





## 40 1. Introduction

41 Human activities have profoundly disrupted the global carbon cycle, leading to a significant  
42 increase in atmospheric CO<sub>2</sub> levels and a consequential alteration of ecosystems' inherent carbon  
43 absorption capacity (Crowther et al., 2016). Soils emerge as the primary carbon reservoirs within  
44 terrestrial ecosystems due to their remarkable capacity to store approximately two to four times  
45 more organic carbon in the top meter (i.e., approximately 1500 Pg C; (Scharlemann et al., 2014).  
46 Furthermore, the residence time of soil organic carbon (SOC) exceeds that of aboveground biomass,  
47 underscoring the crucial role of SOC storage in climate change mitigation strategies (Jobbágy et al.,  
48 2000; Saatchi et al., 2007).

49 Understanding the biogeochemical properties of soils is crucial for assessing and tracking  
50 changes in SOC storage capacity. Soils inherently undergo long-term transformations, with  
51 processes affecting SOC storage varying significantly over time and across different regions. This  
52 complexity and spatiotemporal heterogeneity necessitate studying the biogeochemical properties of  
53 soils at various spatial scales, from local fields to entire landscapes, and at different temporal scales,  
54 from short-term seasonal changes to long-term geological shifts (Stockmann et al., 2015).  
55 Consequently, numerous national and global initiatives have been launched to gather and  
56 standardize empirical data on soil properties accumulated over decades and make them readily  
57 available (Harden et al., 2018; Shangguan et al., 2014).

58 Integrating legacy and collaborative regional and global datasets of soil properties can pose  
59 challenges due to the absence of standardized protocols. For instance, in Spain, numerous  
60 organizations and institutions have collected soil profile data over fifty years, using different  
61 methods, laboratory techniques, standards, scales, and georeferencing systems (Llorente et al.,  
62 2018). Hence, this valuable information is currently scattered and fragmented, requiring substantial  
63 effort to integrate the historical soil databases into a cohesive, harmonized, and geographically well-  
64 defined dataset. Moreover, many of these databases lack the necessary information for accurately  
65 calculating Soil Organic Carbon stocks (SOCs), which depend on data such as SOC concentration  
66 (SOC<sub>c</sub>), bulk density, and coarse fragments (Calvo de Anta et al., 2020; Poeplau et al., 2017). This  
67 deficiency can lead to biased estimates of SOCs across different ecosystems. This challenge is also  
68 exacerbated by the considerable costs and operational complexities associated with soil data  
69 collection (Smith et al., 2020; Vargas et al., 2017). The absence of readily available databases  
70 containing consistent and comprehensive information on soil properties poses a significant  
71 challenge to effective soil monitoring across Spain and other regions of the world. The integration  
72 process is hindered by discrepancies in data formats, resolution, and quality, leading to potential  
73 inaccuracies and gaps in soil information. These challenges underscore the need for systematic  
74 methodologies and collaborative efforts to standardize data collection and processing protocols,  
75 thereby enhancing the reliability and usability of soil data for assessing SOCs and supporting  
76 climate change mitigation strategies.

77 The dynamics of SOC are determined by both soil physical and chemical properties, along  
78 with environmental soil-forming factors (Jenny, 1941). Soil properties exhibit diverse and complex  
79 patterns across scales due to the broad spatial and temporal range of soil formation conditions  
80 (Allen and Starr, 2019). To account for these complexities, SOC modeling has evolved from simple



81 qualitative approaches to sophisticated quantitative estimations and uncertainty assessment through  
82 models such as CLORPT or SCORPAN (Jenny, 1941; McBratney et al., 2003). Rooted in these  
83 theoretical models, digital soil mapping (DSM) has introduced a plethora of empirical models that  
84 estimate SOC<sub>s</sub> as a function of concurrent environmental factors (i.e., explanatory variables) (Chen  
85 et al., 2022), potentially reducing the number of in situ samples required for accurate spatial  
86 predictions (McBratney et al., 2003; Searle et al., 2021). These empirical models usually shows  
87 variability in their SOC<sub>s</sub> estimates, so ensemble methods are used to merge predictions by  
88 leveraging the strengths of each modeling technique (Shangguan et al., 2017; Wang et al., 2018b).  
89 Thus, ensemble modeling is assumed to provide more accurate and robust spatial predictions than  
90 individual models, especially when different models capture distinct aspects of soil dynamics  
91 (Padarian and McBratney, 2020). In conjunction with the spatial predictions of SOC<sub>s</sub>, it is also  
92 crucial to report their associated uncertainty, as it conveys valuable information for the proper  
93 interpretation of these empirically derived estimates (Poggio et al., 2021).

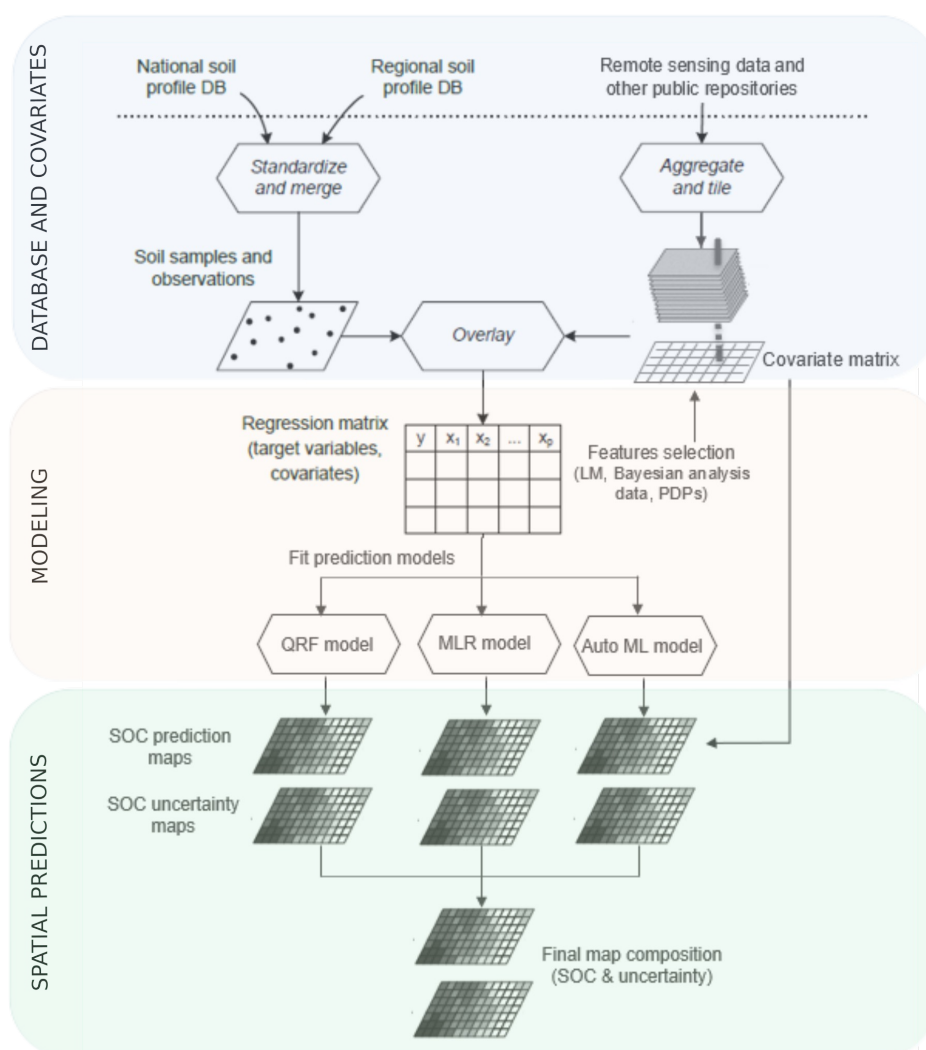
94 In Spain, the aforementioned challenges associated with sampling protocols and stock  
95 calculation procedures have yielded considerable variation in SOC<sub>s</sub> estimates. Local-scale  
96 estimations of SOC<sub>s</sub> have been developed for various ecosystems, including agricultural lands  
97 (Albaladejo et al., 2009; Álvaro-Fuentes et al., 2008; Muñoz-Rojas et al., 2012), forests, and  
98 pastures (Doblas-Miranda et al., 2013). At the national level, SOC<sub>s</sub> estimates within the upper 30  
99 cm depth have shown a notable range, varying from 2.82 Pg C (Rodríguez Martín et al., 2016) to  
100 3.25 Pg C (Calvo de Anta et al., 2020). While SOC<sub>s</sub> are typically standardized to the upper 30 cm,  
101 the subsoil carbon pool (i.e., >30 cm) may contribute up to 50% of the total stock in Mediterranean  
102 soils (Mulder et al., 2016). This discrepancy can lead to an underestimation of a substantial portion  
103 of carbon within the effective soil depth (i.e., the soil depth where most SOC storage occurs).  
104 Therefore, addressing these multifaceted challenges is pivotal in designing more accurate national-  
105 scale assessments of SOC<sub>s</sub> that support effective management strategies for mitigating climate  
106 change.

