultural Soil-
673 Landscape of Germany: An Approach Involving Soil Functional Types With Their Multivariate Parameter
674 Distributions Along the Depth Profile, *Front. Environ. Sci.*, 9, 1–24, <https://doi.org/10.3389/fenvs.2021.692959>,
675 2021.
- 676 Malik, A. A., Puissant, J., Buckeridge, K. M., Goodall, T., Jehmlich, N., Chowdhury, S., Gweon, H. S., Peyton, J.
677 M., Mason, K. E., van Agtmaal, M., Bland, A., Clark, I. M., Whitaker, J., Pywell, R. F., Ostle, N., Gleixner, G.,
678 and Griffiths, R. I.: Land use driven change in soil pH affects microbial carbon cycling processes, *Nat. Commun.*,
679 9, 1–10, <https://doi.org/10.1038/s41467-018-05980-1>, 2018.
- 680 Martin, M. P., Orton, T. G., Llacarce, E., Meersmans, J., Saby, N. P. A., Paroissien, J. B., Jolivet, C., Boulonne,
681 L., and Arrouays, D.: Evaluation of modelling approaches for predicting the spatial distribution of soil organic
682 carbon stocks at the national scale, *Geoderma*, 223–225, 97–107, <https://doi.org/10.1016/j.geoderma.2014.01.005>,
683 2014.
- 684 McBratney, A. B., Mendonça Santos, M. L., and Minasny, B.: On digital soil mapping, 3–52 pp.,
685 [https://doi.org/10.1016/S0016-7061\(03\)00223-4](https://doi.org/10.1016/S0016-7061(03)00223-4), 2003.
- 686 Meersmans, J., Martin, M. P., Llacarce, E., De Baets, S., Jolivet, C., Boulonne, L., Lehmann, S., Saby, N. P. A.,
687 Bispo, A., and Arrouays, D.: A high resolution map of French soil organic carbon, *Agron. Sustain. Dev.*, 32, 841–
688 851, <https://doi.org/10.1007/s13593-012-0086-9>, 2012a.
- 689 Meersmans, J., Martin, M. P., De Ridder, F., Llacarce, E., Wetterlind, J., De Baets, S., Bas, C. Le, Louis, B. P.,
690 Orton, T. G., Bispo, A., and Arrouays, D.: A novel soil organic C model using climate, soil type and management
691 data at the national scale in France, *Agron. Sustain. Dev.*, 32, 873–888, [https://doi.org/10.1007/s13593-012-0085-
692 x](https://doi.org/10.1007/s13593-012-0085-x), 2012b.
- 693 Minasny, B., McBratney, A. B., Malone, B. P., and Wheeler, I.: *Digital Mapping of Soil Carbon*, Elsevier, 1–47
694 pp., <https://doi.org/10.1016/B978-0-12-405942-9.00001-3>, 2013.
- 695 Mulder, V. L., Lacoste, M., Richer-de-Forges, A. C., Martin, M. P., and Arrouays, D.: National versus global
696 modelling the 3D distribution of soil organic carbon in mainland France, *Geoderma*, 263, 16–34,
697 <https://doi.org/10.1016/J.GEODERMA.2015.08.035>, 2016.
- 698 Padarian, J., Minasny, B., and McBratney, A. B.: *Machine learning and soil sciences: A review aided by machine*

- 699 learning tools, 6, 35–52, <https://doi.org/10.5194/soil-6-35-2020>, 2020.
- 700 Pei, T., Qin, C. Z., Zhu, A. X., Yang, L., Luo, M., Li, B., and Zhou, C.: Mapping soil organic matter using the
701 topographic wetness index: A comparative study based on different flow-direction algorithms and kriging
702 methods, *Ecol. Indic.*, 10, 610–619, <https://doi.org/10.1016/j.ecolind.2009.10.005>, 2010.
- 703 Peterson, B., Ulrich, J., and Boudt, K.: Package ‘DEoptim’, <https://doi.org/10.18637/jss.v040.i06>, 2021.
- 704 Poeplau, C., Bolinder, M. A., Eriksson, J., Lundblad, M., and Kätterer, T.: Positive trends in organic carbon storage
705 in Swedish agricultural soils due to unexpected socio-economic drivers, 12, 3241–3251,
706 <https://doi.org/10.5194/bg-12-3241-2015>, 2015.
- 707 Poeplau, C., Jacobs, A., Don, A., Vos, C., Schneider, F., Wittnebel, M., Tiemeyer, B., Heidkamp, A., Prietz, R.,
708 and Flessa, H.: Stocks of organic carbon in German agricultural soils—Key results of the first comprehensive
709 inventory, *J. Plant Nutr. Soil Sci.*, 183, 665–681, <https://doi.org/10.1002/jpln.202000113>, 2020.
- 710 Prechtel, A., Von Łtzw, M., Schneider, B. U., Bens, O., Bannick, C. G., Kögel-Knabner, I., and Hüttl, R. F.:
711 Organic Carbon in soils of Germany: Status quo and the need for new data to evaluate potentials and trends of soil
712 carbon sequestration, *J. Plant Nutr. Soil Sci.*, 172, 601–614, <https://doi.org/10.1002/jpln.200900034>, 2009.
- 713 Ramifehiarivo, N., Brossard, M., Grinand, C., Andriamananjara, A., Razafimbelo, T., Rasolohery, A.,
714 Razafimahatratra, H., Seyler, F., Ranaivoson, N., Rabenarivo, M., Albrecht, A., Razafindrabe, F., and
715 Razakamanarivo, H.: Mapping soil organic carbon on a national scale: Towards an improved and updated map of
716 Madagascar, *Geoderma Reg.*, 9, 29–38, <https://doi.org/10.1016/j.geodrs.2016.12.002>, 2017.
- 717 Rawlins, B. G., Marchant, B. P., Smyth, D., Scheib, C., Lark, R. M., and Jordan, C.: Airborne radiometric survey
718 data and a DTM as covariates for regional scale mapping of soil organic carbon across Northern Ireland, *Eur. J.*
719 *Soil Sci.*, 60, 44–54, <https://doi.org/10.1111/j.1365-2389.2008.01092.x>, 2009.
- 720 Reeves, D. W.: The role of soil organic matter in maintaining soil quality in continuous cropping systems, *Soil*
721 *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.
- 722 Ritchie, J. C., McCarty, G. W., Venteris, E. R., and Kaspar, T. C.: Soil and soil organic carbon redistribution on
723 the landscape, 89, 163–171, <https://doi.org/10.1016/j.geomorph.2006.07.021>, 2007.
- 724 Roßberg, D., Michel, V., Graf, R., Neukampf, R.: Definition von Boden-Klima-Räumen für die Bundesrepublik
725 Deutschland, *Nachrichtenblatt des Dtsch. Pflanzenschutzdienstes*, 59, 155–161, 2007.
- 726 Roßkopf, N., Fell, H., and Zeitz, J.: Organic soils in Germany, their distribution and carbon stocks, 133, 157–170,
727 <https://doi.org/10.1016/j.catena.2015.05.004>, 2015.
- 728 Santos, C. E. da S., Sampaio, R. C., Coelho, L. dos S., Bestarsd, G. A., and Llanos, C. H.: Multi-objective adaptive
729 differential evolution for SVM/SVR hyperparameters selection, *Pattern Recognit.*, 110, 107649,
730 <https://doi.org/10.1016/j.patcog.2020.107649>, 2021.
- 731 Schapire, R. E.: The Boosting Approach to Machine Learning: An Overview, 149–171,
732 https://doi.org/10.1007/978-0-387-21579-2_9, 2003.
- 733 Schneider, F., Amelung, W., and Don, A.: Origin of carbon in agricultural soil profiles deduced from depth
734 gradients of C:N ratios, carbon fractions, $\delta^{13}\text{C}$ and $\delta^{15}\text{N}$ values, *Plant Soil*, 460, 123–148,
735 <https://doi.org/10.1007/s11104-020-04769-w>, 2021.
- 736 Smola, A. J. and Schölkopf, B.: A tutorial on support vector regression, *Stat. Comput.*, 14, 199–222,
737 <https://doi.org/10.1023/B:STCO.0000035301.49549.88>, 2004.
- 738 Storn, R. and Price, K.: Differential Evolution - A Simple and Efficient Heuristic for Global Optimization over
739 Continuous Spaces, *J. Glob. Optim.*, 11, 341–359, <https://doi.org/10.1023/A:1008202821328>, 1997.
- 740 Taghizadeh-Toosi, A., Olesen, J. E., Kristensen, K., Elsgaard, L., Østergaard, H. S., Lægdsmand, M., Greve, M.
741 H., and Christensen, B. T.: Changes in carbon stocks of Danish agricultural mineral soils between 1986 and 2009,
742 *Eur. J. Soil Sci.*, 65, 730–740, <https://doi.org/10.1111/ejss.12169>, 2014.
- 743 Tóth, G., Jones, A., and Montanarella, L.: LUCAS Topsoil Survey: Methodology, Data, and Results,
744 <https://doi.org/10.2788/97922>, 2013.
- 745 Tziachris, P., Aschonitis, V., Chatzistathis, T., Papadopoulou, M., and Doukas, I. J. D.: Comparing machine
746 learning models and hybrid geostatistical methods using environmental and soil covariates for soil pH prediction,

