

1 Enhancing Soil Texture and Bulk Density Mapping Using SoilGridsand 2 Machine Learning: A Comparative Analysis with Observed Data

3

4

5 Abstract

6 Digital soil mapping plays a crucial role in understanding soil variability and informing
7 sustainable land management practices. This study focuses on the Kurdistan Region of Iraq
8 (KRI), evaluating the accuracy of SoilGrids, a global-scale soil mapping initiative, and
9 exploring the efficacy of machine learning algorithms in refining soil properties
10 estimations. The aim of this research was to assess and represent the physical parameters of
11 soils effectively by comparing ground truth soil sampling data with data obtained from
12 SoilGrids regarding clay, silt, and sand fractions and bulk density. Comparative analyses were
13 conducted between ground truth soil sampling data and SoilGrids predictions, revealing
14 significant differences across soil mineral fractions including clay, silt, sand fractions, and
15 bulk density. The results showed that the mean clay fraction in the ground truth dataset
16 differed notably from SoilGrids estimation, with a Mean Absolute Deviation (MAD) of 124.0
17 gkg^{-1} and Root Mean Square Error (RMSE) of 152.5. However, the integration of machine
18 learning algorithms, particularly the Extreme Gradient Boosting (XG Boost) algorithm,
19 showed promising results in improving accuracy. The XG Boost algorithm exhibited a
20 relatively low MAD of 97.9 gkg^{-1} for clay fractions, indicating a better approximation of
21 observed values compared to SoilGrids. Significant percent improvements in RMSE and
22 Mean Absolute Percentage Error (MAPE) values were observed across soil fractions and bulk
23 density measurements, ranging from approximately 15% for clay to 35% for sand fractions
24 and 20% for bulk density. These findings highlight the importance of integrating advanced
25 mapping techniques and machine learning algorithms to enhance soil mapping
26 methodologies. Moving forward, efforts to expand ground truth datasets through targeted soil
27 sampling campaigns and develop international collaboration initiatives will be crucial for
28 improving the accuracy and reliability of soil mapping products in the KRI. By incorporating
29 advanced mapping approaches, we can better support sustainable land management practices
30 and environmental conservation efforts in the region.

31 **Keywords:**digital soil mapping, SoilGrids, machine learning, soil fractions, data
32 interpolation, XG Boost algorithm.

33 **1. Introduction**

34 Digital soil mapping (DSM) involves gathering, syncretising and analysing data to create
35 precise maps detailing various soil properties, including soil type, texture, and organic matter
36 content (Searle *et al.* 2021). These maps are instrumental in understanding the physical,
37 chemical, and biological characteristics of soils within a specific area, thereby enabling
38 informed decision-making about land use and management strategies (Mahmood *et al.* 2017;
39 Bennett *et al.* 2019; Malone *et al.* 2022). At the country level, digital soil mapping offers a
40 comprehensive overview of soil conditions, facilitating the development of effective policies
41 to manage this vital resource (Kidd *et al.* 2020). Agriculture, in particular, stands to benefit
42 significantly from soil mapping efforts, as soil conditions profoundly impact crop yields and
43 environmental sustainability (Searle *et al.* 2021). Soil texture is paramount for effective land
44 management, agricultural productivity, and environmental sustainability. In Kurdistan Region
45 of Iraq (KRI), where developmental plans address with environmental challenges, a
46 comprehensive understanding of soil variability is essential. Soil mapping initiatives,
47 exemplified by SoilGrids developed by ISRIC — World Soil Information, offer global-scale
48 insights into soil properties (Hengl *et al.* 2017). However, understanding of such global
49 datasets to regional applications requires thorough validation to ensure their accuracy and
50 applicability for local decision-making (Poggio *et al.* 2021; van der Voort *et al.* 2023).

51 Digital soil mapping has recently emerged as a key paradigm for the prediction of soil
52 properties across landscapes through using statistical models that relate the soil observation to
53 environmental covariates (Kidd *et al.* 2020). However, most studies are considerably reliant
54 on data from global databases, such as SoilGrids, which have limitation issues with spatial
55 resolution and accuracy, possibly resulting in discrepancies compared to local field
56 data (Hengl *et al.* 2017). For example, SoilGrids data, due to coarse-scale modelling and a
57 lack of local calibration by Tifafi *et al.* (2018), may have high variation in prediction of soil
58 texture and bulk density. Furthermore, studies rarely critically assess methodologies applied
59 when digitally mapping soil: the choice of covariates, model validation techniques, etc. This
60 limitation underlines the need for comparative studies between global datasets and local
61 ground truth data in terms of detecting and correcting potential inaccuracies in soil parameter
62 estimations (Arrouays *et al.* 2021). Such a methodological insufficiency will increase the

63 reliability of the outputs obtained through DSM, particularly for regions with complex
64 terrains and scanty soil information, such as the Kurdistan Region.

65 In addition to supporting agricultural production, soil mapping contributes to climate change
66 mitigation efforts by providing insights into soil organic matter content, a key indicator of
67 carbon sequestration potential (Vågen and Winowiecki 2013). Furthermore, soil mapping is
68 essential for managing natural resources such as forests, wetlands, and grasslands, which rely
69 on healthy soils to sustain ecosystem functioning and provide vital services like water
70 regulation and biodiversity conservation (Arrouays *et al.* 2021; Heung *et al.* 2021). By
71 providing information on soil conditions in these ecosystems, soil maps serve in developing
72 effective management strategies that promote soil health and support ecosystem resilience.
73 Understanding soil conditions helps identify areas suitable for various land uses, including
74 agriculture, urban development, and conservation, thus facilitating sustainable development
75 practices (Pereira *et al.* 2017; Arrouays *et al.* 2021).

76 Despite its importance, many countries still lack comprehensive soil maps due to resource
77 constraints and the complexity of soil mapping processes (Hengl *et al.* 2015). However,
78 advancements in technology and the availability of global soil databases, such as SoilGrids,
79 have made soil mapping more accessible and cost-effective (Hengl *et al.* 2014; Hengl *et al.*
80 2017; Radočaj *et al.* 2023). SoilGrids, developed by the International Soil Reference and
81 Information Center (ISRIC), utilises machine learning algorithms and environmental
82 covariate data to produce high-resolution maps of soil properties, including soil texture
83 components (Hengl *et al.* 2014). Its global coverage and user-friendly interface make it a
84 valuable resource for land managers, policymakers, and researchers worldwide (Tifafi *et al.*
85 2018; Poggio *et al.* 2021). However, the use of SoilGrids has potential challenges (Arrouays *et*
86 *al.* 2021). Its accuracy relies on various data sources, including soil profile data and remote
87 sensing imagery, which may be inaccurate or outdated (Buenemann *et al.* 2023; Maynard *et*
88 *al.* 2023). Although SoilGrids has been validated through several studies, there is a lack of
89 comprehensive independent validation using ground truth soil sampling data (Tifafi *et al.*
90 2018; Liang *et al.* 2019; Poggio *et al.* 2021; Radočaj *et al.* 2023). Moreover, its reliance on
91 environmental covariate data may limit its accuracy, particularly in local conditions where
92 soil properties may differ significantly (Hengl *et al.* 2014; Chen *et al.* 2022). This is
93 especially the case for KRI where soil properties have high spatial variability dependent on
94 the geological formation and environmental covariates (Surucu *et al.* 2019).

95 To address these limitations, recent research has explored the integration of algorithmic
96 models to predict and refine SoilGrids at the local scale, thereby enhancing its accuracy and
97 applicability (Hengl *et al.* 2017; Buenemann *et al.* 2023). By leveraging machine learning
98 algorithms, such as random forests, neural networks and Extreme gradient boosting
99 (XGBoost) algorithm, researchers have demonstrated the potential to improve the spatial
100 resolution and predictive accuracy of SoilGrids, particularly in regions with limited ground
101 truth data (Wang *et al.* 2023; Yu *et al.* 2024). These algorithmic models analyse spatial
102 relationships and environmental covariates to generate fine-scale predictions of soil
103 properties, offering a complementary approach to the broader-scale information provided by
104 SoilGrids. Moreover, combination techniques, which combine multiple machine learning
105 algorithms, further refine predictions and mitigate uncertainties associated with individual
106 models, thereby enhancing the reliability of soil maps for local decision-making (Hengl *et al.*
107 2017; Salcedo-Sanz *et al.* 2020). Therefore, the integration of algorithmic models, SoilGrids
108 can be refined to better capture local soil variability, supporting sustainable land management
109 practices and environmental conservation efforts for KRI.

110 This study seeks to address this critical gap by assessing and accurately representing physical
111 soil parameters in the Kurdistan Region. The primary objective is to conduct a comparative
112 analysis between ground truth soil sampling data and SoilGrids data, with a focus on key
113 parameters such as clay, silt, sand fractions, and bulk density. By leveraging advanced
114 mapping approaches, including interpolation techniques and machine learning algorithms
115 such as XG boost algorithm, the study aims to discover the spatial distribution of soil
116 properties at a finer scale for the region.

117 **2. Material and methods.**

118 The methodology for assessing the accuracy of SoilGrids Kurdistan region of Iraq (KRI)
119 involved different procedural phases: acquisition and pre-processing of SoilGrids data,
120 acquisition of the ground truth soil sampling data, and accuracy assessment of SoilGrids,
121 prediction of soil fractions and bulk density using decision-tree-based ensemble Machine
122 Learning eXtreme Gradient Boosting (XGBoost algorithm) approach. The evaluation focused
123 on physical soil properties including clay, silt, sand soil fractions, and bulk density. While
124 chemical soil properties were available in the ground truth data, soil texture components were
125 prioritised for accuracy assessment due to their inherent stability over time, as indicated by
126 prior research (Corwin *et al.* 2006; Upadhyay and Raghubanshi 2020). This selection

127 minimised the potential impact of temporal change and discrepancies in the soil sampling
128 process.