107 The overall goal of this study was to enhance our understanding of SOC storage and  
108 distribution in peninsular Spain by distinguishing between different carbon variables: SOC<sub>c</sub> (g/kg),  
109 and SOC<sub>s</sub> (tC/ha). To this aim, we addressed two specific goals: 1) to integrate and standardize  
110 disparate soil profile databases developed over the years within peninsular Spain, and 2) to model  
111 and map SOC<sub>c</sub> and SOC<sub>s</sub> at two soil depths using a machine-learning-based ensemble modeling  
112 approach, including their associated uncertainties. Furthermore, we present four 90-meter pixel  
113 resolution SOC maps for peninsular Spain (i.e., SOCM90). These maps outline the spatial estimate  
114 of SOC<sub>c</sub> at depths of 0-30 cm and 30-100 cm, along with the SOC<sub>s</sub> at 0-30 cm and its effective  
115 depth, including their associated uncertainties. This information can serve as a reference point for  
116 effectively estimating and managing soil carbon sinks at the national level, as well as for compiling  
117 the National Greenhouse Gas Emissions Inventory Report (Ministry for Ecological Transition and  
118 the Demographic Challenge, 2024). Furthermore, the insights gained from this study contribute to  
119 global efforts, including the Global Soil Organic Carbon Map and the GlobalSoilMap (FAO and  
120 ITPS, 2018; Arrouays et al., 2014a), as it aligns with the specifications of the Global Soil Organic  
121 Carbon Map consortium (Arrouays et al., 2014b).



## 122 2. Materials and Methods

123 We followed a three-step methodological framework (Fig. 1). First, we collected soil data  
124 from various public sources at different administrative levels and compiled static environmental  
125 predictors. Next, we built and assessed ensemble spatial models for SOCc and SOC<sub>s</sub> using three  
126 distinct supervised learning approaches. The final step involved generating spatial predictions and  
127 evaluating their accuracy.

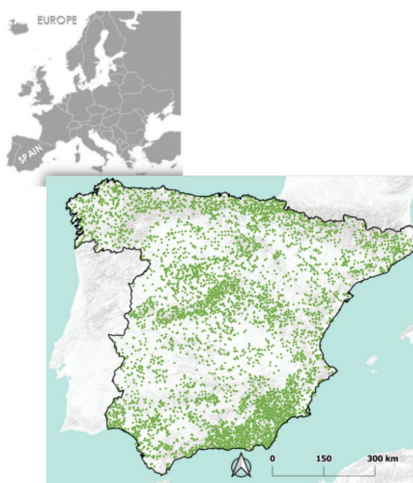


128 **Figure 1.** The three-step methodological framework for this study (adapted from the World Soil  
129 Information Service; Hengl et al., 2017).



130 *Study area*

131 The study area encompassed peninsular Spain, spanning an area of 491,258 km<sup>2</sup> (Fig. 2).  
132 Peninsular Spain is characterized by an intricate topography dominated by rugged mountain  
133 systems, expansive plateaus, and broad watershed depressions. The expansive Central Plateau  
134 covers most of the peninsula, with elevations ranging from 600 to 760 masl. The plateau gently  
135 slopes towards the west, directing the flow of most watercourses towards the Atlantic Ocean.  
136 Surrounding the plateau lie hills and steep mountain ranges, with a maximum altitude of 3,478 masl  
137 (Serrano, 2000).



138 **Figure 2.** Study area (outlined in black) showing the locations of soil samples collected between  
139 1954 and 2018 (points in green; n=8,361).

140 Spain is one of the most diverse countries in Europe in terms of climate, ranging from humid  
141 to semiarid conditions (AEMET IPMA, 2011). The average annual precipitation ranges from 200  
142 mm (in the southeast) to 2200 mm (in the northern and mountainous regions), while the mean  
143 annual temperature spans from below 2.5°C at higher altitudes to over 17.5°C in the southern and  
144 southeastern. The Mediterranean climate is dominant, extending across the inland plateaus  
145 (continental Mediterranean) to the coastal areas (coastal Mediterranean). This climate is  
146 characterized by wet, cold to mild winters and dry, hot, or mild summers, with variable  
147 temperatures and rainfall during autumn and spring. However, these contrasting climatic conditions  
148 are weaker along the coast, transitioning predominantly to arid or semi-desert conditions in the  
149 southeast. Conversely, in the north and northwest, the climate tends to be predominantly oceanic,  
150 characterized by high humidity and mild temperatures.

151 The diverse topography and wide-ranging climatic conditions facilitate a mosaic of land  
152 covers and uses. Despite a decreasing trend in agricultural areas over the past two decades,  
153 agricultural land still occupies approximately 33% of the total land area (MAPA, 2021). The  
154 agricultural landscape boasts diverse systems, encompassing herbaceous crops, primarily dryland  
155 cereals, and orchard crops such as grapes, olives, almonds, and various fruits. Forested areas cover  
156 more than 59% of the peninsula, predominantly comprising natural forests, plantations—mostly



157 found in mountainous regions within humid and subhumid areas—and shrublands. Extensive  
158 grasslands and other herbaceous vegetation thrive, particularly in high-altitude regions and the  
159 northern part of the country. Wetlands and water surfaces cover 0.9% of peninsular Spain's total  
160 area, while artificial surfaces occupy 7.1%.

## 161 2.1 Data compilation

### 162 2.1.1 Soil database

163 The database comprised 8,361 georeferenced soil profiles, containing 27,931 pedogenetic  
164 soil horizons. We collected soil data from public domain resources or were facilitated by national  
165 institutions responsible for the information. Specifically, the Red Carbosol database contributed  
166 78% of the samples, compiled through a collaborative network of Spanish soil experts across  
167 multiple research centers and universities, aggregating data from 635 different sources (Llorente et  
168 al., 2018). The second major source (18% of profiles) was the Consejería de Sostenibilidad, Medio  
169 Ambiente y Economía Azul (Andalusian Government, personal communication). The remaining 4%  
170 of the data were extracted from the LUCDEME database, which was compiled by various regional  
171 institutions, including Región de Murcia (Alias and Ortiz, 1986), the Agrarian Technological  
172 Institute (Junta de Castilla y León), and the University of Castilla La-Mancha (Bravo et al., 2019).  
173 Sampling periods spanned from 1954 to 2018, with most samples collected between 1965 and 2000.

174 To facilitate standardized and reconciled information on soil properties across sampling  
175 units, the compiled data structure was modified to create a unified database. The database included  
176 information on soil properties related to the computation of the variables, i.e., profile and horizon  
177 ID, horizon depth (cm), total carbon content (g/kg), bulk density (g/cm<sup>3</sup>), sand, silt, and clay (%),  
178 and coarse fragment content (> 2 mm; % of total volume).

179 Profile inclusion in the final database was contingent upon meeting quality criteria aligned  
180 with our research objectives and guided by quantitative pedological methodologies (Beaudette et  
181 al., 2013). These criteria encompassed accurate georeferencing, eliminating duplicate information,  
182 handling missing data values, verifying information consistency with the horizon, and fixing format  
183 inconsistencies. Furthermore, soil properties were explored through basic descriptive statistics, such  
184 as minimum, maximum, average, and standard deviation values, to ensure data consistency and  
185 evaluate their variability. Inconsistent data were appropriately reclassified as "no data". Further  
186 database adjustments included log-transforming the original SOC data. This step was taken to  
187 capitalize on the log-normal distribution tendency previously observed in SOC (Yigini et al., 2018).  
188 This transformation aimed to improve the correlation between SOC and its predictive factors,  
189 ultimately enhancing the accuracy of SOC spatial distribution modeling. To prioritize the organic  
190 carbon estimation in mineral soils, profiles or horizons containing more than 20% organic matter  
191 were omitted from the database, which is in line with WRB criteria (WRB-IUSS, 2014). As a result,  
192 Histosol profiles and horizons primarily composed of organic materials (H, O, L) were excluded  
193 from the analysis.

### 194 2.1.2 Standard soil depths and estimation of SOC concentration and stock

195 The soil depth of morphological horizons was standardized to facilitate the integration of the  
196 outcomes. The resulting range was established using widely recognized worldwide criteria for SOC  
197 estimation (Brus et al., 2017). The information concerning SOCc was obtained through analytical





198 measurements and used directly in the estimation process. However, some data were converted into  
199 SOCc from organic matter values, employing a conversion factor of 0.67, based on the assumption  
200 of 58% carbon content in organic matter (Rosell et al., 2001). SOCc values were then discretized  
201 into two standard depths: 0-30 cm and 30-100 cm. For each standard depth, the final value was  
202 determined as the sum of SOCc from morphological horizons, weighted by their original depth.  
203 Subsequently, the SOC<sub>s</sub> within a profile -the total amount of soil carbon per unit area at its effective  
204 depth- were computed using Equation [1]. The effective soil depth (ESD) was defined as the solum,  
205 encompassing surface and subsurface horizons with root presence and biological activity (Baillie,  
206 2001).

$$207 \quad \text{SOC}_s (\text{Kg} \cdot \text{m}^2) = \sum_{i=1}^n \text{SOC} (\text{g/Kg})_i \cdot \text{BD} (\text{Kg} \cdot \text{m}^3)_i \cdot \left[ 1 - \left( \frac{\text{CRFVOL}}{100} \right)_i \right] \cdot \text{HSIZE} (\text{cm}_i) \quad \text{Eq. [1]}$$

208 where  $i$  represents the horizon,  $n$  denotes the total number of horizons in the profile,  $BD$  stands for  
209 bulk density,  $CRFVOL$  is the percentage of coarse fragments (i.e., over 2 mm in diameter), and  
210  $HSIZE$  is the horizon thickness. SOC<sub>s</sub> were estimated for the standard depth of 0-30 cm before  
211 standardizing it to the 0-30 cm range. To prevent the propagation of errors in estimating parameters  
212 that may be absent in certain profiles from the existing soil data, the SOC<sub>s</sub> was computed only for  
213 profiles containing available information on bulk density and coarse fragment content. Thus, the  
214 final number of profiles used for modeling carbon content as a function of depth was 8,332 profiles  
215 for SOCc (g/kg) at a standard depth of 0-30 cm, 6,947 profiles for SOCc (g/kg) at a standard depth  
216 of 30-100 cm, 1,475 profiles for SOC<sub>s</sub> (tC/ha) at a standard depth of 0-30 cm, and 1,499 profiles for  
217 SOC<sub>s</sub> (tC/ha) at its effective depth. Finally, the primary characteristics of SOC in peninsular Spain  
218 were determined through basic descriptive statistics of carbon-related variables, including density  
219 distribution, mean profiles, and spatial autocorrelation. Data was processed using QGIS and the R  
220 packages: *aqp*, *PerformanceAnalytics*, *GSIF*, and *gstat* (R Foundation for Statistical Computing,  
221 2022; Beaudette et al., 2023; Hengl et al., 2020; Pebesma, 2004; Peterson et al., 2014).