- 747 ISPRS Int. J. Geo-Information, 9, <https://doi.org/10.3390/ijgi9040276>, 2020.
- 748 Varma, S. and Simon, R.: Bias in error estimation when using cross-validation for model selection, BMC
749 Bioinformatics, 7, 1–8, <https://doi.org/10.1186/1471-2105-7-91>, 2006.
- 750 Wadoux, A. M. J. C., Minasny, B., and McBratney, A. B.: Machine learning for digital soil mapping: Applications,
751 challenges and suggested solutions, Earth-Science Rev., 210, <https://doi.org/10.1016/j.earscirev.2020.103359>,
752 2020.
- 753 Wang, S., Xu, L., Zhuang, Q., and He, N.: Investigating the spatio-temporal variability of soil organic carbon
754 stocks in different ecosystems of China, Sci. Total Environ., 758, <https://doi.org/10.1016/j.scitotenv.2020.143644>,
755 2021.
- 756 Wang, X., Zhang, Y., Atkinson, P. M., and Yao, H.: Predicting soil organic carbon content in Spain by combining
757 Landsat TM and ALOS PALSAR images, Int. J. Appl. Earth Obs. Geoinf., 92, 102182,
758 <https://doi.org/10.1016/j.jag.2020.102182>, 2020.
- 759 Ward, K. J., Chabrilat, S., Neumann, C., and Foerster, S.: A remote sensing adapted approach for soil organic
760 carbon prediction based on the spectrally clustered LUCAS soil database, Geoderma, 353, 297–307,
761 <https://doi.org/10.1016/j.geoderma.2019.07.010>, 2019.
- 762 Were, K., Bui, D. T., Dick, Ø. B., and Singh, B. R.: A comparative assessment of support vector regression,
763 artificial neural networks, and random forests for predicting and mapping soil organic carbon stocks across an
764 Afromontane landscape, Ecol. Indic., 52, 394–403, <https://doi.org/10.1016/j.ecolind.2014.12.028>, 2015.
- 765 Wiesmeier, M., Spörlein, P., Geuß, U., Hangen, E., Haug, S., Reischl, A., Schilling, B., von Lützw, M., and
766 Kögel-Knabner, I.: Soil organic carbon stocks in southeast Germany (Bavaria) as affected by land use, soil type
767 and sampling depth, Glob. Chang. Biol., 18, 2233–2245, <https://doi.org/10.1111/j.1365-2486.2012.02699.x>, 2012.
- 768 Wright, M. N. and Ziegler, A.: Ranger: A fast implementation of random forests for high dimensional data in C++
769 and R, J. Stat. Softw., 77, 1–17, <https://doi.org/10.18637/jss.v077.i01>, 2017.
- 770 Yu, D., Hu, F., Zhang, K., Liu, L., and Li, D.: Available water capacity and organic carbon storage profiles in soils
771 developed from dark brown soil to boggy soil in Changbai Mountains, China, Soil Water Res., 16, 11–21,
772 <https://doi.org/10.17221/150/2019-SWR>, 2021.
- 773 Zhang, J., Niu, Q., Li, K., and Irwin, G. W.: Model selection in SVMs using Differential Evolution, IFAC, 14717–
774 14722 pp., <https://doi.org/10.3182/20110828-6-IT-1002.00584>, 2011.
- 775 Zhong, Z., Chen, Z., Xu, Y., Ren, C., Yang, G., Han, X., Ren, G., and Feng, Y.: Relationship between soil organic
776 carbon stocks and clay content under different climatic conditions in Central China, 9, 1–14,
777 <https://doi.org/10.3390/f9100598>, 2018.
- 778 Zhou, T., Geng, Y., Ji, C., Xu, X., Wang, H., Pan, J., Bumberger, J., Haase, D., and Lausch, A.: Prediction of soil
779 organic carbon and the C:N ratio on a national scale using machine learning and satellite data: A comparison
780 between Sentinel-2, Sentinel-3 and Landsat-8 images, Sci. Total Environ., 755,
781 <https://doi.org/10.1016/j.scitotenv.2020.142661>, 2021.
- 782