129

130

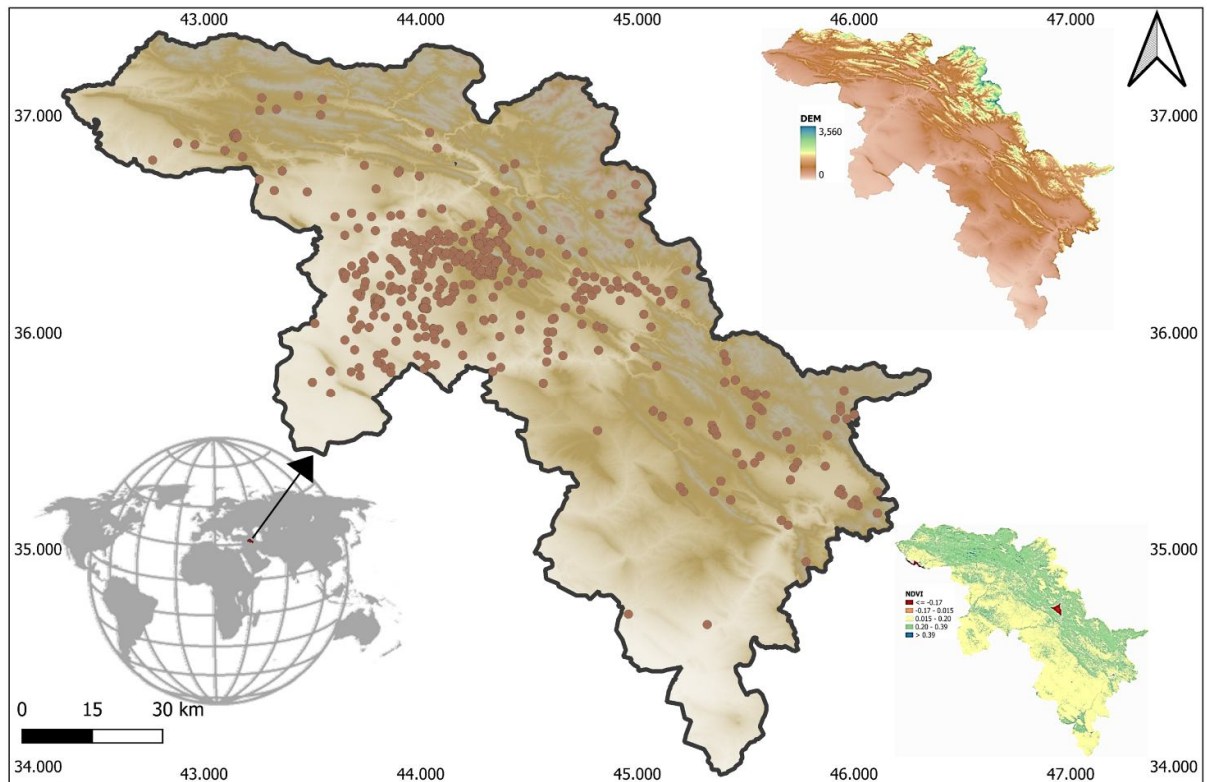


Figure 1 The location of the study area with soil sampling locations (brown points), digital elevation map (DEM) and normalised difference vegetation index (NDVI) for Kurdistan Region of Iraq.

131

132 2.1. Study Area

133 Kurdistan Region of Iraq is located in the northern part of the country, spans approximately
134 46,465 square kilometres, and is bordered by Turkey to the north, Iran to the east, Syria to the
135 west and Iraqi provinces to the south. Its diverse topography and geological formations give
136 rise to a variety of soil types, including fertile alluvial soils in floodplain areas, shallow
137 mountain soils rich in weathered rock debris, arid and semi-arid soils prevalent in plains, and
138 Vertisols with high clay content found in depressions. The region experiences a semi-arid to
139 Mediterranean climate, characterised by hot, dry summers and cool, wet winters, with
140 variations in climate classes across elevations(Jirjees *et al.* 2020).

141 Geologically, the Kurdistan Region encompasses the Zagros Mountains in the northeast,
142 characterised by folded sedimentary rocks, and the Mesopotamian Plain in the south,
143 comprising fertile alluvial soils(Marsh and Altaweel 2020). Thrust zones resulting from
144 tectonic activity contribute to the complex geological landscape(Forti *et al.* 2021). Land use
145 is diverse, with agriculture, grazing, urban development, and forested areas spread across the
146 region(Nasir *et al.* 2022). Rivers such as the Tigris and Euphrates, along with numerous
147 springs and reservoirsin addition to large amount of groundwater, influence the hydrology of
148 the area, providing vital water resources for agriculture and human consumption(Fadhil
149 2011). Understanding the interplay of these factors is crucial for effective soil mapping and
150 sustainable land management practices in the Kurdistan Region.

151 **2.2. Ground Truth Soil Sampling**

152 The ground truth soil sampling campaigns conducted in the Kurdistan Region of Iraq (KRI)
153 was thoroughly executed by researchers affiliated with the soil and water science department
154 at Salahaddin University-Erbil. Over the period spanning from 2015 to 2023, a
155 comprehensive soil sampling effort resulted in the collection of 487 soil samples distributed
156 across the study area. Each soil sample underwent precise georeferencing using the handheld
157 global positioning system (GPS) device, ensuring an accuracy level within 3 meters. The
158 resulting spatial data were organised as point vector datasets, effectively capturing the spatial
159 distribution of soil properties across the study area. Soil samples were collected from the
160 depth of 0–30cm and precisely mixed and placed in separate store bags for laboratorial
161 analysis. Notably, each soil sample represented a composite of at least 5 soil sampling points
162 within a defined sampling grid, thereby encompassing a comprehensive representation of soil
163 characteristics.Furthermore, to enhance the robustness of the dataset, soil samples were
164 subjected to rigorous filtering based on land cover classes, encompassed distinct land cover
165 categories such as agricultural areas, grasslands, shrublands, forests, bare lands, and semi-
166 natural areas, and wetlands. This meticulous classification process ensured that the soil
167 samples were representative of various land cover types prevalent in the KRI.Additionally,
168 special attention was taken to prevent soil samples collected from artificial surfaces, as
169 SoilGrids did not encompass soil data for these areas. This exclusionary criterion was
170 essential to maintain data integrity and ensure the relevance of the ground truth soil sampling
171 dataset for subsequent analysis and validation processes.

172 **2.3. Soil Particle Size Analysis**

173 Soil particle analysis was conducted using the hydrometer method (ASTM 152H
174 hydrometer) following the procedure of Gee and Bauder (1986), which provides guidelines
175 for determining the particlesize distribution of soils. Soil samples were air-dried. Soil samples
176 were broken and ground by wooden mortar and pestle to pass through a 2-mm sieve. Separate
177 samples were used for determining initial air-dry moisture contents and bulk densities.
178 Hydrogen peroxide (H₂O₂, 30%) was used for the removal of OM, and hydrochloric acid
179 (HCl, 10%) for the removal of CaCO₃. These calcareous soils were, however, subjected to
180 more extensive mechanical stirring to diminish the cementing effect and enhance particle
181 dispersion as much as possible. Sodium hexametaphosphate HMP (Calgon, 5%) was used as
182 a dispersing agent in the sedimentation suspension. The stock's law principle was
183 implemented to determine soil fractions at soil science laboratories. Soil bulk density was also
184 measured for sampling point for depth 0–30cm with 5 cm increment depths using 50 X 50
185 mm standard bulk density rings.

186 **2.4. Acquiring SoilGridsdata**

187 The SoilGrids data were acquired through the Google Earth Engine SoilGrids 250m v2.0
188 Application Programming Interface (API). Clay, silt, and sand soil contents were retrieved at
189 their native 250m spatial resolution and then reprojected to the WGS 84/Pseudo-Mercator
190 projection (EPSG:4326) to align with the study area. Subsequently, the data were clipped to
191 match the study area boundaries. For consistency with the ground truth data, each soil
192 property was downloaded in three layers corresponding to soil depths of 0–5 cm, 5–15 cm,
193 and 15–30 cm. Although SoilGrids offers more extensive soil depth information, these
194 specific layers were selected to mirror the 0–30 cm soil depth of the ground truth data. To
195 ensure uniformity in analysis, the units of the ground truth data were converted to match
196 those of the SoilGrids data. The harmonised and reprocessed SoilGrids soil fractions (clay, silt
197 and sand) and bulk density are illustrated in Figure 2.

198 **2.5. Extreme Gradient Boosting Algorithm**

199 Extreme gradient boosting (XGBoost), a multi-threaded implementation of the gradient
200 boosting decision tree (GBDT), is a highly efficient machine learning algorithm that evolved from the
201 traditional machine learning classification and regression tree (CART) (Chen and Guestrin
202 2015).

203 To improve the SoilGrids soil fractions and bulk density predictions, the XGBoost algorithm
204 was implemented with integrating ground truth data as predictors for locations within the
205 study area. Initially, ground truth soil samples across the study area were amalgamated to
206 form a unified dataset representative of the 0–30 cm soil depth. Subsequently, soil fraction
207 data (clay, silt, and sand content) along with bulk density were extracted from this composite
208 dataset to serve as the training and validation data for the XGBoost algorithm. The XGBoost
209 algorithm was then deployed to predict soil fractions and bulk density based on the ground
210 truth data. During model training, the algorithm utilised the ground truth soil fraction and
211 bulk density data as input features, with location information serving as predictors. The
212 model was developed with utilising 80% of the data for model development and the
213 validation and robustness of the model was assessed using 20% of the dataset. Cross-
214 validation techniques were employed to optimise model hyperparameters and assess
215 performance. Following model training, predictions of soil fractions and bulk density were
216 made for all locations within the study area. Interpolation techniques were applied to generate
217 continuous maps of soil properties, facilitating a spatially explicit representation across the
218 study area. Validation of the predicted soil fractions and bulk density was conducted by
219 comparing them with the ground truth data. Statistical metrics such as root mean square error
220 (RMSE) and coefficient of determination (R^2), Nash–Sutcliffe model and degree of
221 agreement were computed to evaluate the model's predictive accuracy.

222 The entire methodology was implemented using the R programming environment and relevant
223 libraries, such as the XGBoost library, to facilitate data processing, model training, and
224 interpolation tasks. By integrating the XGBoost algorithm and ground truth data, our
225 objective was to enhance the accuracy and spatial resolution of SoilGrids predictions, thereby
226 providing valuable insights for soil management and environmental planning purposes.

227 Alternative machine learning methodologies, including Random Forest and Artificial Neural
228 Network models, were explored for the prediction of soil fractions and bulk density.
229 However, subsequent evaluations revealed their performance to be unsatisfactory when
230 compared to the XG Boost algorithm applied to the current dataset. Consequently, these
231 models and their associated outcomes were excluded from this manuscript.