### 222 2.1.3 Representativeness of the soil database

223 The spatial representativeness of the compiled soil database was evaluated using a  
224 probability distribution approach based on the maximum entropy method (Maxent), which has  
225 extensively been applied in modeling spatial point patterns (Phillips et al., 2006) and  
226 representativeness of environmental monitoring networks (Villarreal et al., 2018). We used the  
227 Maxent approach to estimate the relationship between the number of soil samples and soil-forming  
228 factors. Thus, the Maxent model was used to identify distinct representative areas based on  
229 pertinent environmental covariates chosen through covariate selection methods (see “Feature  
230 selection” subsection below). The model’s predictive accuracy was evaluated using the Area Under  
231 the Curve (AUC) metric for the training data (see section 2.3 below). This metric was compared  
232 against the AUC expected for a random model. Thus, AUC values lower than 0.5 would indicate  
233 that the model’s predictive performance was worse than random estimation (Fielding and Bell,  
234 1997). We built a model-based predictive map depicting the similarity of environmental predictor  
235 variables, thereby minimizing relative entropy between them and locations containing sampled soil  
236 data (Elith et al., 2011). This method effectively conveyed information on the spatial  
237 representativeness of soil samples across different environmental factors in peninsular Spain.



#### 238 2.1.4 Environmental Covariates

239 We identified the environmental variabcorporated into the SOC model following the  
240 SCORPAN conceptual spatial inference model. This conceptual model categorizes soil property  
241 predictions based on seven forming factors, encompassing soil properties (S), climatic variables  
242 (C), biota (O), relief (R), parent material (P), time, age (A), and spatial location (N). These  
243 covariates were grouped by static and dynamic variables as follows:

##### 244 Static variables

245 Relief Factor. Geomorphometric variables depicting terrain characteristics were assessed.  
246 Topographic relief was evaluated through geomorphometry and feature extraction derived from the  
247 Geomorpho90m global dataset at a resolution of 90 meters under the WGS84 geodetic datum. This  
248 dataset comprised 26 fully standardized geomorphometric variables derived from the MERIT-  
249 Digital Elevation Model (DEM), encompassing layers that depict (i) the rate of change across the  
250 elevation gradient using first and second derivatives, (ii) ruggedness, and (iii) geomorphological  
251 forms (Amatulli et al., 2020). Details regarding the source and resolution of these products are  
252 outlined in Table 1.

253 Human Factor. Land cover and land use information (IGN, 2012) were reclassified into 13  
254 classes, i.e., non-irrigated arable land, permanently irrigated land, heterogeneous agricultural areas,  
255 agro-forestry areas, broad-leaved forest, coniferous forest, mixed forest, sclerophyllous vegetation,  
256 pastures, moors and heathland, sparsely vegetated areas, and transitional woodland/shrub.

257 Parent Material Factor. Lithological classes were derived from the lithological map of Spain  
258 1M, which comprises 22 hierarchical levels (IGME, 1995).

259 Soil Properties. Soil information was obtained from the Digital District Soil Atlas (USDA,  
260 1987). The soil map was digitized based on the 1:2,000,000 map in the Atlas Nacional de España  
261 (Soil Science) published by the CSIC/IRNAS (De la Rosa et al., 2001).

##### 262 Dynamic variables

263 Climate Factor. Precipitation and temperature climatic variables were obtained from  
264 Ninyerola et al. (2005).

265 Biota Factor. We computed a set of ecosystem functioning attributes derived from remotely  
266 sensed indices. These attributes are associated with the carbon cycle, water cycle, and energy or  
267 heat balance. Functional attributes relate to each index's quantity, seasonality, and phenology. The  
268 satellite products spanned a sufficient time interval to ensure their stability over time (2000-2019).  
269 The complete remotely sensed set comprised 172 ecosystem functioning attributes as candidate  
270 predictors for SOC. Details regarding the sources and resolution of the satellite products are  
271 provided in Table 1. The Google Earth Engine platform (Gorelick et al., 2017) was employed to  
272 derive these attributes, including the annual mean (serving as a proxy for annual total amount),  
273 annual maximum and minimum (indicating annual extremes), seasonal standard deviation  
274 (describing seasonality), and sine and cosine of the dates of maximum and minimum (indicating  
275 phenology) (Alcaraz-Segura et al., 2017).

276 In summary, 254 environmental covariates, capturing the diverse forming factors across  
277 peninsular Spain, were computed as an initial step preceding the spatial modeling of SOC.





278 **Table 1.** Description of predictors for spatial modeling of SOC.

Category	Variables	Source	Scale/ Resolution
<b>Topographic</b>			
	(i) Slope, Aspect, Aspect cosine, Aspect sine, Eastness, Northness, Convergence, Compound topographic index, Stream power index, East-West first order partial derivative, North-South first order partial derivative	Geomorpho90m (Amatulli et al., 2020)	90m
	(ii) Profile curvature, Tangential curvature, East-West second order partial derivative, North-South second order partial derivative, Second order partial derivative		
	(iii) Elevation standard deviation, Terrain ruggedness index, Roughness, Vector ruggedness measure, Topographic position index, Maximum multiscale deviation (dev-magnitude), Scale of the maximum multiscale deviation, Maximum multiscale roughness, Scale of the maximum multiscale roughness		
	(iv) Geomorphon		
<b>Climate</b>			
	Mean annual precipitation (mm). Mean, minimum, and maximum annual temperature (°C). Radiation (kW/m <sup>2</sup> ). Period 1951-1999.	University of Barcelona (Ninyerola et al., 2005)	20 m
<b>Land features</b>			
	Soil types	Proyecto SEIS.net (MIMAM- CSIC)	1:100000
	Lithology	SGE-IGME (Spain)	1:200000
	Land use/cover	IGN-Corine Land Cover	1:100000
<b>Remotely sensed indices</b>			
<b>Carbon cycle</b>	Normalized Difference Vegetation Index (NDVI), Enhanced Vegetation Index (EVI)	MOD13Q1	250 m
<b>Water cycle</b>	Precipitation, Normalized Difference Water Index (NDWI)	CHIRPS	1km
	Evapotranspiration (ET)	MCD43A4	500 m
<b>Radiative balance</b>	Albedo	MCD43B3	500 m
<b>Sensible heat</b>	Land Surface Temperature (LST)	MOD11A2	1 km
<b>Ecosystem Functional attributes (inter-annual and monthly mean):</b>		(Alcaraz-Segura et al., 2017)	
<b>Amount:</b> mean, maximum, minimum			
<b>Seasonality:</b> standard deviation, coefficient of variation, range, relative range			
<b>Phenology:</b> sine and cosine of the dates of maximum and minimum			

279 (i) First order derivative, (ii) Second order derivative, (iii) Ruggedness, (iv) Geomorphological forms.

280 Finally, the environmental covariates were compiled into a covariate matrix with spatially  
 281 explicit information. A method alternative to reprojection and rescaling was employed to  
 282 accommodate the diverse formats of coordinate reference systems (CRS) and spatial resolutions of  
 283 the covariates. This approach aimed to reduce geometric distortion and mitigate computational  
 284 limitations stemming from the substantial volume of data (Bauer-Marschallinger et al., 2014).  
 285 Firstly, a reference matrix was created based on the pixel center locations (x, y) of the most detailed  
 286 resolution layer, i.e., MERIT-DEM (90 m), with WGS84 (EPSG 4326) as the CRS. Subsequently,



287 geoprocessing techniques extracted the covariate values for each location in the reference matrix.  
288 When applicable, the reference matrix coordinates were then reprojected to match the CRS of each  
289 covariate. The resulting matrix comprised the value of each covariate (columns) extracted for every  
290 point in the study area at intervals of 90 meters. This matrix was organized into a tiling system to  
291 enhance computing processing time during geoprocessing analyses. Categorical variables were  
292 rasterized. To do that, only categories with a sufficiently representative number of soil samples (i.e.,  
293 over 100 data points) were considered. These categories were subsequently transformed into a  
294 binary variable, indicating the presence or absence of the specific category (Yigini et al., 2018). The  
295 covariance matrix was generated using the R packages *sp*, *rgdal*, and *raster* (Hijmans, 2024; Keitt  
296 et al., 2012; Pebesma and Bivand, 2005).

## 297 2.2 Modeling

298 The data matrix comprised the log transformations of SOCc and SOCs and the covariate  
299 point data extracted at the same locations as the soil profiles. To mitigate geometric distortions, the  
300 coordinates of the profiles were reprojected to match the CRS of each covariate, where applicable.  
301 This organized structure facilitated analyzing the relationship between carbon and the covariates,  
302 ensuring accurate spatial alignment for meaningful interpretation.