1 Spatial prediction of organic carbon in German agricultural topsoil
2 using machine learning algorithms~~Performance of three machine~~
3 ~~learning algorithms for predicting soil organic carbon in German~~
4 ~~agricultural soil~~

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

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

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

8

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

10 **Supplements**

11 **Table S1: The range of parameters for tuning in full dataset (AP1 and AP1L) and mineral and organic soil subsets**
12 **(AP2 and AP2L). Tuning parameter ranges corresponding to the models trained by different algorithms. The ranges**
13 **differ considering the one-model (full dataset, AP1 & AP1L) or two-model approach (mineral and organic data**
14 **subset, AP2 & AP2L). BRT = boosted regression trees, RF = random forest, and SVR = support vector regression.**

<u>Algorithm</u>	<u>Tuning parameter</u>	<u>Full dataset</u>	<u>Mineral soil data subset</u>	<u>Organic soil data subset</u>
SVR	C	1-100	1-50	1-200
	epsilon	0-5	0-1	0-5
	gamma	0.001-1	0.001-1	0.001-1
RF	mtry	3-13	3-13	3-13
BRT	number of trees	100-3000	100-3000	100-3000
	shrinkage	0.001-0.1	0.001-0.1	0.001-0.1
	interaction depth	1-5	1-5	1-5
	bag fraction	0.5-0.9	0.5-0.9	0.5-0.9

15 **Table S2: Error metrics of the algorithms**
16 **Predictive model performance of the models trained with different machine**
17 **learning algorithms and datasets: A) built on the German Agricultural Soil Inventory, B) including LUCAS data in**
the training set. BRT = boosted regression trees, RF = random forest, and SVR = support vector regression.

	<u>Algorithm</u>	<u>RMSE</u>	<u>MAE</u>	<u>%MAPE</u>	<u>%Bias</u>	<u>AIC</u>	<u>BIC</u>	<u>Approach</u>
A	<u>BRT</u>	<u>32.9</u>	<u>12.4</u>	<u>50.9</u>	<u>-32</u>	<u>14865</u>	<u>14889</u>	<u>AP1</u>
	<u>RF</u>	<u>33.2</u>	<u>12.3</u>	<u>48.6</u>	<u>-30</u>	<u>14913</u>	<u>14919</u>	<u>AP1</u>
	<u>SVR</u>	<u>31.6</u>	<u>12.3</u>	<u>47.4</u>	<u>-20</u>	<u>14643</u>	<u>14661</u>	<u>AP1</u>
	<u>BRT</u>	<u>9.5</u>	<u>6.2</u>	<u>35.9</u>	<u>-20</u>	<u>7500</u>	<u>7524</u>	<u>Mineral</u>
	<u>RF</u>	<u>9.1</u>	<u>5.9</u>	<u>34</u>	<u>-20</u>	<u>7288</u>	<u>7294</u>	<u>Mineral</u>
	<u>SVR</u>	<u>9.2</u>	<u>5.8</u>	<u>31.8</u>	<u>-10</u>	<u>7331</u>	<u>7349</u>	<u>Mineral</u>
	<u>BRT</u>	<u>107</u>	<u>90.4</u>	<u>48.5</u>	<u>-26</u>	<u>757</u>	<u>768</u>	<u>Organic</u>
	<u>RF</u>	<u>106.1</u>	<u>89.3</u>	<u>48.2</u>	<u>-28</u>	<u>750</u>	<u>753</u>	<u>Organic</u>
	<u>SVR</u>	<u>101.7</u>	<u>86.9</u>	<u>45.6</u>	<u>-22</u>	<u>746</u>	<u>754</u>	<u>Organic</u>
	<u>BRT</u>	<u>22</u>	<u>9.1</u>	<u>36.3</u>	<u>-20</u>	<u>12578</u>	<u>12602</u>	<u>AP2</u>
	<u>RF</u>	<u>21.7</u>	<u>8.8</u>	<u>34.5</u>	<u>-20</u>	<u>12496</u>	<u>12502</u>	<u>AP2</u>
	<u>SVR</u>	<u>21</u>	<u>8.6</u>	<u>32.3</u>	<u>-10</u>	<u>12310</u>	<u>12328</u>	<u>AP2</u>

	<u>Algorithm</u>	<u>RMSE</u>	<u>MAE</u>	<u>%MAPE</u>	<u>%Bias</u>	<u>AIC</u>	<u>BIC</u>	<u>Approach</u>
<u>B</u>	<u>BRT</u>	<u>31.3</u>	<u>11.8</u>	<u>47.4</u>	<u>-30</u>	<u>14568</u>	<u>14592</u>	<u>AP1L</u>
	<u>RF</u>	<u>32.5</u>	<u>12.1</u>	<u>46.8</u>	<u>-30</u>	<u>14754</u>	<u>14759</u>	<u>AP1L</u>
	<u>SVR</u>	<u>32.6</u>	<u>12.3</u>	<u>46.4</u>	<u>-20</u>	<u>14775</u>	<u>14792</u>	<u>AP1L</u>
	<u>BRT</u>	<u>9.4</u>	<u>6.2</u>	<u>35.6</u>	<u>-20</u>	<u>7429</u>	<u>7453</u>	<u>Mineral</u>
	<u>RF</u>	<u>9.1</u>	<u>6</u>	<u>34.6</u>	<u>-20</u>	<u>7268</u>	<u>7274</u>	<u>Mineral</u>
	<u>SVR</u>	<u>9.1</u>	<u>5.8</u>	<u>31.7</u>	<u>-10</u>	<u>7275</u>	<u>7293</u>	<u>Mineral</u>
	<u>BRT</u>	<u>105.4</u>	<u>88.4</u>	<u>45</u>	<u>-20</u>	<u>754</u>	<u>765</u>	<u>Organic</u>
	<u>RF</u>	<u>104.1</u>	<u>86.2</u>	<u>43.5</u>	<u>-20</u>	<u>745</u>	<u>748</u>	<u>Organic</u>
	<u>SVR</u>	<u>100.2</u>	<u>81.7</u>	<u>40.2</u>	<u>-12</u>	<u>741</u>	<u>749</u>	<u>Organic</u>
	<u>BRT</u>	<u>21.7</u>	<u>9</u>	<u>36</u>	<u>-20</u>	<u>12486</u>	<u>12510</u>	<u>AP2L</u>
	<u>RF</u>	<u>21.4</u>	<u>8.7</u>	<u>34.9</u>	<u>-20</u>	<u>12379</u>	<u>12385</u>	<u>AP2L</u>
	<u>SVR</u>	<u>20.7</u>	<u>8.4</u>	<u>31.9</u>	<u>-10</u>	<u>12191</u>	<u>12209</u>	<u>AP2L</u>