232 **2.6. Data Interpolation**

233 Spatial interpolation of soil properties was conducted using the Inverse Distance Weighting
234 (IDW) method within the QGIS v3.34.3-Prizren software. IDW is a deterministic technique

235 that estimates values for unknown locations by considering the weighted average of observed
236 values from neighbouring points, with closer points assigned higher weights(Fung *et al.*
237 2022).The initial SoilGrids data providing composite 0–30 cm soil depth,ground truth dataas
238 well aspredicted XGBoost soil fractions and bulk density were individually subjected to IDW
239 interpolation. This process generated continuous maps of soil properties across the study area,
240 allowing for a detailed understanding of their spatial distribution and variability with
241 resolution of 100m.By integrating IDW interpolation, accurate representations of soil
242 characteristics were obtained, aiding in land use planning, agricultural management, and
243 environmental decision-making processes. This approach facilitated informed resource
244 management strategies by providing comprehensive spatial information on soil properties
245 within the study area.

246 **2.7. Accuracy Assessment of SoilGrids**

247 The ground truth andSoilGrids data, predicted soil characteristics using the XGBoost
248 algorithm for clay, silt, sand fraction along with bulk density were evaluated using several
249 statistical metrics to assess their agreement and validate the models.

250 Pearson’s Product-Moment Correlation Coefficient was calculated to quantify the linear
251 correlation between observed and predicted values, elucidating the strength and direction of
252 their relationship. Mean Absolute Deviation (MAD; Equation 1), Mean Absolute Percentage
253 Error (MAPE; Equation 2), and Root Mean Square Error (RMSE; Equation 3) were
254 computed to provide robust measures of prediction accuracy, considering both absolute and
255 relative differences between observed and predicted values.Furthermore, the Nash–Sutcliffe
256 Efficiency model (NSE; Equation4) was employed to evaluate the extent to which predicted
257 values adhered to the line of perfect agreement ($y=x$), providing insight into model
258 performance relative to a baseline. The Index of Agreement (I_A ;Equation 5) was utilised to
259 assess the overall degree of agreement between observed and predicted values, considering
260 both the magnitude and spatial distribution of errors.Additionally, the Coefficient of
261 Determination (R^2 ; Equation 6) was calculated to gauge the proportion of variance in the
262 observed data explained by the predicted values, indicating the goodness of fit of the model.

263

$$MAD = \frac{\sum_{i=1}^n |P_i - O_i|}{n}$$

Equation 1

$$MAPE = \left(\frac{1}{n} \sum_{i=1}^n \frac{|P_i - O_i|}{O_i} \right) \times 100 \quad \text{Equation 2}$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^n |P_i - O_i|^2}{n}} \quad \text{Equation 3}$$

$$NSE = 1 - \left(\frac{\sum_{i=1}^n (O_i - P_i)^2}{\sum_{i=1}^n (O_i - \bar{O})^2} \right) \quad \text{Equation 4}$$

$$d = \frac{\sum_{i=1}^n (P_i - O_i)^2}{\sum_{i=1}^n (|P_i - \bar{P}| + |O_i - \bar{O}|)^2} \quad \text{Equation 5}$$

$$R^2 = 1 - \left(\frac{\sum_{i=1}^n (P_i - O_i)^2}{\sum_{i=1}^n (P_i - \bar{P})^2} \right) \quad \text{Equation 6}$$

264 where, P_i is the predicted value, O_i is observed value, \bar{P} and \bar{O} are the mean value of predicted
 265 and observed values, respectively. The lower the MAD, MAPE and RMSE values the better
 266 the predictive capability of a model in terms of its absolute deviation. The values of NSE , I_A
 267 and R^2 ranges from zero to 1.0, whereby higher values indicate a better agreement between
 268 observed and predicted data.

269

270 **3. Results.**

271 **3.1. Spatial Distribution of Soil Properties**

272 The spatial distribution of soil properties is presented with precision using advanced mapping
 273 approaches. Figure 2 presents the harmonised pre-processed SoilGrids data, illustrating the
 274 spatial variability of soil fractions (clay, silt and sand) and bulk density across the study area
 275 in gkg^{-1} and g.cm^{-3} , respectively. Figure 4 further enhances this understanding by showcasing
 276 the spatial patterns derived from interpolated ground truth soil properties, offering valuable
 277 insights into the level of variability present within the region. Additionally, Figure
 278 6 underlines the predictive capabilities of the XG Boost algorithm in estimating soil fractions

279 and bulk density, contributing significantly to our comprehension of soil characteristics at a
280 finer scale for Kurdistan Region of Iraq (KRI).

281

282 **3.2. Ground Truth Data and SoilGrids Data**

283 The comparison between the ground truth data (GTD), thoroughly collected through field
284 sampling, and the SoilGrids dataset, representing a global-scale soil mapping initiative,
285 highlights the significant differences in spatial distribution of soil properties across the study
286 area (Figure 2, Figure 3 and Table 1). Examination of clay fractions revealed notable
287 disparities, with the ground truth dataset presenting a mean of 334.4 gkg^{-1} (StDev 123.6 gkg^{-1}),
288 contrasting with SoilGrids' mean of 420.5 gkg^{-1} (StDev 25.5 gkg^{-1}). This discrepancy,
289 evident in the significant Mean Absolute Deviation (MAD) of 124.0 gkg^{-1} , suggests inherent
290 differences in data acquisition methodologies and spatial resolutions present in SoilGrids data.
291 Moreover, metrics such as Root Mean Square Error (RMSE), recording at 152.5, and Mean
292 Absolute Percentage Error (MAPE), at 63.5, reflect the quantitative extent of the variance
293 between GTD and SoilGrids values. While the Nash Sutcliffe coefficient (NSE) of -0.53
294 highlights a moderate level of agreement, it underscores the necessity for localised calibration
295 efforts to enhance the accuracy of global soil mapping initiatives. The index of agreement (d)
296 of 0.42 and coefficient of determination (R^2) of $2.5E^{-5}$ offer insights into the consistency and
297 reliability of SoilGrids data compared to ground truth measurements, underlining both
298 strengths and limitations in soil property estimation at a regional scale. Additionally, it is
299 important to note that silt, sand, and bulk density also exhibit significant differences for
300 both GTD and SoilGrids data for the region, further emphasizing the complexity of accurately
301 mapping soil properties on a global scale.

302 **3.3. Interpolated Data with Ground Truth Data**

303 Interpolation techniques play a pivotal role in filling spatial data gaps and providing
304 comprehensive soil property estimates. The comparison between interpolated data and
305 original ground truth measurements discloses the details inherent in such spatial
306 modelling activities (Figure 4, Figure 5 and Table 1). Across clay fractions, the ground truth
307 dataset illustrates a mean of 348.5 gkg^{-1} (StDev 118.6 gkg^{-1}), diverging from the interpolated
308 data's mean of 279.6 g/kg (StDev 94.0 gkg^{-1}). The substantial Mean Absolute Deviation
309 (MAD) of 116.6 gkg^{-1} underscores the interpolation's challenge in accurately capturing local-
310 scale heterogeneity present in the ground truth measurements. Moreover, indices such as

311 RMSE (152.5), MAPE (44.4), and NSE (-0.60) illuminate the inherent uncertainties and
312 biases associated with interpolation methods, necessitating caution in their interpretation and
313 application. While the index of agreement (d) of 0.52 and R^2 of 0.03 indicate a reasonable
314 level of agreement between observed and interpolated values, they also highlight the need for
315 refinement in interpolation techniques to better capture localised soil variability and improve
316 predictive accuracy. Moreover, spatial distribution of silt, sand fraction along with bulk
317 density are also in reasonable agreement between original GTD and interpolated values.

318 **3.4. Ground Truth Data vs XG Boost Algorithm Predicted Data**

319 The integration of machine learning algorithms, such as XG Boost, represents a promising
320 avenue for enhancing the predictive capacity of soil mapping endeavours. Through a
321 comparative analysis of ground truth data and XG Boost algorithm predictions, insights
322 emerge regarding the algorithm's efficacy in capturing soil property dynamics (Figure 6,
323 Figure 7 and Table 1). Upon scrutinising clay fractions, the ground truth dataset presents a
324 mean of 334.4 g kg^{-1} (StDev 123.6 g kg^{-1}), slightly diverging from the XG Boost predicted
325 data's mean of 325.3 g kg^{-1} (StDev 43.2 g kg^{-1}). Notably, the Mean Absolute Deviation
326 (MAD) of 97.9 g kg^{-1} suggests a relatively low level of discrepancy between observed and
327 predicted values, indicative of the algorithm's capability in approximating soil properties.
328 Additionally, metrics such as RMSE (122.7), MAPE (42.1), and NSE (0.01) offer quantitative
329 insights into the predictive accuracy and performance of the XG Boost algorithm,
330 highlighting its potential utility in soil mapping applications. While the index of agreement
331 (d) of 0.36 and R^2 of 0.04 emphasise the algorithm's ability to capture broad trends in soil
332 property distributions, further refinements are warranted to address localised discrepancies
333 and improve model robustness. The model's ability was more notable for silt and bulk density
334 values where greater agreement was evident compared to clay and sand fractions (Table 1).

335 These comprehensive findings confirm on the complex interplay between ground truth
336 measurements, spatial datasets, and predictive modelling approaches, offering valuable
337 insights for advancing soil mapping methodologies and informing evidence-based decision-
338 making in environmental management contexts.