### 303 2.2.1 Covariate selection

304 We analyzed covariate importance (CVI) to discern which covariates had the most  
305 significant impact on soil carbon models, thus reducing the substantial number of covariates (254)  
306 and mitigating the risk of potential overfitting (Gregorutti et al., 2017). CVI was evaluated using  
307 three selection methods: multiple linear regression, Bayesian analysis, and projection pursuit  
308 regression models (coupled with partial dependence plots).

309 Multiple linear regression is an easily interpretable method where the dependent variable  
310 (here SOCc and SOCs) is represented as a linear combination of regression coefficients, predictor  
311 variables, and a random error term. This error term accounts for the variation in the dependent  
312 variable that cannot be explained by the linear relationship with the predictors and their coefficients.  
313 Two fundamental assumptions are considered: a) the existence of a linear relationship between the  
314 response variable and the predictors (environmental covariates), and b) the absence of  
315 multicollinearity among the predictors (Yan and Su, 2009). To evaluate the importance of each  
316 covariate, we utilized the absolute value of the t-statistic. The t-statistic was computed by dividing  
317 the regression coefficient associated with each predictor by its standard error, serving as a measure  
318 of predictor accuracy (Greenwell and Boehmke, 2020). This analysis was conducted using the *vip* R  
319 package (Greenwell and Boehmke, 2020).

320 Bayesian analysis is an inferential approach rooted in the probability distribution of  
321 parameters derived from observed data and additional available information. Unlike traditional  
322 multilinear models, Bayesian probability models treat parameters as random variables and integrate  
323 data and prior information on the parameter distribution through a likelihood function to form a new  
324 posterior distribution (Gelman et al., 2013). The resulting parameter estimates are conditional on the  
325 observed data, which, following the rules of probability theory, ensures a consistent posterior  
326 distribution interpretation (McElreath, 2018). The Bayesian models were fitted using the Markov  
327 Chain Monte Carlo (MCMC) sampling methodology through an iterative selection of the most  
328 significant covariates. The final covariate selection was determined by evaluating the models based



329 on their Watanabe–Akaike information criteria (WAIC) scores, with preference given to the model  
330 exhibiting the lowest WAIC value. This criterion serves as a measure of model performance,  
331 balancing goodness of fit with model complexity. We implemented the Bayesian approach using the  
332 *rethinking* package in R (McElreath, 2020).

333 Projection Pursuit Regression (PPR) involves forming linear combinations of non-  
334 parametric functions of the predictor variables, enabling the exploration of nonlinear relationships  
335 within the data (Friedman and Stuetzle, 1981). Thus, PPR captures complex relationships between  
336 predictor variables and the dependent variable. Partial dependence plots (PDPs) were constructed  
337 for each covariate of interest, illustrating how changes in each covariate affect the predicted  
338 outcome while keeping other covariates fixed. The flatness of each PDP was assessed to determine  
339 the strength of the association between the covariate and the predicted outcome. CVI scores were  
340 computed based on the variability in the partial dependence values, supporting the identification of  
341 the most influential covariates. These scores captured the variability in the partial dependence  
342 values for each main effect by calculating the standard deviation of the y-axis values for each PDP.  
343 Covariates were ranked based on their CVI scores and the shape of their PDPs, guiding the  
344 selection of covariates for inclusion in the final model. The PPR models and PDPs were constructed  
345 using the *stats* and *pdp* packages in R (Greenwell, 2017).

346 Finally, covariates exhibiting the highest CVI scores, consistently identified across all three  
347 selection techniques, were integrated into the final dataset. Finally, expert criteria guided the  
348 decision-making process for the final selection of variables. This selection procedure was conducted  
349 individually for each dependent variable (i.e., SOCc and SOC<sub>s</sub>) at both the 0–30 cm and 30–100 cm  
350 standard depths (SOCc), and at both the 0–30 cm and the effective depth (SOC<sub>s</sub>).

### 351 2.2.2 Predictive models

352 There are multiple algorithms for predicting SOC using DSM and arguably there is not a  
353 single ideal one for predicting SOC across large geographical areas (Guevara et al., 2018). Given  
354 the intricate and often nonlinear relationship between SOC and environmental variables, the  
355 integration of multi-model ensemble methods with machine learning (ML) algorithms has been  
356 adopted to predict SOC spatial variability along with associated uncertainties (Gray et al., 2015;  
357 Shangquan et al., 2017; Wang et al., 2018a). Ensemble learning, a branch of ML, combines multiple  
358 base ML models, homogeneous (e.g., a combination of multiple decision trees) or heterogeneous  
359 (e.g., a combination of decision trees with support vector machines), to enhance predictive  
360 performance by mitigating errors between observed and predicted data (Zhang and Ma, 2012).  
361 Here, we employed three ensemble modeling approaches for predicting SOC.

362 Quantile regression forest (QRF), unlike traditional regression methods, QRF can handle  
363 sparse legacy data effectively without the need for kriging interpolation of residuals (Meinshausen,  
364 2006). The QRF algorithm can predict SOC values at various quantiles without relying on specific  
365 assumptions about the; therefore, it extends the capabilities of the random forest by providing  
366 accurate estimates across the entire distribution of the response variable. We leverage the  
367 *quantregForest* package in R (Meinshausen, 2006) to implement the QRF method. Validation  
368 statistics for QRF models were computed using out-of-bag error estimation. This method evaluates  
369 the model's performance by measuring the prediction error on data points not included in the  
370 samples used to train each decision tree in the ensemble. Thus, out-of-bag error estimation provides



371 an unbiased estimate of the model's predictive accuracy, helping to measure its performance in  
372 generalizing to new data.

373 Ensemble Machine Learning (MLR) combines linear model regression predictors with non-  
374 parametric models using bagging and boosting algorithms (Dietterich, 2000). Bagging and boosting  
375 are ensemble learning techniques used to improve the performance of ML models by combining  
376 multiple base models. Overall, bagging involves creating multiple subsets of the training data  
377 through bootstrapping (random sampling with replacement). Then, a base model (e.g., a decision  
378 tree) is trained on each subset independently. Finally, predictions from all base models are  
379 combined, typically by averaging, to make the final prediction. Conversely, Boosting works  
380 sequentially by training a series of weak models (i.e., models that perform slightly better than  
381 random models) and giving more weight to mispredicted data in subsequent iterations. Each new  
382 model focuses on the observations mispredicted by the previous models, thus gradually improving  
383 the model's performance. The final prediction is typically a weighted sum of predictions from all  
384 weak models. Both bagging and boosting aim to reduce overfitting and improve the overall  
385 performance and robustness of the model. We constructed a stacked model by integrating  
386 predictions from five modeling techniques: linear model regression, random forest, deep learning,  
387 cubist, and weighted k-nearest neighbor classifier. The predictions were used as features for a  
388 stacked model trained to compress the predictions from the base models. Each base model  
389 underwent independent assessment through ten-fold cross-validation. To incorporate spatial  
390 information, the dataset was split using spatial partitioning by k-means clustering, considering two  
391 classification layers: the XY locations of soil data and the Köppen climate classification (Kottek et  
392 al., 2006). The *MLR* R package was used to implement this ensemble approach (Bischl et al., 2016).

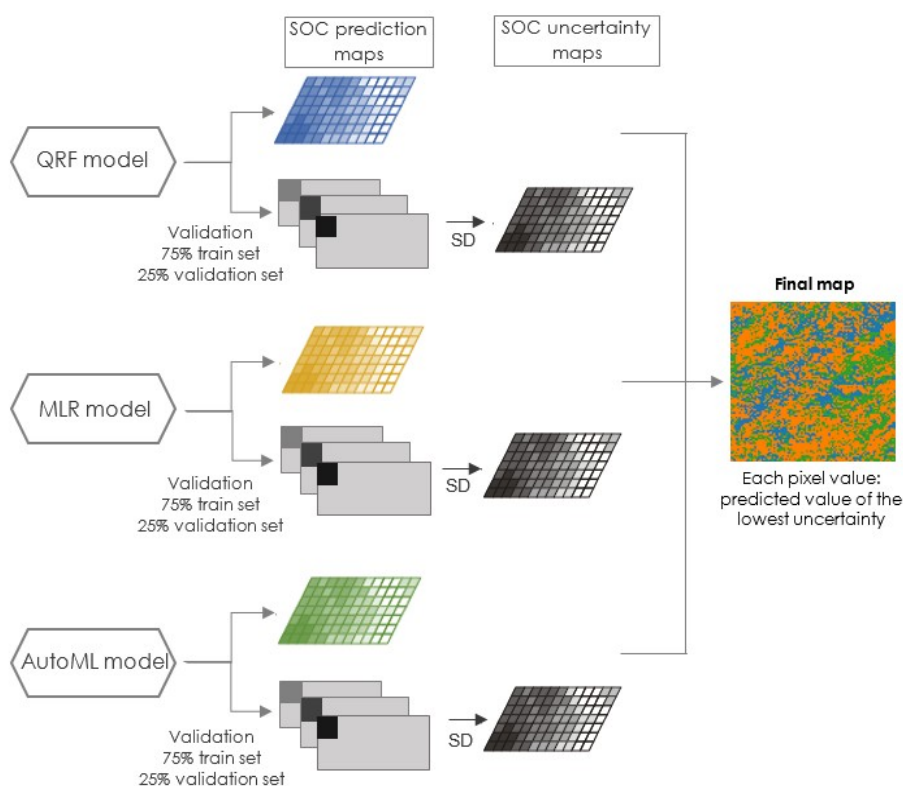
393 Auto-machine learning (AutoML) automates building ML models by efficiently selecting  
394 algorithms, tuning hyperparameters, and optimizing computational resources. With the vast array of  
395 ML algorithms available and the complexities of hyperparameter tuning, manual selection, and  
396 optimization can be challenged. We used the *H2O* package in R (Fryda T; Erin LeDell, 2024) to  
397 address this issue, which employs various Gradient Boosting Machine (GBM) algorithms such as  
398 generalized linear models, distributed random forests, deep neural networks, XGBoost, and gradient  
399 boosting machines. This diversity of models allows stacked ensembles to produce robust final  
400 predictions. The models were evaluated using 10-fold cross-validation and ranked based on their  
401 root-mean-square error (RMSE).