18

19 **Table S3: Percent change in predictive model performance comparing models trained with different machine learning**
20 **algorithms and data sets: A) and B) comparison of models trained by using data from the –German Agricultural Soil**
21 **Inventory, C) and D) comparison of models trained by using data from the German Agricultural Soil Inventory and**
22 **LUCAS, A) and C) comparison of models trained with comparison with regards to the different machine learning**
23 **algorithms, B) and D) comparison of the one-model approach (AP1) to the two-model approach (AP2), E) comparison**
24 **of the approaches before and after including LUCAS, BRT = boosted regression trees, RF = random forest, and SVR =**
25 **support vector regression.**

	<u>Algorithm</u>	<u>RMSE (%)</u>	<u>MAE (%)</u>	<u>MAPE (%)</u>	<u>Approach</u>
<u>A</u>	<u>BRT to RF</u>	<u>0.9</u>	<u>-0.8</u>	<u>-4.5</u>	<u>AP1</u>
	<u>RF to SVR</u>	<u>-4.8</u>	<u>0.0</u>	<u>-2.5</u>	<u>AP1</u>
	<u>BRT to SVR</u>	<u>-4.0</u>	<u>-0.8</u>	<u>-6.9</u>	<u>AP1</u>
	<u>BRT to RF</u>	<u>-4.2</u>	<u>-4.8</u>	<u>-5.3</u>	<u>Mineral</u>
	<u>RF to SVR</u>	<u>1.1</u>	<u>-1.7</u>	<u>-6.5</u>	<u>Mineral</u>
	<u>BRT to SVR</u>	<u>-3.2</u>	<u>-6.5</u>	<u>-11.4</u>	<u>Mineral</u>
	<u>BRT to RF</u>	<u>-0.8</u>	<u>-1.2</u>	<u>-0.6</u>	<u>Organic</u>
	<u>RF to SVR</u>	<u>-4.1</u>	<u>-2.7</u>	<u>-5.4</u>	<u>Organic</u>
	<u>BRT to SVR</u>	<u>-5.2</u>	<u>-4.0</u>	<u>-6.0</u>	<u>Organic</u>
	<u>BRT to RF</u>	<u>-1.4</u>	<u>-3.3</u>	<u>-5.0</u>	<u>AP2</u>
	<u>RF to SVR</u>	<u>-3.2</u>	<u>-2.3</u>	<u>-6.4</u>	<u>AP2</u>
	<u>BRT to SVR</u>	<u>-4.5</u>	<u>-5.5</u>	<u>-11.0</u>	<u>AP2</u>
	<u>Algorithm</u>	<u>RMSE</u>	<u>MAE</u>	<u>MAPE</u>	<u>Approach</u>
<u>B</u>	<u>BRT</u>	<u>-33.1</u>	<u>-26.6</u>	<u>-28.7</u>	<u>AP1 to AP2</u>
	<u>RF</u>	<u>-34.6</u>	<u>-28.5</u>	<u>-29.0</u>	<u>AP1 to AP2</u>
	<u>SVR</u>	<u>-33.5</u>	<u>-30.1</u>	<u>-31.9</u>	<u>AP1 to AP2</u>
	<u>Algorithm</u>	<u>RMSE</u>	<u>MAE</u>	<u>MAPE</u>	<u>Approach</u>