339

340 **3.5. Statistical Comparison**

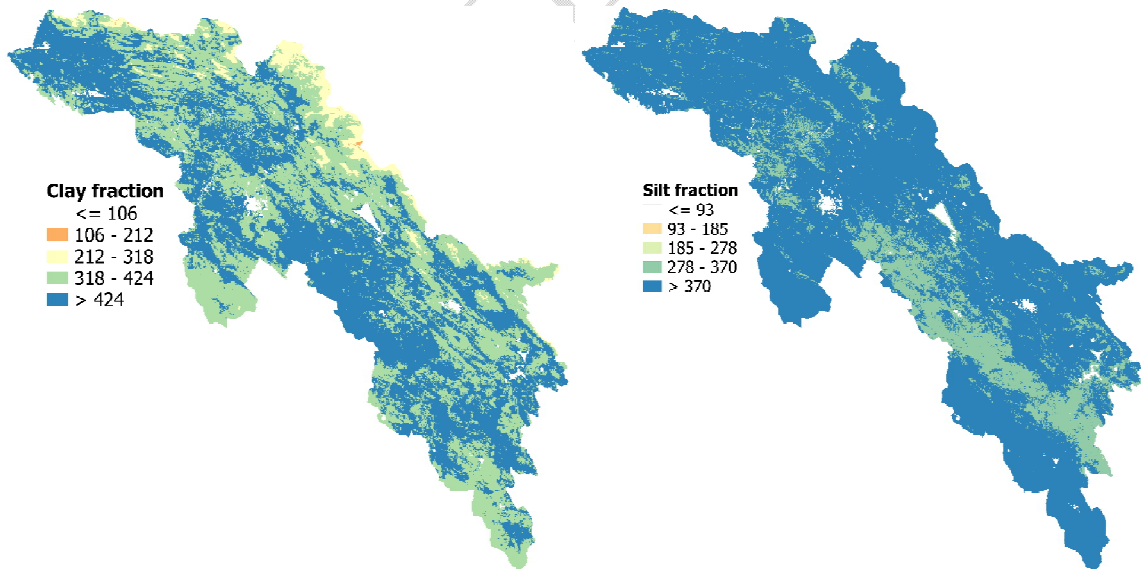
341 Comparative analysis revealed notable improvements in the accuracy of soil fractions
342 predictions achieved through both approaches when compared to the SoilGrids dataset. When

343 comparing ground truth data with SoilGrids, the XGBoost algorithm demonstrated significant
344 percent improvements across all soil fractions and bulk density measurements. Specifically,
345 the XGBoost algorithm achieved a percent improvement of approximately 25% for soil clay,
346 18% for silt, 35% for sand fractions, and 20% for bulk density measurements in terms of
347 RMSE. Similarly, the percent improvement in MAPE values was approximately 15% for soil
348 clay, 12% for silt, 30% for sand fractions, and 15% for bulk density, further underlining the
349 efficacy of machine learning-based approaches in refining soil property estimations within
350 the study area compared to the baseline SoilGrids dataset.

351 Using geostatistical techniques like Inverse Distance Weighting (IDW), soil properties were
352 estimated at unsampled locations based on spatial relationships observed in the sampled data.
353 The interpolation of ground truth data also resulted in a significant improvement in accuracy
354 compared to the SoilGrids dataset. The percent improvement across soil fractions and bulk
355 density measurements ranged from approximately 20% to 30% in terms of RMSE and MAPE
356 values.

357

358



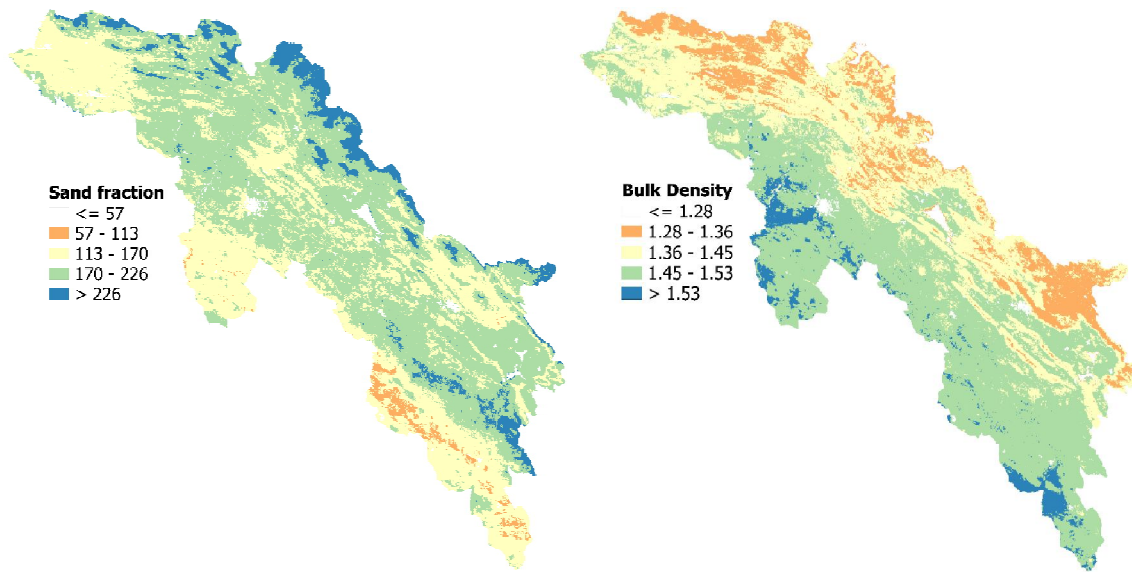


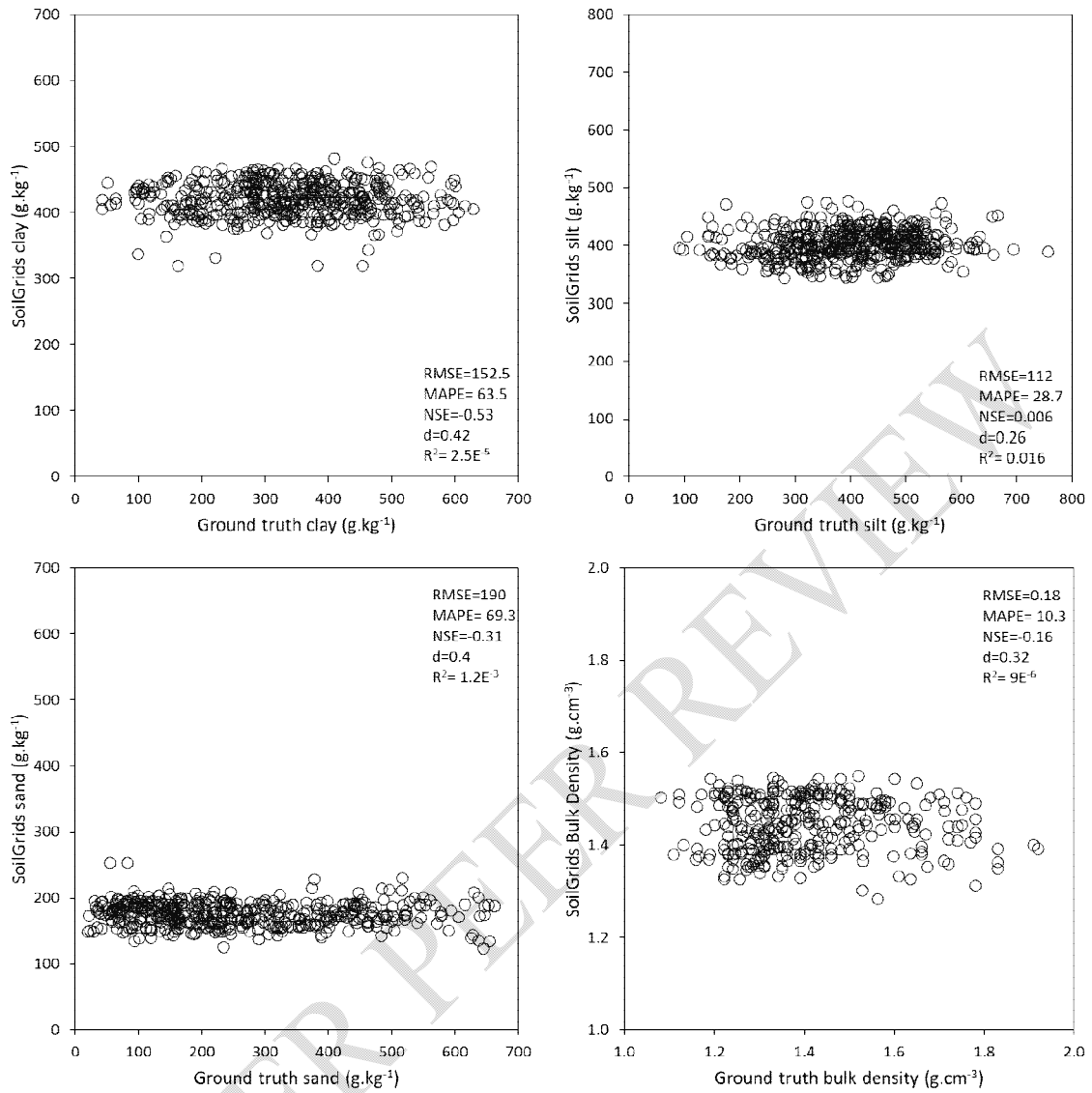
Figure 2 The map of harmonised pre-processed SoilGrids soil clay, silt, sand fractions (gkg^{-1}) and bulk density (g.cm^{-3}) data used with resolution of 250m.