### 402 2.3 Spatial Prediction

403 We built model-based predictive maps for SOC<sub>c</sub> at depths of 0-30 cm and 30-100 cm, and  
404 for SOC<sub>s</sub> at depths of 0-30 cm and the effective depth, with a pixel resolution of 90 meters. To do  
405 that, we first generated three prediction maps using each ensemble modeling approach (i.e., QRF,  
406 MLR, and AutoML) and calculated the standard deviation for each pixel in these maps. Then, we  
407 assigned the predicted value from the most accurate map -determined by the lowest standard  
408 deviation value (i.e., the highest agreement among models)- to each pixel (Fig. 3). This approach  
409 ensures that the final predictive maps reflect the most reliable estimates of SOC content at each  
410 pixel, incorporating the collective insights from multiple modeling techniques and accounting for  
411 the associated uncertainty. This approach has been successfully applied in digital soil mapping



412 (Varón-Ramírez et al., 2022; Arroyo-Cruz et al., 2017) and in reducing uncertainty when modeling  
413 ecosystem-related variables (Gavilán-Acuña et al., 2021).



414 **Figure 3.** Scheme illustrating the generation of predictive maps of soil organic carbon using an  
415 ensemble modeling approach to reduce model uncertainty. SD (standard deviation).

416 The metrics and information criteria used to evaluate the models' performance generated by  
417 the three ensemble modeling approaches included coefficient of determination ( $R^2$ ), concordance  
418 correlation coefficient (CC), which measures the concordance level between predicted and observed  
419 values, root mean square error (RMSE), and mean absolute error (MAE). The dataset was randomly  
420 split into calibration (75%) and validation (25%) data for each modeling approach. Conditional  
421 quantile plots were generated to further assess model performance across the entire distribution of  
422 observed SOC values (Wilks, 2019). These plots evenly divide the predicted values and identify  
423 corresponding observation values, including the median, 25th/75th, and 10th/90th percentiles. By  
424 doing so, they offer insights into the alignment between predictions and observations across the  
425 entire range of values. The conditional quantile plots were generated using the *openair* package in R  
426 (Carslaw and Ropkins, 2012).



### 427 3. Soil database overview

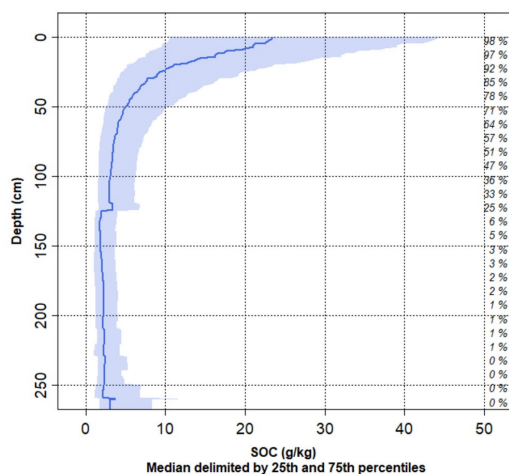
428 The harmonized soil database comprised 8,332 soil profiles with 25,370 morphological  
 429 horizons distributed across peninsular Spain. Only profiles with complete information were used for  
 430 SOC content calculation to prevent error propagation. This resulted in 1,499 profiles, slightly  
 431 exceeding the 1,475 used for SOC estimation at the standard depth of 0-30 cm. SOCc exhibited  
 432 high variability across soil horizons, ranging from 0.02 g/kg to 296.9 g/kg, with a mean of 16.53  
 433 g/kg (Table 2). The mean SOCc at the standard depth of 0-30 cm was 20.7 g/kg; at the 30-100 cm  
 434 depth, it was 5.8 g/kg, representing 35% of the soil profile mean. Similarly, SOCs varied from 0.006  
 435 kg/m<sup>2</sup> to 87.1478 kg/m<sup>2</sup>, with a mean of 2.965 kg/m<sup>2</sup>. SOCs at the standard depth of 0-30 cm  
 436 accounted for 65% of the soil profile mean.

437 **Table 2.** Statistical summary of SOCc and SOCs at different soil depths.

Variable	Depth	Number of profiles	Minimum	1st Quantile	Median	Mean	3rd Quantile	Maximum
SOCc (g/kg)	0-30	8,332	0.017	7.148	14.008	20.691	27.098	257.95
	30-100	6,947	0.017	1.700	3.371	5.833	6.814	185.743
SOCs (kg/m <sup>2</sup> )	0-30	1,475	0.119	2.066	3.738	5.31	6.95	39.967
	ESD <sup>(1)</sup>	1,499	0.119	3.000	5.300	8.198	10.260	93.892

438 (1) Effective soil depth. Soil organic carbon concentration (SOCc). Soil organic carbon stock (SOCs).

439 Both SOCc and SOCs followed a normal distribution with a right-skew after transforming  
 440 the original values to a natural log (Fig. S1 in Supplementary material). The highest SOCc values  
 441 were concentrated at the upper layers (0-30 cm), decreasing rapidly with depth (Fig. 4). Mean  
 442 profile values ranged from 23 g/kg in the upper horizon (0-5 cm) to 3 g/kg in the deepest horizon  
 443 (>200 m). In general, peninsular Spain's soils showed shallow depths, with only 35% of horizons  
 444 extending over 100 cm and decreasing to 3% at depths exceeding 150 cm.



445 **Figure 4.** Average Soil Organic Carbon concentration (SOCc; g/kg) at different depths.

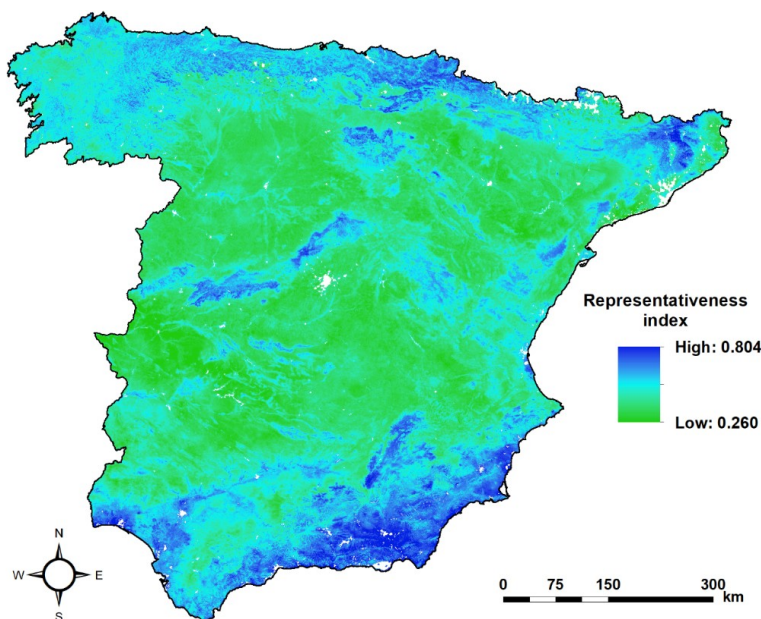




446       Nugget-to-sill ratio (NSR) analysis indicated weak spatial autocorrelation of SOC with an  
447 NSR >75% for the log-transformed SOCc. The large nugget effect (i.e., 1.3; Fig. S2 in  
448 Supplementary material) suggested that the data did not capture a significant portion of fine-scale  
449 SOCc variation. Despite the weak autocorrelation evidence, SOCc kriging predictions aligned with  
450 the expected distribution and demonstrated relatively low standard errors (Fig. S2 in Supplementary  
451 material). To delve deeper into the spatial dependence of SOCc with soil depth, autocorrelation  
452 analysis was conducted at six different depths, ranging from 5 cm to depths exceeding 100 cm  
453 (Table S1 in Supplementary material). A moderate dependence was observed, with the best spatial  
454 correlation (NSR = 35%) at the shallowest horizon (0-5 cm), gradually decreasing with depth to an  
455 NSR of 72% at 2 m.

#### 456 **4. Representativeness of the soil database**

457       The representativeness of a soil type in the database (i.e., the probability of a soil type being  
458 sampled) as a function of soil-forming factors is illustrated in Fig. 5. The representativeness values  
459 ranged from 0 (low) to 1 (high). The Maxent model yielded an AUC value of 0.611, indicating a  
460 predictive capacity greater than that yielded by a random model (i.e., 0.5). The predictive map  
461 showed variations in representativeness across different ecosystems. Overall, mountainous regions  
462 such as the Central Mountain System in central Spain, Sierra Morena and the Betic System in the  
463 south, as well as the Sierra de la Demanda, Cantabrian Range, and the Pyrenees in the north,  
464 showed high representativeness, as depicted by the blue areas in Fig. 5. In contrast, the Central  
465 Plateau, Ebro Depression, and Tajo and Guadiana Basins (i.e., interior regions in both the south and  
466 the north) were among the least represented areas (areas in green).