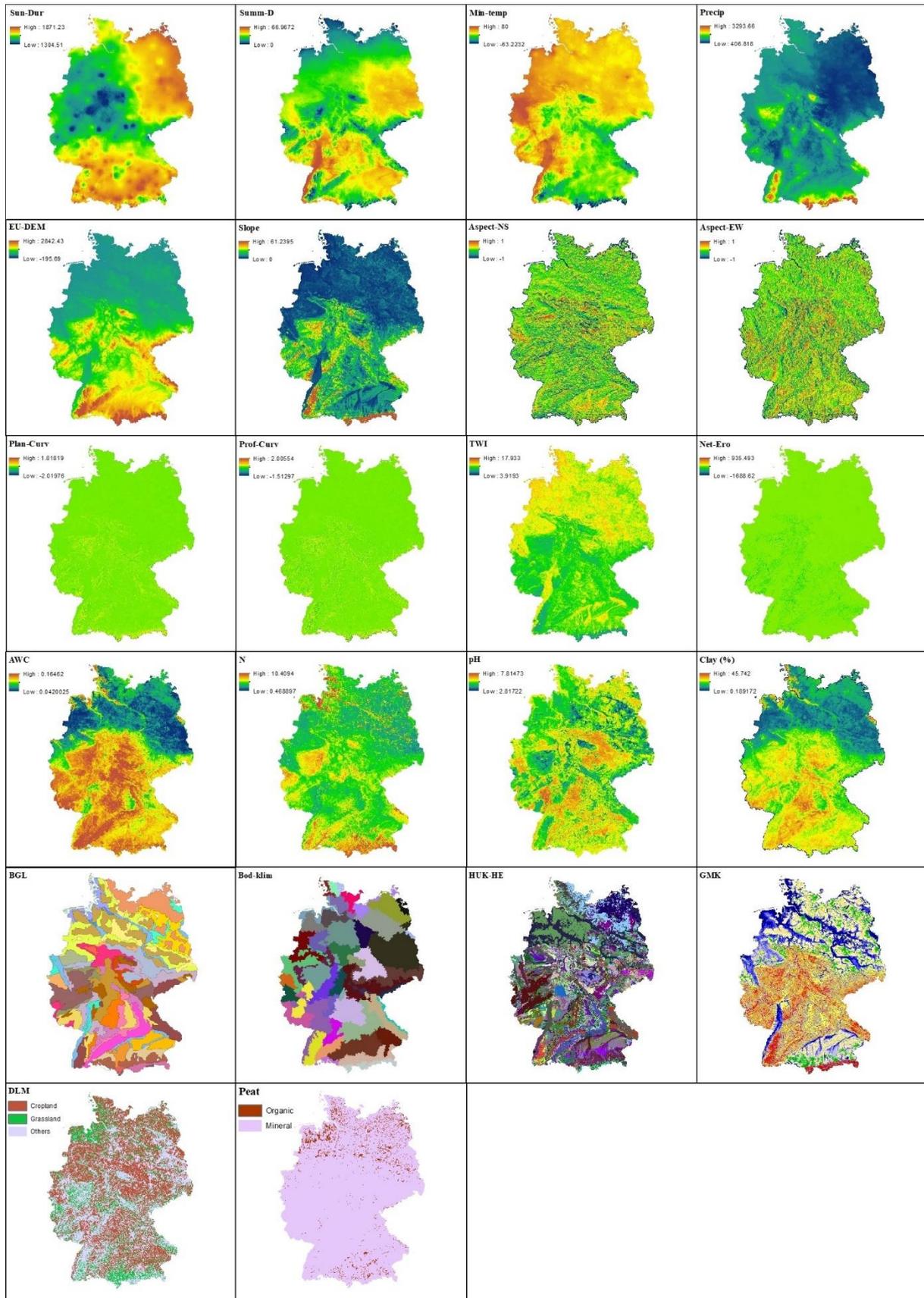
C	<u>BRT to RF</u>	<u>3.8</u>	<u>2.5</u>	<u>-1.3</u>	<u>AP1L</u>
	<u>RF to SVR</u>	<u>0.3</u>	<u>1.7</u>	<u>-0.9</u>	<u>AP1L</u>
	<u>BRT to SVR</u>	<u>4.2</u>	<u>4.2</u>	<u>-2.1</u>	<u>AP1L</u>
	<u>BRT to RF</u>	<u>-3.2</u>	<u>-3.2</u>	<u>-2.8</u>	<u>Mineral</u>
	<u>RF to SVR</u>	<u>0.0</u>	<u>-3.3</u>	<u>-8.4</u>	<u>Mineral</u>
	<u>BRT to SVR</u>	<u>-3.2</u>	<u>-6.5</u>	<u>-11.0</u>	<u>Mineral</u>
	<u>BRT to RF</u>	<u>-1.2</u>	<u>-2.5</u>	<u>-3.3</u>	<u>Organic</u>
	<u>RF to SVR</u>	<u>-3.7</u>	<u>-5.2</u>	<u>-7.6</u>	<u>Organic</u>
	<u>BRT to SVR</u>	<u>-5.2</u>	<u>-8.2</u>	<u>-10.7</u>	<u>Organic</u>
	<u>BRT to RF</u>	<u>-1.4</u>	<u>-3.3</u>	<u>-3.1</u>	<u>AP2L</u>
	<u>RF to SVR</u>	<u>-3.3</u>	<u>-3.4</u>	<u>-8.6</u>	<u>AP2L</u>
	<u>BRT to SVR</u>	<u>-4.6</u>	<u>-6.7</u>	<u>-11.4</u>	<u>AP2L</u>
	<u>Algorithm</u>	<u>RMSE</u>	<u>MAE</u>	<u>MAPE</u>	<u>Approach</u>
D	<u>BRT</u>	<u>-30.7</u>	<u>-23.7</u>	<u>-24.1</u>	<u>AP1L to AP2L</u>
	<u>RF</u>	<u>-34.2</u>	<u>-28.1</u>	<u>-25.4</u>	<u>AP1L to AP2L</u>
	<u>SVR</u>	<u>-36.5</u>	<u>-31.7</u>	<u>-31.3</u>	<u>AP1L to AP2L</u>
	<u>Algorithm</u>	<u>RMSE</u>	<u>MAE</u>	<u>MAPE</u>	<u>Approach</u>
E	<u>BRT</u>	<u>-4.9</u>	<u>-4.8</u>	<u>-6.9</u>	<u>AP1 to AP1L</u>
	<u>RF</u>	<u>-2.1</u>	<u>-1.6</u>	<u>-3.7</u>	<u>AP1 to AP1L</u>
	<u>SVR</u>	<u>3.2</u>	<u>0.0</u>	<u>-2.1</u>	<u>AP1 to AP1L</u>
	<u>BRT</u>	<u>-1.1</u>	<u>0.0</u>	<u>-0.8</u>	<u>Mineral</u>
	<u>RF</u>	<u>0.0</u>	<u>1.7</u>	<u>1.8</u>	<u>Mineral</u>
	<u>SVR</u>	<u>-1.1</u>	<u>0.0</u>	<u>-0.3</u>	<u>Mineral</u>
	<u>BRT</u>	<u>-1.5</u>	<u>-2.2</u>	<u>-7.2</u>	<u>Organic</u>
	<u>RF</u>	<u>-1.9</u>	<u>-3.5</u>	<u>-9.8</u>	<u>Organic</u>
	<u>SVR</u>	<u>-1.5</u>	<u>-6.0</u>	<u>-11.8</u>	<u>Organic</u>
	<u>BRT</u>	<u>-1.4</u>	<u>-1.1</u>	<u>-0.8</u>	<u>AP2 to AP2L</u>
	<u>RF</u>	<u>-1.4</u>	<u>-1.1</u>	<u>1.2</u>	<u>AP2 to AP2L</u>
	<u>SVR</u>	<u>-1.4</u>	<u>-2.3</u>	<u>-1.2</u>	<u>AP2 to AP2L</u>

26

27 **Table S43:** List of covariates, their abbreviations and ~~reference~~ their SCORPAN ID.