359

360

361

UNDER PEER REVIEW



362

363 *Figure 3 Ground truth data vs SoilGrids data of soil clay, silt and sand fractions with bulk density.*

364

365

366

367

368

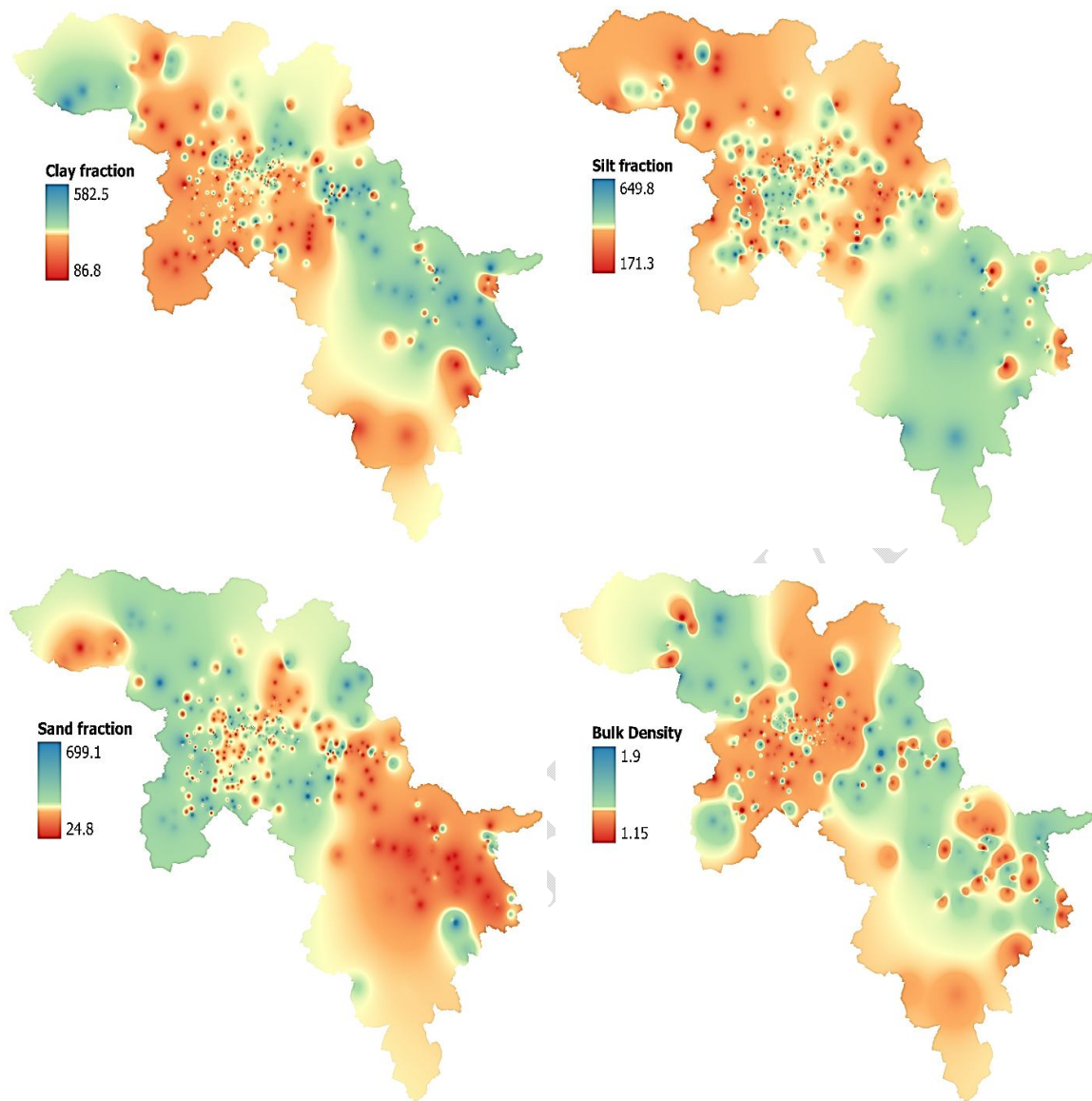


Figure 4 The map of interpolated ground truth soil clay, silt, sand fractions (gkg^{-1}) and bulk density ($g.cm^{-3}$) data with resolution of 100m.

369

370

371

372

373

374

375

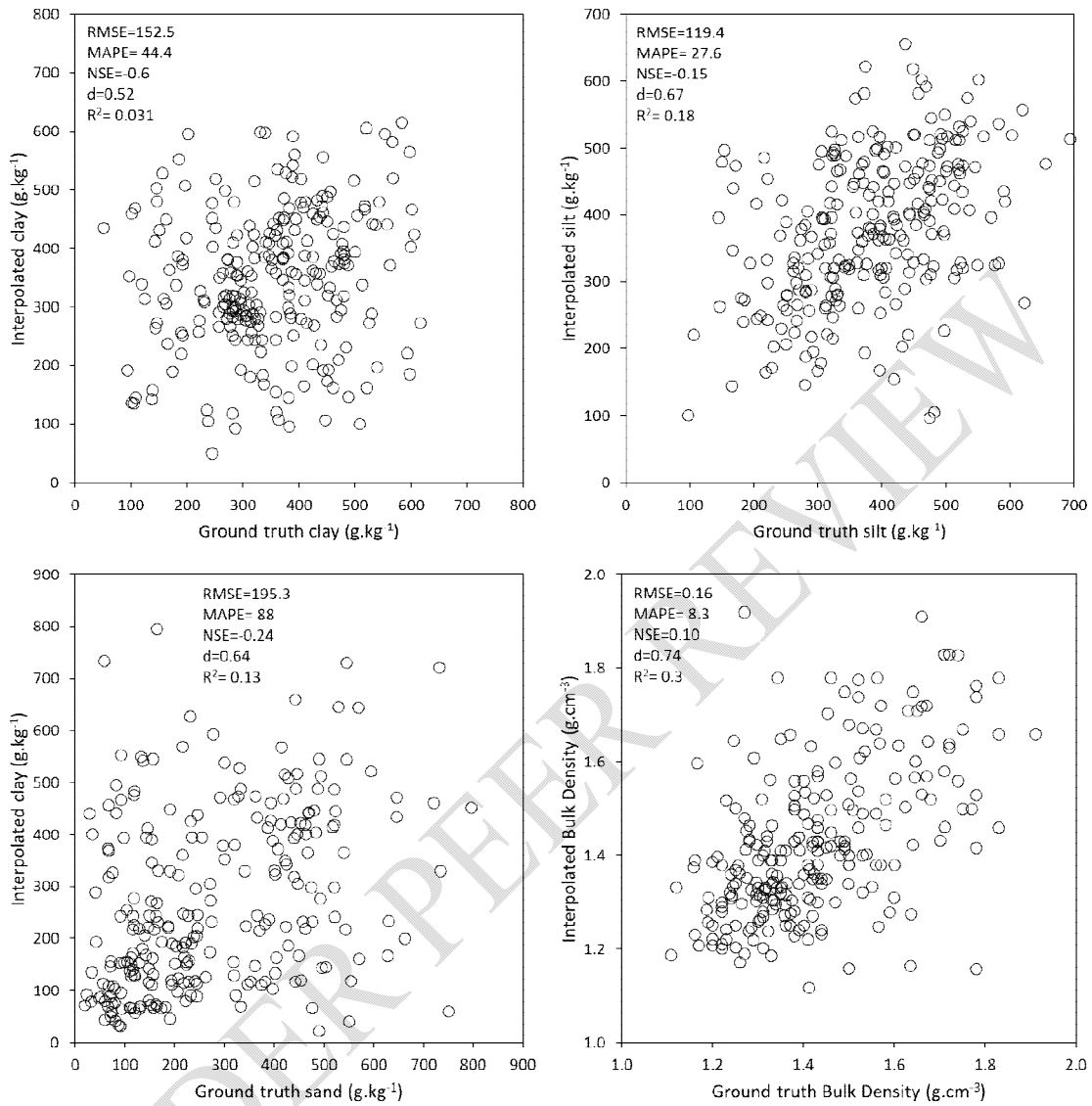


Figure 5 Ground truth data vs extracted data from interpolated ground truth of soil clay, silt and sand fractions with bulk density.

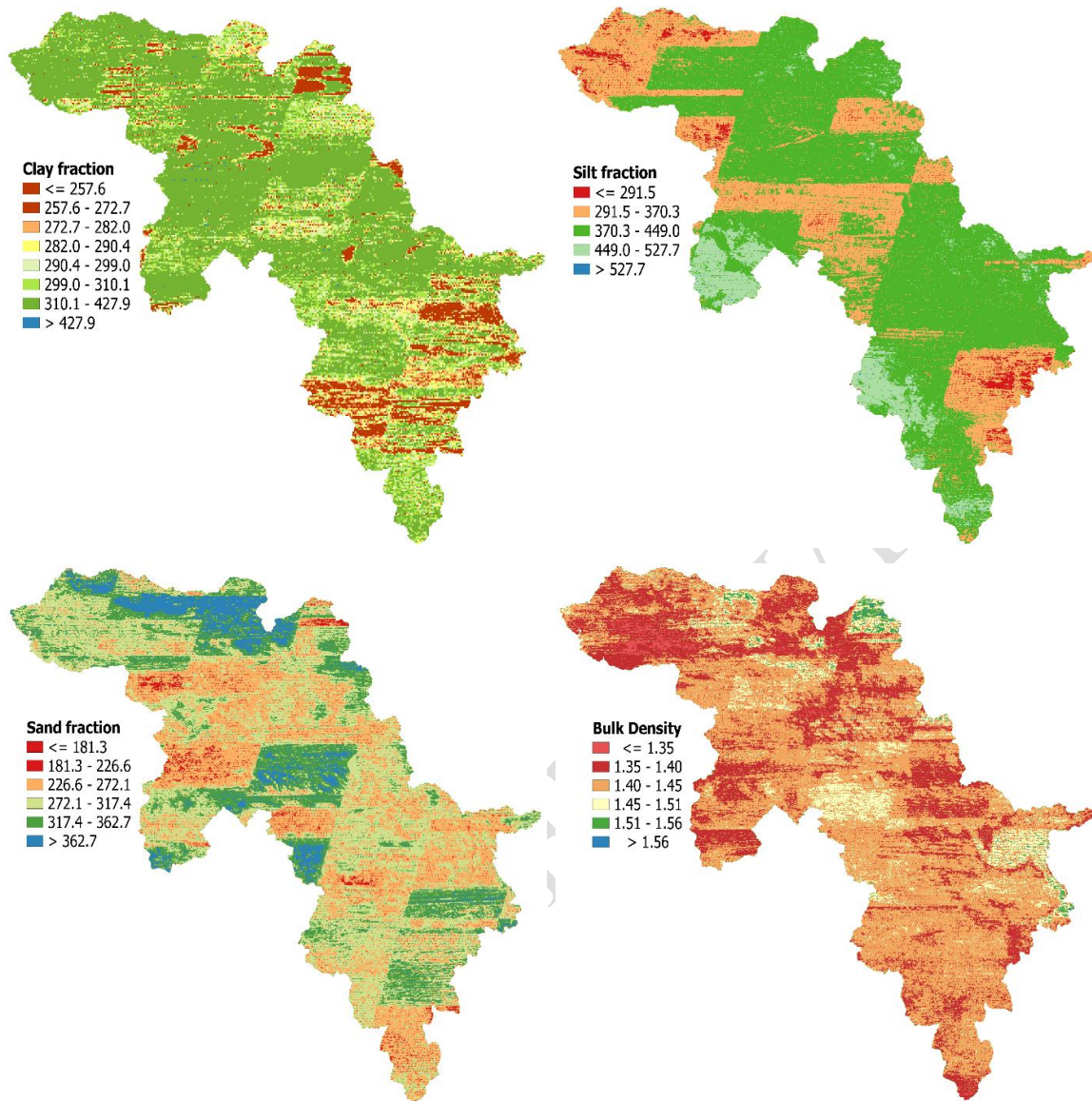


Figure 6 The map of predicted soil clay, silt, sand fractions (gkg^{-1}) and bulk density ($g.cm^{-3}$) data using XG Boost algorithm with resolution of 100m.