467 **Figure 5.** Maxent model-based representativeness of soil types sampled across peninsular Spain.



## 468 5. SOC modeling and prediction

469 The CVI analysis resulted in a highly reduced covariate space across selection methods for  
 470 SOC<sub>s</sub> and SOC<sub>c</sub>. Table 3 shows the selected covariates at different depths. Seventeen and nineteen  
 471 covariates were selected for modeling SOC<sub>c</sub> at the 0-30 cm and 30-100 cm depth, respectively. For  
 472 modeling SOC<sub>s</sub>, thirteen covariates were selected for the 0-30 cm depth and the effective depth.  
 473 Overall, annual precipitation, and the mean and minimum spring temperatures, were identified as  
 474 the most influential climatic factors. Similarly, among the remotely sensed variables related to the  
 475 carbon cycle, the annual mean and maximum NDVI and the monthly mean EVI (in March) were  
 476 found relevant. Moreover, annual and monthly indices linked to the water cycle, such as ET and  
 477 NDWI, exhibited notable importance. Regarding topographic covariates, indices derived from  
 478 terrain roughness, including maximum rough-magnitude, dev-scale related to topographic position,  
 479 and slope, were identified as particularly relevant. Additionally, soil covariates, such as lithology,  
 480 soil type, and land use/cover, were included in all SOC models due to their fundamental role in  
 481 modeling soil properties (Jenny, 1941).

482 **Table 3.** Covariates selection for modeling soil organic carbon concentration (SOC<sub>c</sub>) and soil  
 483 organic carbon stock (SOC<sub>s</sub>) at different standard depths.

Frequency	SOC <sub>c</sub> (0-30 cm)	SOC <sub>c</sub> (30-100 cm)	SOC <sub>s</sub> (0-30 cm/ESD) <sup>(1)</sup>
Climate variables <sup>(2)</sup>			
Annual	- Min Temp	- Mean Pp	- Mean Pp - Mean Temp
Monthly	- Mean Pp for May - Mean Temp for May - Min Temp for May	- Max Temp for Feb - Mean Temp for May - Min Temp for March	- Max Temp for April - Min Temp for May
Remotely sensed indices <sup>(2)</sup>			
Annual	- Max Albedo - Mean NDVI - Max NDWI	- Mean ET - Max LST - Max NDVI - Max NDWI	- Max ET - Mean NDVI
Monthly	- Mean Albedo for August - Mean ET for March - Mean LST for March - Mean NDVI for June	- Mean EVI for March - Mean ET for May - Mean LST for July - Mean NDWI for July	- Mean EVI for March
Topographic variable <sup>(2)</sup>			
Static	- dev-magnitude - dev-scale - rough-magnitude	- dev-magnitude - dev-scale - rough-magnitude - slope	- elev-stdev - rough-magnitude - slope
Land features			
Static	- Lithology	- Soil types	- Land use/cover
Total number	17	19	13



484 <sup>(1)</sup> ESD: effective soil depth. <sup>(2)</sup> Max: maximum; Min: minimum; dev-magnitude: Max multiscale deviation; dev-scale: Scale of the Max multiscale;  
 485 rough-magnitude: Max multiscale roughness; elev-stdev: Elevation standard deviation; Pp: precipitation; Temp: temperature; NDVI: Normalized  
 486 Difference Vegetation Index; NDWI: Normalized Difference Water Index; ET: Evapotranspiration, LST: Land Surface Temperature.

487 Analysis of residuals revealed  $R^2$  values of 0.68 for SOCc and 0.54 for SOCs in the upper  
 488 30 cm, with lower values in deeper horizons. Table 4 summarizes SOC models' performance at  
 489 different depths. Model validation included computing the CC for both calibration (CCcal) and  
 490 validation (CCval) data and the normalized root mean square error (nRMSE) and normalized mean  
 491 absolute error (nMAE) for validation data. The inconsistency between CCcal and CCval was  
 492 minimal, except in the AutoML model, where the GBM family of algorithms notably contributed to a  
 493 higher CCcal. Based on CCval and nMAE metrics, a decline in accuracy with depth was observed for  
 494 both SOC variables. For instance, at the 0-30 cm depth, CCval was 0.583, and nMAE was 0.441,  
 495 while at the 30-100 cm depth, CCval was 0.351, and nMAE was 0.668 for SOCc. Comparing the  
 496 same depth intervals, SOCc demonstrated greater accuracy than SOCs across all models in the upper  
 497 30 cm, whereas it exhibited the lowest accuracy at 30-100 cm. Although the disparities in nMAE  
 498 among the three predictive ensemble approaches were minimal, there were noticeable differences in  
 499 CC values, indicating a decreasing trend of higher performance in the AutoML, QRF, and MLR  
 500 models.

501 **Table 4.** Average model performance for soil organic carbon concentration (SOCc) and stock (SOCs)  
 502 at different depths.

Parameter	CC <sub>cal</sub>	CC <sub>val</sub>	nRMSE	nMAE	CC <sub>cal</sub>	CC <sub>val</sub>	nRMSE	nMAE
<b>SOC concentration</b>								
<b>Predictive model</b>	<b>0-30 cm</b>				<b>30-100 cm</b>			
	AutoML	0.825	0.583	0.684	0.433	0.629	0.351	1.296
QRF	0.485	0.472	0.733	0.458	0.354	0.307	0.680	0.610
MLR	0.434	0.350	0.775	0.474	0.278	0.227	0.737	0.578
<b>SOC stock</b>								
<b>Predictive model</b>	<b>0-30 cm</b>				<b>Effective soil depth</b>			
	AutoML	0.672	0.417	0.548	0.441	0.563	0.378	0.651
QRF	0.445	0.381	0.545	0.504	0.358	0.295	0.646	0.552
MLR	0.401	0.270	0.628	0.535	0.323	0.232	0.711	0.570

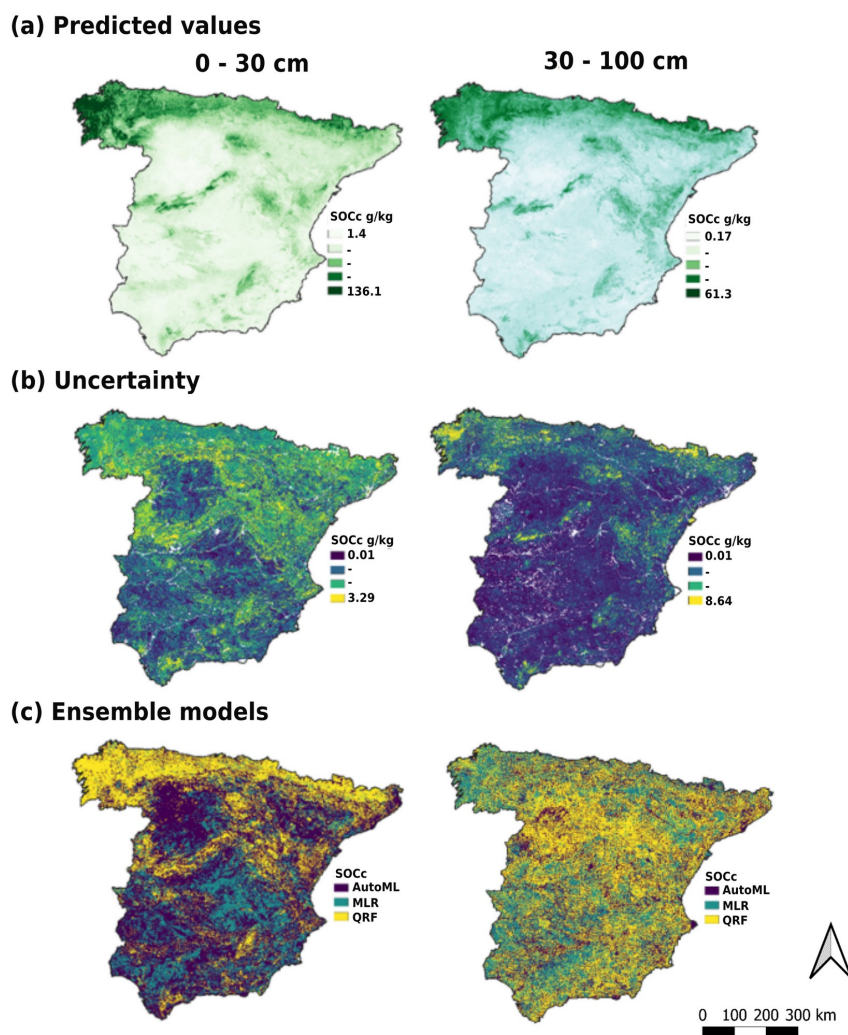
503 Normalized root mean square error (nRMSE) and Normalized mean absolute error (nMAE) in parts per unit.

504 The conditional quantile plots revealed that the predicted SOC values did not encompass  
 505 the entire range of observed values (Fig. S3 in Supplementary material). The highest predicted  
 506 SOCc at 0-30 cm reached a maximum of 136.1 g/kg (compared to the maximum observed value of  
 507 257.95 g/kg), while the highest predicted SOCs value at the effective soil depth (ESD) reached  
 508 38.33 kg/m<sup>2</sup> (compared to the observed value of 93.892 kg/m<sup>2</sup>) (Fig. S3 in Supplementary  
 509 material).