SCORPAN ID	Covariates	Abbreviation
S	Net erosion	Net-Ero
	Available water capacity	AWC
	Total nitrogen	TN
	pH	pH
	Soil organic map	Peat

	Clay content	Clay
C	Multi-annual grid of annual sunshine duration over Germany	Sun-Dur
	Multi-annual grids of number of summer days over Germany	Summ-D
	Multi-annual grids of monthly averaged daily minimum air temperature (2m) over Germany	Min-temp
	Multi-annual grids of precipitation height over Germany	Precip
O	Landuse	DLM
R	Digital elevation model	EU-DEM
	Slope	Slope
	Aspect north south direction	Aspect-NS
	Aspect east west direction	Aspect-EW
	Plan Curvature	Plan-Curv
	Profile curvature	Prof-Curv
	Topographic wetness index	TWI
	Geomorphographic map	GMK
P	Large-scale landscape unit map (Bodengrosslandschaft)	BGL
	Large-scale soil climate region map (Bodenklima)	Bod-klim
	Hydrological unit	HUK-HE
N	X coordination	x
	Y coordination	y

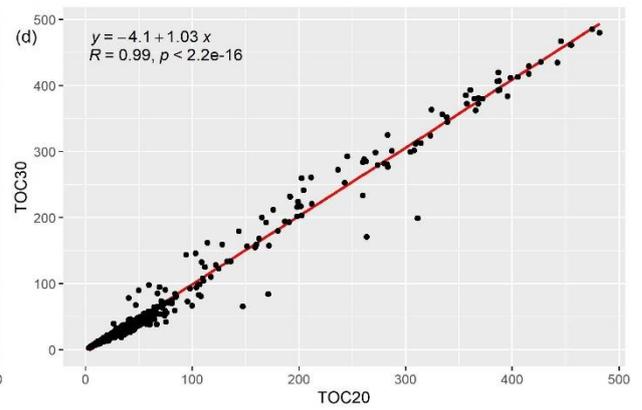
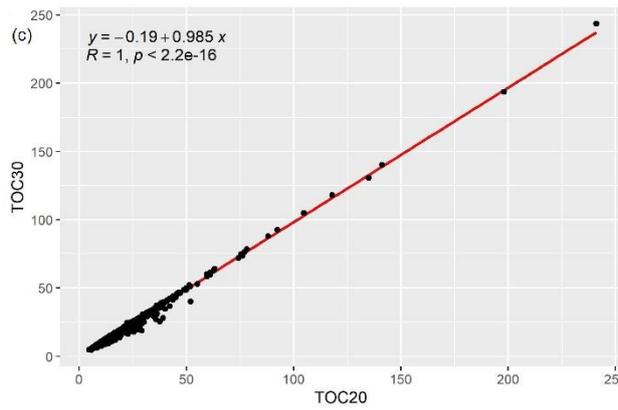
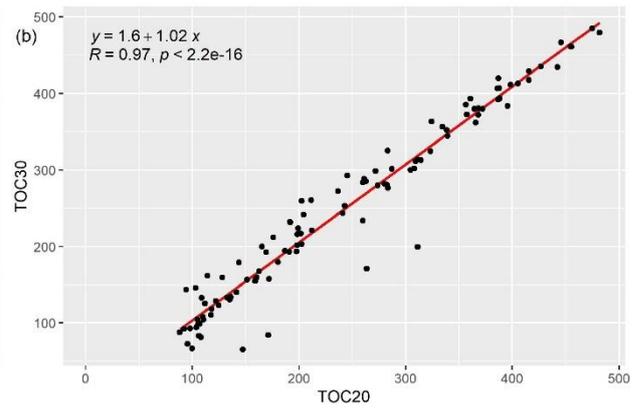
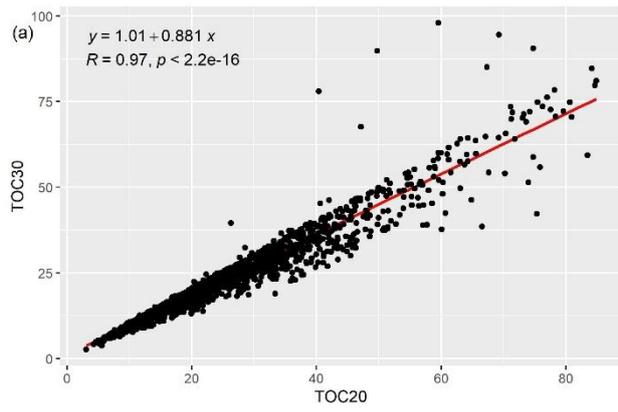


28

29 **Figure S1: Selected covariates: Sun-Dur** sunshine duration (DWD, 2017), **Summ-D** summer days (DWD, 2018b), **Min-**
 30 **temp** minimum temperature (DWD, 2018a), **Precip** precipitation (DWD, 2018c), **EU-DEM** digital elevation model
 31 (European Union Copernicus Land Monitoring Service, 2016), **Net-Ero** net soil erosion and deposition rates (Borrelli
 32 et al., 2018), **AWC** available water capacity (Ballabio et al., 2016), **N** total nitrogen (Ballabio et al., 2019), **pH** map of
 33 **pH** (Ballabio et al., 2019), **%Clay** % Clay (Ballabio et al., 2016), **BGL** soil scapes unit (BGR, 2008) [Legend], **Bod-**

34 Klim) soil-climate region (Roßberg et al., 2007), HUK-HE) hydrogeological unit of hydrogeological map-(BGR, SDG,
35 2019), GMK) geomorphographic map of Germany (BGR, 2007) [Legend], DLM) Land use (BKG, 2019), Peat) Organic
36 soils (Roßkopf et al., 2015).