378

379

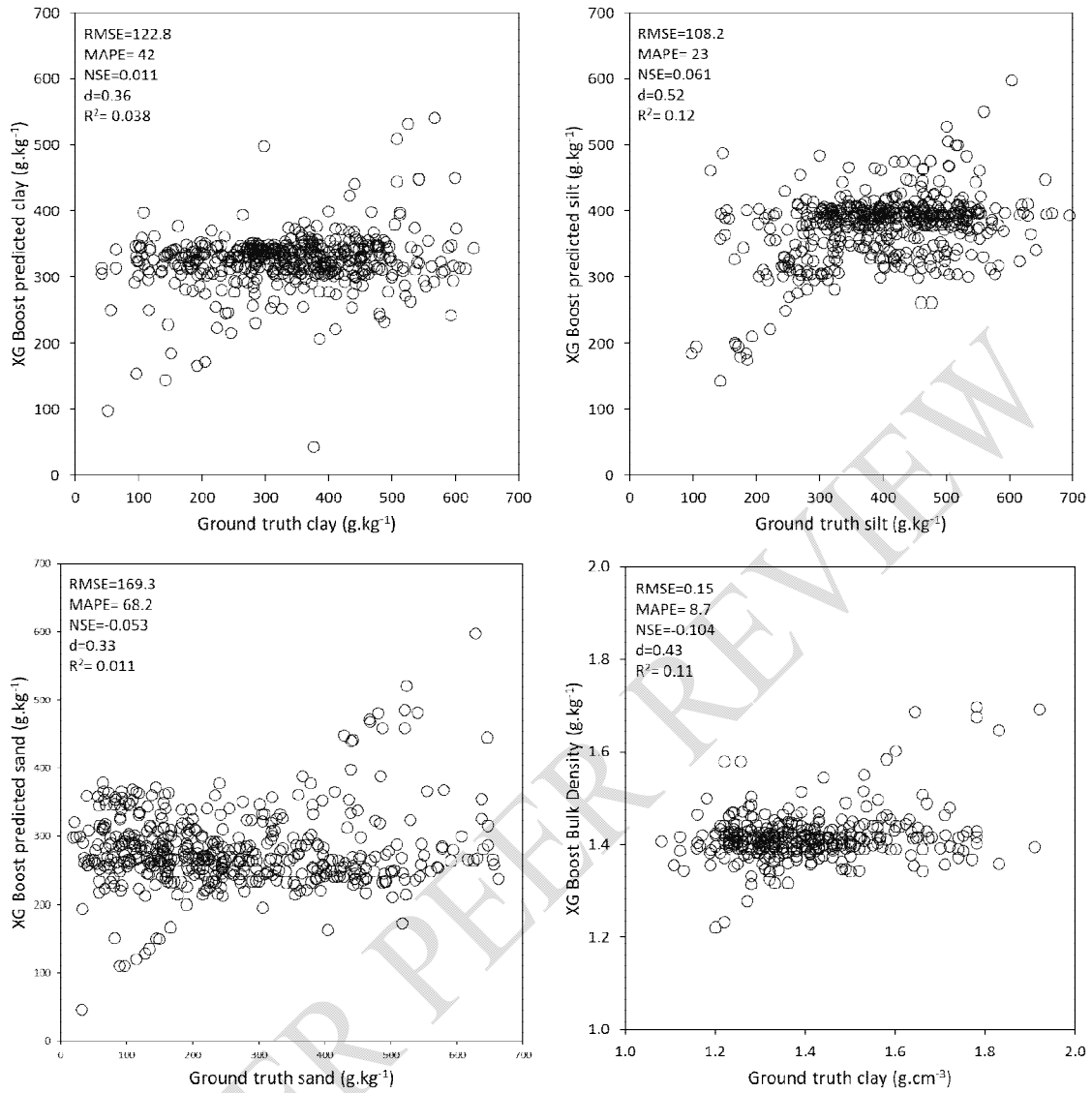


Figure 7 Ground truth data vs XG Boost algorithm predicted data of soil clay, silt and sand fractions with bulk density.

380

381 Table 1:

382

383

384

385

386

387 **Table 1** The accuracy of SoilGrids layers according to the ground truth soil sampling data, interpolated and
 388 XG Boost algorithm predicted. Where GTD: Ground truth data, StDev= Standard deviation, MAD: Mean
 389 Absolute Deviation, RMSE: Root Mean Square Error, MAPE: Mean absolute Percentage Error, NSE: Nash
 390 Sutcliffe coefficient, *d*: Index of agreement (Willmott), *R*: Correlation, *R*²: Coefficient of Determination.

Soil Properties	Ground truth data vs. SoilGrids data								
	GTD mean (StDev)	SoilGrids mean (StDev)	MAD	RMSE	MAPE	NSE	<i>d</i>	<i>R</i>	<i>R</i> ²
Clay	334.4(123.6)	420.5(25.5)	124.0	152.5	63.5	-0.53	0.42	5E ⁻³	2.5E ⁻⁵
Silt	399.8(112.4)	403.0(25.3)	89.9	112.0	28.7	0.006	0.26	0.13	0.016
Sand	265.9(167.1)	176.5(17.8)	144.2	190.6	69.3	-0.30	0.40	-0.03	1.2E ⁻³
Bulk density	1.41(0.17)	1.45(0.06)	0.144	0.18	10.3	-0.16	0.32	3E ⁻³	9E ⁻⁶
Soil Properties	Ground truth data vs. Interpolated data								
	GTD mean (StDev)	SoilGrids mean (StDev)	MAD	RMSE	MAPE	NSE	<i>d</i>	<i>R</i>	<i>R</i> ²
Clay	334.4(123.6)	279.6(94.0)	116.6	152.4	44.4	-0.60	0.52	0.18	0.03
Silt	399.8(112.4)	379.3(110.9)	90.8	119.4	27.6	-0.15	0.67	0.42	0.18
Sand	265.9(167.1)	267.9(169.6)	142.3	195.3	88.8	-0.24	0.64	0.35	0.13
Bulk density	1.41(0.17)	1.434(0.170)	0.121	0.164	8.35	0.10	0.74	0.55	0.30
Soil Properties	Ground truth data vs. XG Boost algorithm predicted data								
	GTD mean (StDev)	SoilGrids mean (StDev)	MAD	RMSE	MAPE	NSE	<i>d</i>	<i>R</i>	<i>R</i> ²
Clay	334.4(123.6)	325.3(43.2)	97.9	122.7	42.1	0.01	0.36	0.20	0.04
Silt	399.8(112.4)	376.4(52.1)	83.2	108.2	23.1	0.06	0.52	0.35	0.13
Sand	265.9(167.1)	281.2(56.2)	137.5	169.3	68.2	-0.05	0.33	0.1	0.01
Bulk density	1.41(0.17)	1.417(0.06)	0.122	0.155	8.67	0.11	0.43	0.33	0.12

391

392 **4. Discussion.**

393 **4.1. Accuracy assessment of SoilGrids**

394 Accurate assessment of soil properties is crucial for various environmental applications,
 395 ranging from land use planning to climate change mitigation(Montanarella and Vargas 2012).
 396 The evaluation of SoilGrids, a global soil mapping initiative, and the potential for future soil
 397 parameter prediction products are critical activities in advancing our understanding of soil
 398 variability and informing evidence-based decision-making. Cross-validation and independent
 399 validation are fundamental methodologies employed to assess the accuracy of soil mapping
 400 products. Cross-validation techniques, well-documented in scientific literature, have been
 401 instrumental in evaluating the performance of SoilGrids at different spatial resolutions,
 402 including 1km and 250m versions (Hengl *et al.* 2014; Hengl *et al.* 2017; Poggio *et al.* 2021).
 403 Previous studies have reported relatively high *R*² values (0.64 to 0.83) and comparable RMSE
 404 values (9.5–10.9) for physical soil parameters, indicating promising predictive capabilities
 405 (Hengl *et al.*, 2014). An independent study for the assessment of SoilGrids soil fractions in

406 Croatia reported lower R^2 values of 0.27, 0.039 and 0.039 for clay, silt and sand fractions,
407 respectively (Radočaj *et al.* 2023). However, the results from this study's independent
408 evaluation of SoilGrids revealed discrepancies ($R^2 \leq 0.016$), particularly in clay, sand
409 fractions, and bulk density, suggesting potential limitations in capturing local-scale
410 variability. This aligns with Radočaj *et al.* (2023), who noted lower R^2 values in independent
411 assessments, emphasizing the need for ground truth validation across diverse geographic
412 contexts.

413 Independent validation, essential for unbiased accuracy estimation, requires the use of ground
414 truth soil sampling data not utilised in the creation of soil mapping products. While efforts
415 were made to ensure the representativeness of ground truth data in this study, challenges such
416 as mismatched soil depths and landscape heterogeneity could have influenced accuracy
417 assessment outcomes as this study evaluated 0–30cm as a single soil depth rather than
418 SoilGrids soil depth increments (0–5cm, 5–15cm and 15–30cm). Furthermore, the absence of
419 comprehensive global soil sampling programs poses limitations to independent validation
420 efforts, highlighting the need for enhanced data collection programs.