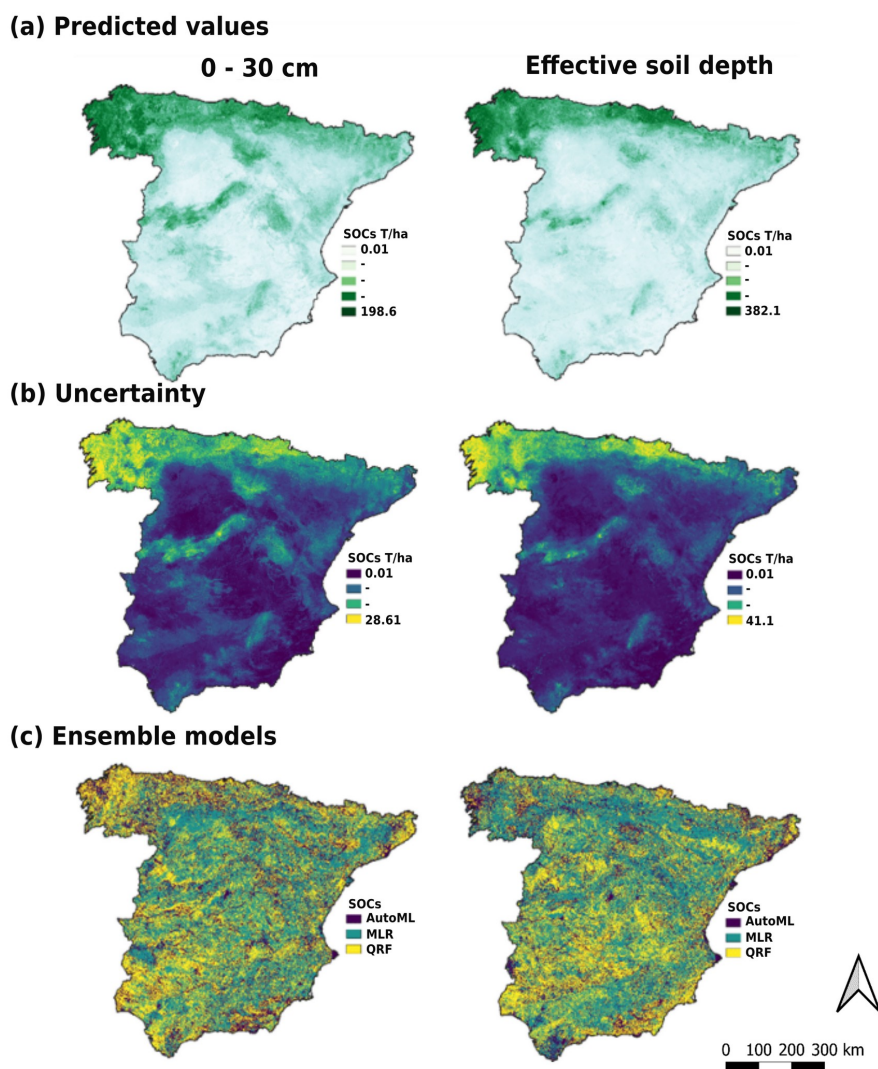


510 Model-based predictive maps of SOCc and SOC<sub>s</sub> in peninsular Spain are shown in Figs. 6  
511 and 7, respectively. Overall, SOCc and SOC<sub>s</sub> exhibited similar and consistent spatial patterns across  
512 depths (Fig. 6a and 7a). Higher SOC predicted values correlated with both climatic and topographic  
513 features. For example, the northwest and north regions, characterized by a humid climate, exhibited  
514 the highest SOC content, gradually decreasing toward the east. Similarly, major mountainous areas,  
515 such as the Central Mountain System, Iberian System, and Subbetic System, also achieved  
516 noticeable SOC predictions. In contrast, low SOC predictions were obtained in extensive  
517 agricultural regions, including the Central Plateau, southwest Guadiana alluvial plain, and the  
518 Guadalquivir depression, along with arid zones in the southeast. The uncertainty around SOC  
519 predictions revealed distinct spatial patterns for SOC predictions, with notable disparities between  
520 the northern and southern regions of peninsular Spain. These disparities positively correlated with  
521 SOC content. North areas with higher SOC predicted values also exhibited greater uncertainty  
522 (Figs. 6b and 7b). Although the spatial distribution of SOC uncertainty remained consistent between  
523 concentration and stock, SOCc exhibited larger areas with higher uncertainty. Regarding standard  
524 depths, regions with higher uncertainty were more prevalent in upper horizons (0-30 cm) than in  
525 deeper horizons.

526 The spatial distribution of the ensemble modeling approach used for predicting SOC  
527 content varied depending on the combination of SOCc and SOC<sub>s</sub> and depth (Figs 6c and 7c).  
528 Regarding SOCc at 0-30 cm, all three models exhibited a uniform spatial distribution, except for  
529 regions with higher SOCc in the north and mountainous areas, where the QRF model was dominant.  
530 In contrast, predictions for SOCc at 30-100 cm showed a reduced contribution from the AutoML  
531 approach, with MLR and QRF models dominating distinct areas. Models' capacity to predict SOC<sub>s</sub>  
532 was spatially consistent at different depths, with no discernible spatial pattern.



533 **Figure 6.** (a) Soil organic carbon concentration (SOCc) maps in peninsular Spain at a 90-meter  
534 pixel resolution. (b) Uncertainty maps are based on the standard deviation of predictions obtained  
535 through the three ensemble modeling approaches. (c) Ensemble algorithm employed for pixel-level  
536 SOC predictions. The visualization used the cumulative pixel count cut method, with a default  
537 range from 2% to 98%.



538 **Figure 7.** (a) Soil organic carbon stock (SOCs) maps in peninsular Spain at a 90-meter pixel  
539 resolution. (b) Uncertainty maps are based on the standard deviation of predictions obtained through  
540 the three ensemble modeling approaches. (c) Ensemble algorithm employed for pixel-level SOC  
541 predictions. The visualization used the cumulative pixel count cut method, with a default range from  
542 2% to 98%.





## 543 **6. Data availability**

544 The soil organic carbon concentration (g/kg) maps for the 0-30 cm and 30-100 cm standard  
545 depths, along with the soil organic carbon stock (tC/ha) maps for the 0-30 cm standard depth and the  
546 effective soil depth, including their associated uncertainties, —all at a 90-meter pixel resolution—  
547 (SOCM90) are freely available at  
548 <https://doi.org/10.6073/pasta/48edac6904eb1aff4c1223d970c050b4> (Durante et al., 2024).

## 549 **7. Further considerations**

550 Our study seeks to enhance the availability and reliability of soil information essential for  
551 informed decision-making regarding SOC management and climate change mitigation strategies. By  
552 integrating disparate soil profile databases and employing advanced ensemble modeling techniques,  
553 we aimed to provide comprehensive and standardized SOC maps for peninsular Spain, facilitating  
554 access to critical soil data at the national scale. As part of the broader effort to enhance soil data  
555 accessibility and usability, our methodology demonstrates the transformation of previously  
556 inaccessible soil information into actionable insights for spatial variability studies and carbon stock  
557 assessments.

558 By establishing a systematic approach to organizing national soil information, we mitigate  
559 potential errors and discrepancies in future data generation processes, ensuring the reliability and  
560 consistency of soil carbon estimates. Furthermore, the enhanced spatial modeling approach of soil  
561 information in peninsular Spain supports ongoing global soil information initiatives, including the  
562 Global Soil Organic Carbon Map, a project of FAO, the Global Soil Partnership, and the  
563 GlobalSoilMap.Net project. It enables informed decision-making regarding land use planning,  
564 agricultural practices, and environmental conservation efforts.

565 The SOCM90 integrated information on more than eight thousand profiles for peninsular  
566 Spain soils. Despite these advancements, it is essential to acknowledge the existence of data gaps in  
567 certain areas and incentivize future soil survey programs to increase sampling efforts in  
568 underrepresented regions. By expanding soil monitoring networks and improving spatial coverage,  
569 the SOCM90 can contribute to more comprehensive assessments of SOC content and inform targeted  
570 soil management strategies.

## 571 **8. Author contribution**

572 Conceptualization: all; Data curation: PD and CO; Formal analysis: PD; Funding acquisition: PD and  
573 CO; Methodology: PD, RV, MG, DA, and CO; Supervision: RV, MG, DA, and CO; Validation: PD,  
574 RV, MG, and CO; Writing – original draft preparation: JMRM, PD, CO; Writing – review & editing:  
575 JMRM, PD, RV, MG, and CO.

## 576 **9. Competing interests**

577 The authors declare that they have no conflict of interest.



## 578 10. Acknowledgements

579 PD acknowledges support from the pre-doctoral grant [DI-15-08093] awarded by the  
580 ‘National Programme for the Promotion of Talent and Its Employability’ of the Ministry of Economy,  
581 Industry, and Competitiveness, which are partially funded by the European Social Fund (ESF) from  
582 the European Commission. JMRM was funded by the University of Almería through the Spanish  
583 Ministry of Universities (María Zambrano Program) [grant number RR\_C\_2021\_09]; the University  
584 of Almería’s programme for research and knowledge transfer [grant number  
585 P\_FORT\_GRUPOS\_2023/26]. RV was supported by the NASA Carbon Monitoring System grant  
586 80NSSC21K0964. MG was funded by UNESCO-IGCP (grant no. 765) and Conahcyt (grant no. CF-  
587 2023-I-1846). DA acknowledges support from the project “Plan Complementario de I+D+i en el área  
588 de Biodiversidad (PCBIO)” funded by the European Union within the framework of the Recovery,  
589 Transformation and Resilience Plan - NextGenerationEU and by the Regional Government of  
590 Andalucía, by the EarthCul project (PID2020-118041GB-I00 Spanish National Research and  
591 Innovation Plan 2020). CO acknowledges support from the projects INTEGRATYON3 (PID2020-  
592 117825GB-C21 and C22), both funded by MCIN/AEI/10.13039/501100011033.