37

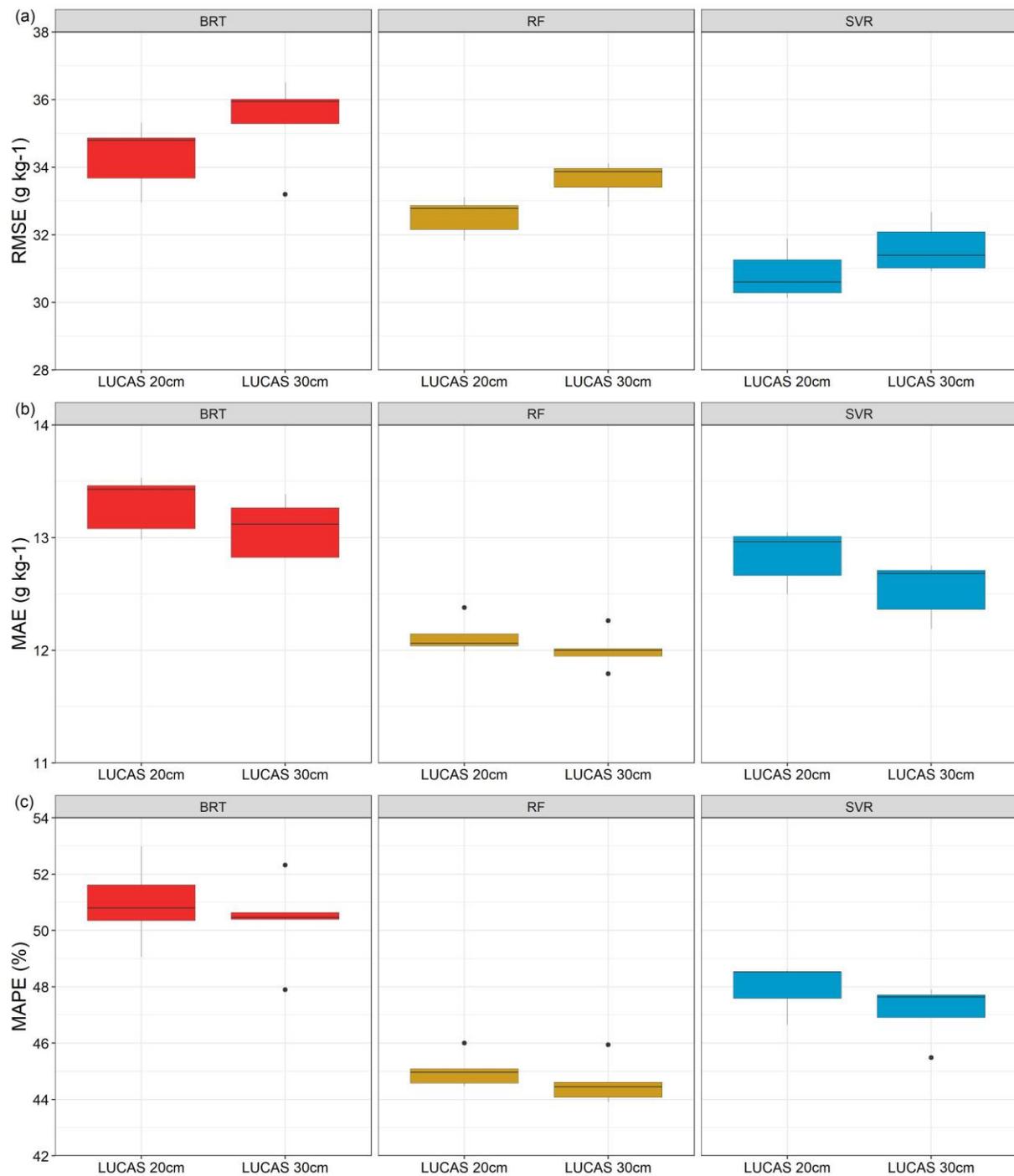


38

39

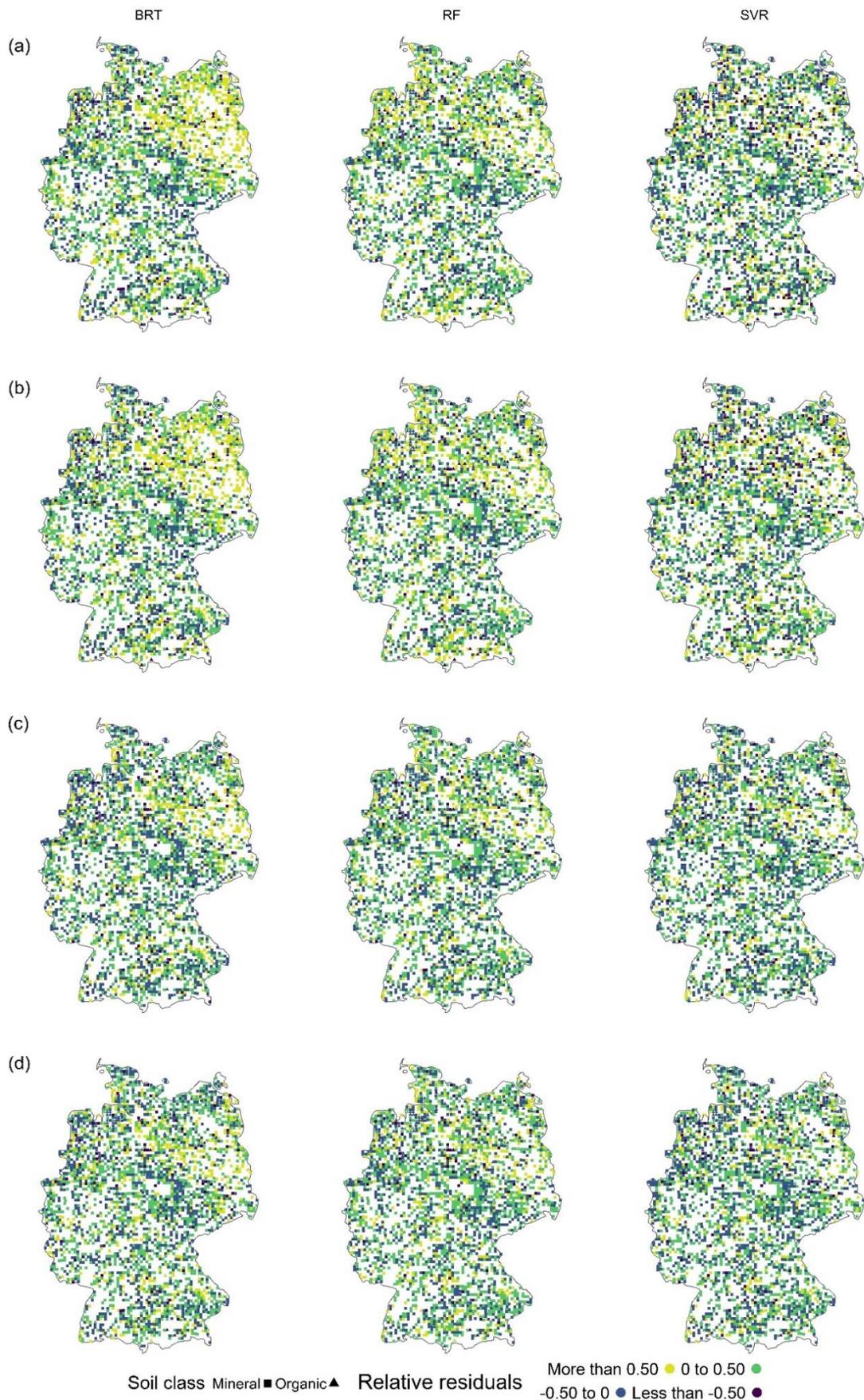
40

Figure S2: Regression plot for SOC depth extrapolation in A) Mineral soils, B) Organic soils, C) Cropland, D) Grassland.