421 SoilGrids' reliability is important for its applicability in various environmental studies and
422 management practices (Liang *et al.* 2019). Acknowledging the limitations of this study, future
423 research should focus on evaluating SoilGrids' accuracy across different spatial scales,
424 considering factors such as soil types, bio-geolocations, climate classes, and land cover types.
425 Furthermore, conducting digital soil mapping at a national level using comprehensive ground
426 truth soil sampling data could serve as a complementary approach to enhance the reliability
427 of SoilGrids, particularly in regions not adequately represented in its training dataset (Kidd *et al.*
428 *et al.* 2020; Radočaj *et al.* 2023).. While SoilGrids offers valuable insights into soil variability
429 on a global scale, its accuracy remains contingent upon the availability and quality of ground
430 truth data. Continued efforts in refining accuracy assessment methodologies and expanding
431 ground truth data coverage will contribute to enhancing the reliability and usability of
432 SoilGrids in addressing various environmental challenges. Acknowledging that factors such as
433 spatial resolution, data coverage, and the representativeness of ground truth data play
434 significant roles in determining the reliability of soil mapping products (Heung *et al.* 2021),
435 addressing these challenges requires collaborative efforts between researchers, policymakers,
436 and data providers to improve data quality and enhance the accuracy of soil mapping
437 initiatives (Hengl *et al.* 2017; Kidd *et al.* 2020).

438 The reliability of SoilGrids is crucial for informing land use decisions, agricultural practices,
439 and climate change mitigation strategies in a local scale (Hengl *et al.* 2017). Therefore,
440 ensuring the accuracy of soil mapping products is essential for effective resource
441 management and sustainable development in Kurdistan region where the region requires
442 sustainable and productive projects in agriculture, manufacture, and infrastructure sectors. By
443 advancing our understanding of soil properties and their spatial distribution, we can better
444 inform policy decisions and implement sustainable land management practices to mitigate
445 environmental degradation and ensure the long-term health of ecosystems in the region.

446 **4.2. Independent Validation of SoilGrids**

447 Independent validation serves as a critical step in assessing the accuracy and reliability of
448 SoilGrids products, providing an unbiased estimate of model performance (Hengl *et al.* 2014;
449 Poggio *et al.* 2021). In this study, independent validation involved comparing SoilGrids
450 predictions with ground truth soil sampling data and predicted soil fractions and bulk density
451 using XG Boost algorithm that were not used during the model training process (Skidmore *et*
452 *al.*, 2002). Challenges arise when conducting independent validation, particularly in regions
453 where SoilGrids was created based on zero soil samples (Hengl *et al.* 2014; Maynard *et al.*
454 2023). For instance, the absence of soil sampling data in the study area, as documented in the
455 ISRIC WoSIS Soil Profile Database, poses difficulties in accurately validating SoilGrids
456 predictions (Hengl *et al.* 2017). Therefore, acknowledging these challenges is essential for
457 transparency in the validation process.

458 The selection of appropriate ground truth data is paramount to ensure the representativeness
459 of soil variability within the study area (Skidmore 1999; Das *et al.* 2022). Furthermore,
460 discrepancies in spatial resolution between SoilGrids and ground truth data may hinder the
461 assessment of local-scale variations in soil properties (Bogunovic *et al.* 2017; Poggio *et al.*
462 2021; Xu *et al.* 2023). Then, innovative approaches to address these challenges, such as
463 leveraging supplementary environmental variables or integrating data from alternative soil
464 mapping initiatives could better serve the accuracy of soil mapping (Heung *et al.* 2021; Lu *et*
465 *al.* 2022). Additionally, efforts to expand ground truth datasets through targeted soil sampling
466 campaigns could enhance the accuracy and reliability of independent validation efforts
467 (Mahmood *et al.* 2017; Tifafi *et al.* 2018). Despite challenges associated with data availability
468 and spatial resolution discrepancies, innovative approaches and expanded ground truth
469 datasets can enhance the reliability of independent validation efforts, ultimately improving

470 our understanding of soil variability and supporting informed decision-making in
471 environmental management for Kurdistan region.

472

473 **4.3. Independent Validation and Algorithm Prediction**

474 Independent validation is crucial for assessing the accuracy of soil mapping models. Our
475 study used ground truth soil sampling data that were not part of the SoilGrids model training
476 process. This approach ensured an unbiased estimate of model performance (Skidmore 1999).
477 However, challenges can rise due to the absence or insufficient number of soil sampling data
478 in the study area, posing difficulties in accurately validating SoilGrids predictions (Hengl *et al.*
479 *2014*). To address these challenges, machine learning algorithms were employed,
480 particularly the Extreme Gradient Boosting (XG Boost) algorithm, to predict soil fractions
481 and bulk density based on ground truth data. The XG Boost algorithm demonstrated
482 promising results in approximating soil properties, with relatively greater agreement between
483 observed and predicted values (Chen and Guestrin 2016). Despite the challenges associated
484 with data availability and spatial resolution discrepancies, the XG Boost algorithm enhanced
485 the accuracy of SoilGrids predictions, particularly in regions with limited ground truth data
486 (Figure 6 and Figure 7). By integrating machine learning algorithms and independent
487 validation techniques, we enhance our understanding of soil variability and support evidence-
488 based environmental management strategies (Wang *et al.* 2020; Lu *et al.* 2022)(Wang *et al.*,
489 2020; Lu *et al.*, 2022).

490

491 The findings were consistent with previous studies assessing the accuracy of SoilGrids
492 products (Hengl *et al.* 2017; Poggio *et al.* 2021). Cross-validation techniques and
493 independent validation demonstrated improvements in prediction accuracy, highlighting the
494 utility of machine learning algorithms in refining soil property estimations (Table 1). While
495 challenges remain, such as spatial resolution disparities and data availability issues,
496 innovative approaches and expanded ground truth datasets can mitigate these limitations
497 (Montanarella 2015; Tifafi *et al.* 2018). However, accurate soil mapping is essential for
498 informed decision-making in environmental management contexts (Montanarella and Vargas
499 2012). This study provides valuable insights into the reliability of SoilGrids predictions,
500 particularly in regions with limited ground truth data. By integrating machine learning
501 algorithms and independent validation techniques, we enhance our understanding of soil

502 variability and support evidence-based environmental management strategies (Wang *et al.*
503 2020; Lu *et al.* 2022). However, incorporating auxiliary environmental covariates and
504 integrating data from alternative soil mapping initiatives can further enhance prediction
505 accuracy (Bogunovic *et al.* 2017; Radočaj *et al.* 2023). Additionally, development of
506 international collaboration and data-sharing initiatives along with efforts to conduct targeted
507 soil sampling campaigns and assess soil properties at multiple spatial scales can facilitate
508 access to high-quality soil data and improve the accuracy of global soil mapping efforts
509 (Mahmood *et al.* 2017; Searle *et al.* 2021). Hence, incorporating advanced machine learning
510 algorithms, integrating global and local multi-scale datasets, and expanding ground truth data
511 coverage are essential steps towards enhancing the accuracy and reliability of soil mapping
512 products (Geng 2020).

513

514

515

516 **5. Conclusion**

517 The study assessed digital soil mapping methods in the Kurdistan Region of Iraq, comparing
518 ground truth soil samples with SoilGrids data and using advanced mapping techniques like
519 interpolation and machine learning to improve soil variability understanding. The results
520 revealed significant disparities across soil mineral fractions such as clay, silt, sand fractions,
521 and bulk density. For instance, the mean clay fraction in the ground truth dataset differed
522 notably from SoilGrids' estimation, with a MAD of 124.0 gkg⁻¹ and RMSE of 152.5. Similar
523 disparities were observed for other soil properties, underscoring the limitations of global soil
524 mapping initiatives at capturing regional-scale variability accurately. However, the integration
525 of machine learning algorithms, particularly the Extreme Gradient Boosting (XG Boost)
526 algorithm, showed promising results in improving the accuracy of soil property estimations.
527 The XG Boost algorithm exhibited a relatively lower MAD of 97.9 gkg⁻¹ for clay fractions,
528 indicating a better approximation of observed values compared to SoilGrids. Additionally,
529 significant percent improvements in RMSE and MAPE values across soil fractions and bulk
530 density measurements underscored the efficacy of machine learning-based approaches in
531 refining soil property estimations within the study area. These findings highlight the
532 importance of leveraging advanced mapping techniques and integrating machine learning

533 algorithms to enhance the accuracy and reliability of soil mapping methodologies. This study
534 further provides valuable insights for informing evidence-based decision-making in
535 environmental management contexts. By incorporating advanced mapping approaches and
536 leveraging machine learning algorithms, we can better support sustainable land management
537 practices and environmental conservation efforts in the region.

538 **Highlights:**

- 539 • Significant disparities between ground truth data and SoilGrids in Kurdistan Region,
540 Iraq.
- 541 • Integration of XG Boost algorithm improves accuracy of soil property predictions.
- 542 • Machine learning shows promise in refining estimations of soil mineral fractions.
- 543 • Novel insights into spatial variability of soil properties help land management.

544

545

546 **Disclaimer (Artificial intelligence)**

547 **Option 1:**

548 **Author(s) hereby declare that NO generative AI technologies such as Large Language Models**
549 **(ChatGPT, COPILOT, etc.) and text-to-image generators have been used during the writing or editing**
550 **of this manuscript.**

551

552

553

554 **6. References**

555 Arrouays, D, Mulder, VL, Richer-de-Forges, AC (2021) Soil mapping, digital soil mapping
556 and soil monitoring over large areas and the dimensions of soil security–A review.
557 *Soil Security***5** 100018.

558 Bennett, JM, McBratney, A, Field, D, Kidd, D, Stockmann, U, Liddicoat, C, Grover, S (2019)
559 Soil Security for Australia. *Sustainability***11** (12), 3416.

560 Bogunovic, I, Trevisani, S, Seput, M, Juzbasic, D, Durdevic, B (2017) Short-range and
561 regional spatial variability of soil chemical properties in an agro-ecosystem in eastern
562 Croatia. *Catena***154** 50-62.

563 Buenemann, M, Coetzee, ME, Kutuahupira, J, Maynard, JJ, Herrick, JE (2023) Errors in soil
564 maps: The need for better on-site estimates and soil map predictions. *PLoS One***18** (1),
565 e0270176.

- 566 Chen, S, Arrouays, D, Mulder, VL, Poggio, L, Minasny, B, Roudier, P, Libohova, Z,
567 Lagacherie, P, Shi, Z, Hannam, J (2022) Digital mapping of GlobalSoilMap soil
568 properties at a broad scale: A review. *Geoderma***409** 115567.
- 569 Chen, T, Guestrin, C 'Xgboost: Reliable large-scale tree boosting system, Proceedings of the
570 22nd SIGKDD Conference on Knowledge Discovery and Data Mining, San
571 Francisco, CA, USA'. 2015. Available
- 572 Chen, T, Guestrin, C 'Xgboost: A scalable tree boosting system, Proceedings of the 22nd acm
573 sigkdd international conference on knowledge discovery and data mining'. 2016.
574 Available
- 575 Corwin, D, Lesch, S, Oster, J, Kaffka, S (2006) Monitoring management-induced spatio-
576 temporal changes in soil quality through soil sampling directed by apparent electrical
577 conductivity. *Geoderma***131** (3-4), 369-387.
- 578 Das, B, Murgakonkar, D, Navyashree, S, Kumar, P (2022) Novel combination artificial neural
579 network models could not outperform individual models for weather-based cashew
580 yield prediction. *International Journal of Biometeorology***66** (8), 1627-1638.
- 581 Fadhil, AM (2011) Drought mapping using Geoinformation technology for some sites in the
582 Iraqi Kurdistan region. *International Journal of Digital Earth***4** (3), 239-257.
- 583 Forti, L, Perego, A, Brandolini, F, Mariani, GS, Zebari, M, Nicoll, K, Regattieri, E, Barbaro,
584 CC, Bonacossi, DM, Qasim, HA (2021) Geomorphology of the northwestern
585 Kurdistan Region of Iraq: landscapes of the Zagros Mountains drained by the Tigris
586 and Great Zab Rivers. *Journal of Maps***17** (2), 225-236.
- 587 Fung, KF, Chew, KS, Huang, YF, Ahmed, AN, Teo, FY, Ng, JL, Elshafie, A (2022)
588 Evaluation of spatial interpolation methods and spatiotemporal modeling of rainfall
589 distribution in Peninsular Malaysia. *Ain Shams Engineering Journal***13** (2), 101571.
- 590 Gee, G, Bauder, J (1986) Particle-size analysis. In 'Methods of soil analysis. Part 1. Physical
591 and mineralogical methods'.(Ed. A Klute) pp. 383-411. *Soil Science Society of
592 America: Madison, WI*
- 593 Geng, X (2020) Development of operational methods to predict soil classes and properties in
594 Canada using machine learning. Carleton University.
- 595 Hengl, T, De Jesus, JM, MacMillan, RA, Batjes, NH, Heuvelink, GB, Ribeiro, E, Samuel-
596 Rosa, A, Kempen, B, Leenaars, JG, Walsh, MG (2014) SoilGrids1km—global soil
597 information based on automated mapping. *PLoS One***9** (8), e105992.
- 598 Hengl, T, Heuvelink, GB, Kempen, B, Leenaars, JG, Walsh, MG, Shepherd, KD, Sila, A,
599 MacMillan, RA, Mendes de Jesus, J, Tamene, L (2015) Mapping soil properties of
600 Africa at 250 m resolution: Random forests significantly improve current predictions.
601 *PLoS One***10** (6), e0125814.
- 602 Hengl, T, Mendes de Jesus, J, Heuvelink, GB, Ruiperez Gonzalez, M, Kilibarda, M, Blagotić,
603 A, Shanguan, W, Wright, MN, Geng, X, Bauer-Marschallinger, B (2017)
604 SoilGrids250m: Global gridded soil information based on machine learning. *PLoS
605 One***12** (2), e0169748.

- 606 Heung, B, Saurette, D, Bulmer, CE (2021) Digital Soil Mapping. *Digging into Canadian*
607 *Soils*
- 608 Jirjees, S, Seeyan, S, Fatah, K (2020) Climatic analysis for Pirmam area, Kurdistan Region,
609 Iraq. *The Iraqi Geological Journal* 75-92.
- 610 Kidd, D, Searle, R, Grundy, M, McBratney, A, Robinson, N, O'Brien, L, Zund, P, Arrouays,
611 D, Thomas, M, Padarian, J (2020) Operationalising digital soil mapping–Lessons
612 from Australia. *Geoderma Regional* **23** e00335.
- 613 Liang, Z, Chen, S, Yang, Y, Zhou, Y, Shi, Z (2019) High-resolution three-dimensional
614 mapping of soil organic carbon in China: Effects of SoilGrids products on national
615 modeling. *Science of the total environment* **685** 480-489.
- 616 Lu, L, Li, S, Wu, R, Shen, D (2022) Study on the Scale Effect of Spatial Variation in Soil
617 Salinity Based on Geostatistics: A Case Study of Yingdaya River Irrigation Area.
618 *Land* **11** (10), 1697.
- 619 Mahmood, F, Khan, I, Ashraf, U, Shahzad, T, Hussain, S, Shahid, M, Abid, M, Ullah, S
620 (2017) Effects of organic and inorganic manures on maize and their residual impact
621 on soil physico-chemical properties. *Journal of soil science and plant nutrition* **17** (1),
622 22-32.
- 623 Malone, B, Stockmann, U, Glover, M, McLachlan, G, Engelhardt, S, Tuomi, S (2022) Digital
624 soil survey and mapping underpinning inherent and dynamic soil attribute condition
625 assessments. *Soil Security* **6** 100048.
- 626 Marsh, A, Altaweel, M (2020) The search for hidden landscapes in the Shahrizor: Holocene
627 land use and climate in Northeastern Iraqi Kurdistan. *New Agendas in Remote Sensing*
628 *and Landscape Archaeology in the Near East* 7.
- 629 Maynard, JJ, Yeboah, E, Owusu, S, Buenemann, M, Neff, JC, Herrick, JE (2023) Accuracy of
630 regional-to-global soil maps for on-farm decision-making: are soil maps “good
631 enough”? *Soil* **9** (1), 277-300.
- 632 Montanarella, L 'The global soil partnership, IOP Conference Series: Earth and
633 Environmental Science'. 2015. (IOP Publishing. Available
- 634 Montanarella, L, Vargas, R (2012) Global governance of soil resources as a necessary
635 condition for sustainable development. *Current opinion in environmental*
636 *sustainability* **4** (5), 559-564.
- 637 Nasir, SM, Kamran, KV, Blaschke, T, Karimzadeh, S (2022) Change of land use/land cover in
638 kurdistan region of Iraq: A semi-automated object-based approach. *Remote Sensing*
639 *Applications: Society and Environment* **26** 100713.
- 640 Pereira, P, Brevik, EC, Muñoz-Rojas, M, Miller, BA, Smetanova, A, Depellegrin, D, Misiune,
641 I, Novara, A, Cerdà, A (2017) Soil mapping and processes modeling for sustainable
642 land management. In 'Soil mapping and process modeling for sustainable land use
643 management'. pp. 29-60. (Elsevier:

- 644 Poggio, L, De Sousa, LM, Batjes, NH, Heuvelink, GB, Kempen, B, Ribeiro, E, Rossiter, D
645 (2021) SoilGrids 2.0: producing soil information for the globe with quantified spatial
646 uncertainty. *Soil*7 (1), 217-240.
- 647 Radočaj, D, Jurišić, M, Rapčan, I, Domazetović, F, Milošević, R, Plaščak, I (2023) An
648 Independent Validation of SoilGrids Accuracy for Soil Texture Components in
649 Croatia. *Land*12 (5), 1034.
- 650 Salcedo-Sanz, S, Ghamisi, P, Piles, M, Werner, M, Cuadra, L, Moreno-Martínez, A,
651 Izquierdo-Verdiguier, E, Muñoz-Marí, J, Mosavi, A, Camps-Valls, G (2020) Machine
652 learning information fusion in Earth observation: A comprehensive review of
653 methods, applications and data sources. *Information Fusion*63 256-272.
- 654 Searle, R, McBratney, A, Grundy, M, Kidd, D, Malone, B, Arrouays, D, Stockman, U, Zund,
655 P, Wilson, P, Wilford, J (2021) Digital soil mapping and assessment for Australia and
656 beyond: A propitious future. *Geoderma Regional*24 e00359.
- 657 Skidmore, AK (1999) Accuracy assessment of spatial information. In 'Spatial statistics for
658 remote sensing'. pp. 197-209. (Springer:
- 659 Surucu, A, Ahmed, TK, Gunal, E, Budak, M (2019) Spatial Variability of Some Soil
660 Properties in an Agricultural Field of Halabja City of Sulaimania Governorate, Iraq.
661 *Fresenius Environment Bulletin*28 (1), 193-206.
- 662 Tifafi, M, Guenet, B, Hatté, C (2018) Large differences in global and regional total soil
663 carbon stock estimates based on SoilGrids, HWSD, and NCSCD: Intercomparison
664 and evaluation based on field data from USA, England, Wales, and France. *Global
665 Biogeochemical Cycles*32 (1), 42-56.
- 666 Upadhyay, S, Raghubanshi, A (2020) Determinants of soil carbon dynamics in urban
667 ecosystems. In 'Urban ecology'. pp. 299-314. (Elsevier:
- 668 Vågen, T-G, Winowiecki, LA (2013) Mapping of soil organic carbon stocks for spatially
669 explicit assessments of climate change mitigation potential. *Environmental Research
670 Letters*8 (1), 015011.
- 671 van der Voort, TS, Verweij, S, Fujita, Y, Ros, GH (2023) Enabling soil carbon farming:
672 presentation of a robust, affordable, and scalable method for soil carbon stock
673 assessment. *Agronomy for Sustainable Development*43 (1), 22.
- 674 Wang, F, Wei, Y, Yang, S (2023) Enhanced Understanding of Key Soil Properties in Northern
675 Xinjiang Using Water-Heat-Spectral Datasets Based on Bioclimatic Guidelines.
676 *Land*12 (9), 1769.
- 677 Wang, Z, Shi, W, Zhou, W, Li, X, Yue, T (2020) Comparison of additive and isometric log-
678 ratio transformations combined with machine learning and regression kriging models
679 for mapping soil particle size fractions. *Geoderma*365 114214.
- 680 Xu, C, Torres-Rojas, L, Vergopolan, N, Chaney, NW (2023) The Benefits of Using
681 State-Of-The-Art Digital Soil Properties Maps to Improve the Modeling of Soil
682 Moisture in Land Surface Models. *Water resources research*59 (4), e2022WR032336.

683 Yu, W, Zhou, W, Wang, T, Xiao, J, Peng, Y, Li, H, Li, Y (2024) Significant Improvement in
684 Soil Organic Carbon Estimation Using Data-Driven Machine Learning Based on
685 Habitat Patches. *Remote Sensing***16** (4), 688.

686

UNDER PEER REVIEW