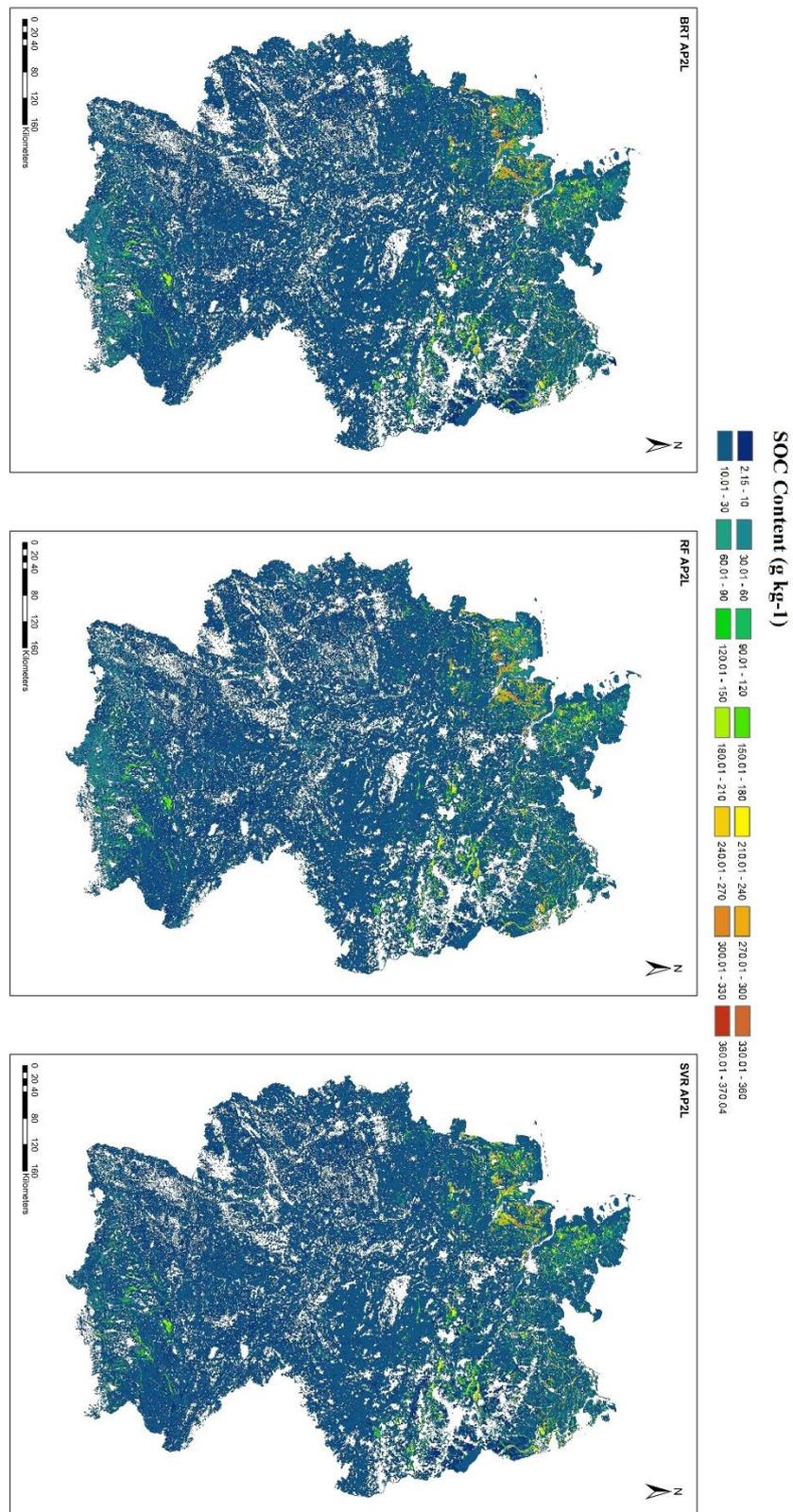
41

42 **Figure S31:** Boxplots comparing algorithm-model performance with regards to the three machine learning algorithms
 43 considering- LUCAS at the original sampling depth (20 cm) versus LUCAS with depth extrapolated (30 cm): A) RMSE
 44 (g kg⁻¹), B) MAE (g kg⁻¹) and C) MAPE (%). **BRT** = boosted regression trees, **RF** = random forest, and **SVR** = support
 45 vector regression.



46

47 **Figure S4: Spatial distribution of relative residuals from the models trained with the different machine learning**
 48 **algorithms. A) AP1 approach, B) AP1L approach, C) AP2 approach and D) AP2L approach. BRT = boosted**
 49 **regression trees, RF = random forest, and SVR = support vector regression.**



51

52 **Figure S5: Spatial prediction of SOC content (g kg⁻¹) of German agricultural soils based on the two-model approach**
 53 **for the three algorithms (BRT AP2L, RF AP2L, SVR AP2L). BRT = boosted regression trees, RF = random forest,**
 54 **and SVR = support vector regression.**

55 Disclaimer: It is important to note that the provided spatial prediction of SOC content must not be used to identify
56 the organic soils of Germany or to determine their spatial distribution. One reason is low sample size of organic
57 soils and the systematic underestimation of their SOC content, which leads to an underestimation of their spatial
58 extent. Furthermore, organic soils might have been mixed with mineral soil, i.e. due to deep ploughing, or feature
59 a mineral soil cover. Thus, organic soils might be present despite having a mineral topsoils. Therefore, this study
60 cannot nor intends to identify or classify organic soils.