## 593 11. References

- 594 AEMET IPMA: Atlas climático ibérico/Iberian climate atlas, Agencia Estatal de  
595 Meteorología (España) ; Instituto de Meteorología (Portugal), 2011.
- 596 Albaladejo, J., Martínez-Mena, M., Almagro, M., Ruiz-Navarro, A., Ortiz, R., and  
597 Albaladejo Montoro, Juan; Martínez-Mena García, María; Almagro Costa, Mercedes; Ruiz-  
598 Navarro, Ana; Ortiz Silla, R.: Factores de control en la dinámica del Carbono Orgánico de los  
599 suelos de la Región de Murcia, Murcia, 155–158 pp., 2009.
- 600 Alcaraz-Segura, D., Lomba, A., Sousa-Silva, R., Nieto-Lugilde, D., Alves, P., Georges, D.,  
601 Vicente, J. R., and Honrado, J. P.: Potential of satellite-derived ecosystem functional attributes to  
602 anticipate species range shifts, *International Journal of Applied Earth Observation and*  
603 *Geoinformation*, 57, 86–92, <https://doi.org/10.1016/j.jag.2016.12.009>, 2017.
- 604 Alias, L. and Ortiz, R.: Memorias y mapas de suelos de las hojas del MTN a escala  
605 1:100.000, 1986.
- 606 Allen, T. F. H. and Starr, T. B.: *Hierarchy: perspectives for ecological complexity*, University  
607 of Chicago Press, 2019.
- 608 Álvaro-Fuentes, J., López, M. V., Cantero-Martínez, C., and Arrúe, J. L.: Tillage Effects on  
609 Soil Organic Carbon Fractions in Mediterranean Dryland Agroecosystems, *Soil Science Society of*  
610 *America Journal*, 72, 541–547, <https://doi.org/10.2136/sssaj2007.0164>, 2008.
- 611 Amatulli, G., McInerney, D., Sethi, T., Strobl, P., and Domisch, S.: Geomorpho90m,  
612 empirical evaluation and accuracy assessment of global high-resolution geomorphometric layers,  
613 *Sci Data*, 7, 162, <https://doi.org/10.1038/s41597-020-0479-6>, 2020.
- 614 Arrouays, D., Grundy, M. G., Hartemink, A. E., Hempel, J. W., Heuvelink, G. B. M., Hong,  
615 S. Y., Lagacherie, P., Lelyk, G., McBratney, A. B., McKenzie, N. J., Mendonca-Santos, M. d. L.,



- 616 Minasny, B., Montanarella, L., Odeh, I. O. A., Sanchez, P. A., Thompson, J. A., and Zhang, G.-L.:  
617 GlobalSoilMap, 93–134, <https://doi.org/10.1016/B978-0-12-800137-0.00003-0>, 2014a.
- 618 Arrouays D., McBratney A.B., Minasny B., Hempel J.W., Heuvelink G.B.M., MacMillan  
619 R.A., Hartemink A.E., Lagacherie P., and McKenzie N.J.: The GlobalSoilMap project  
620 specifications. In: Arrouays D., McKenzie N., Hempel J., de Forges A.R., McBratney A.B. (Eds.),  
621 GlobalSoilMap: Basis of the Global Spatial Soil Information System, CRC Press, London, 2014b.
- 622 C.E. Arroyo-Cruz, Larson, J., Guevara, M.: A machine learning approach for mapping soil  
623 properties in Mexico using legacy data, climate and terrain covariates at a coarse scale 23–27. In:  
624 Arrouays, D., Savin, I., Leenaars, J., McBratney, A.B. (Eds.) GlobalSoilMap - Digital Soil Mapping  
625 from Country to Globe: Proceedings of the Global Soil Map 2017 Conference, July 4–6, 2017,  
626 Moscow, Russia. CRC Press, London. <https://doi.org/10.1201/9781351239707>, 2017.
- 627 Baillie, I. C.: Soil survey staff 1999, soil taxonomy: a basic system of soil classification for  
628 making and interpreting soil surveys, agricultural handbook 436, Natural Resources Conservation  
629 Service, USDA, Washington DC, USA, pp. 869, 2001.
- 630 Beaudette, D. E., Roudier, P., and O’Geen, A. T.: Algorithms for quantitative pedology: A  
631 toolkit for soil scientists, *Comput Geosci*, 52, 258–268,  
632 <https://doi.org/10.1016/J.CAGEO.2012.10.020>, 2013.
- 633 Beaudette, D. E., Skovlin, J., Brown, A. G., Roudier, P., and Roecker, S. M.: Algorithms for  
634 Quantitative Pedology, in: *Geopedology: An Integration of Geomorphology and Pedology for Soil  
635 and Landscape Studies*, Springer, 201–222, 2023.
- 636 Bischl, B., Lang, M., Kothhoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G.,  
637 and Jones, Z. M.: mlr: Machine Learning in R, *Journal of Machine Learning Research*, 1–5 pp.,  
638 2016.
- 639 Bravo, S., García-Ordiales, E., García-Navarro, F. J., Amorós, J. Á., Pérez-de-los-Reyes, C.,  
640 Jiménez-Ballesta, R., Esbrí, J. M., García-Noguero, E. M., and Higuera, P.: Geochemical  
641 distribution of major and trace elements in agricultural soils of Castilla-La Mancha (central Spain):  
642 finding criteria for baselines and delimiting regional anomalies, *Environmental Science and  
643 Pollution Research*, 26, 3100–3114, <https://doi.org/10.1007/s11356-017-0010-6>, 2019.
- 644 Brus, D., Hengl, T., Heuvelink, G., Kempen, B., Mulder, T. V. L., Olmedo, G. F., Poggio, L.,  
645 Ribeiro, E., Thine Omuto, C., Yigini, Y., and others: Soil organic carbon mapping: GSOC map  
646 cookbook manual, FAO, 2017.
- 647 Calvo de Anta, R., Luís, E., Febrero-Bande, M., Galiñanes, J., Macías, F., Ortíz, R., and  
648 Casás, F.: Soil organic carbon in peninsular Spain: Influence of environmental factors and spatial  
649 distribution, *Geoderma*, 370, 114365, <https://doi.org/10.1016/j.geoderma.2020.114365>, 2020.
- 650 Carlsaw, D. C. and Ropkins, K.: Openair - An R package for air quality data analysis,  
651 *Environmental Modelling & Software*, 27–28, 52–61, <https://doi.org/10.1016/j.envsoft.2011.09.008>,  
652 2012.
- 653 Chen, S., Arrouays, D., Leatitia Mulder, V., Poggio, L., Minasny, B., Roudier, P., Libohova,  
654 Z., Lagacherie, P., Shi, Z., Hannam, J., Meersmans, J., Richer-de-Forges, A. C., and Walter, C.:



- 655 Digital mapping of GlobalSoilMap soil properties at a broad scale: A review, *Geoderma*, 409,  
656 115567, <https://doi.org/10.1016/j.geoderma.2021.115567>, 2022.
- 657 Crowther, T. W., Todd-Brown, K. E. O., Rowe, C. W., Wieder, W. R., Carey, J. C.,  
658 Machmuller, M. B., Snoek, B. L., Fang, S., Zhou, G., Allison, S. D., Blair, J. M., Bridgman, S. D.,  
659 Burton, A. J., Carrillo, Y., Reich, P. B., Clark, J. S., Classen, A. T., Dijkstra, F. A., Elberling, B.,  
660 Emmett, B. A., Estiarte, M., Frey, S. D., Guo, J., Harte, J., Jiang, L., Johnson, B. R., Kröel-Dulay,  
661 G., Larsen, K. S., Laudon, H., Lavallee, J. M., Luo, Y., Lupascu, M., Ma, L. N., Marhan, S.,  
662 Michelsen, A., Mohan, J., Niu, S., Pendall, E., Peñuelas, J., Pfeifer-Meister, L., Poll, C., Reinsch, S.,  
663 Reynolds, L. L., Schmidt, I. K., Sistla, S., Sokol, N. W., Templer, P. H., Treseder, K. K., Welker, J.  
664 M., and Bradford, M. A.: Quantifying global soil carbon losses in response to warming, *Nature*,  
665 540, 104–108, <https://doi.org/10.1038/nature20150>, 2016.
- 666 Dieterich, T. G.: Ensemble Methods in Machine Learning, 1–15, [https://doi.org/10.1007/3-667-540-45014-9\\_1](https://doi.org/10.1007/3-667-540-45014-9_1), 2000.
- 668 Doblas-Miranda, E., Rovira, P., Brotons, L., Martínez-Vilalta, J., Retana, J., Pla, M., and  
669 Vayreda, J.: Soil carbon stocks and their variability across the forests, shrublands and grasslands of  
670 peninsular Spain, *Biogeosciences*, 10, 8353–8361, <https://doi.org/10.5194/bg-10-8353-2013>, 2013.
- 671 Elith, J., Phillips, S. J., Hastie, T., Dudík, M., Chee, Y. E., and Yates, C. J.: A statistical  
672 explanation of MaxEnt for ecologists, *Divers Distrib*, 17, 43–57, [https://doi.org/10.1111/j.1472-4642.2010.00725.x](https://doi.org/10.1111/j.1472-673-4642.2010.00725.x), 2011.
- 674 FAO and ITPS: Global Soil Organic Carbon Map (GSOCmap), 2018.
- 675 Fielding, A. H. and Bell, J. F.: A review of methods for the assessment of prediction errors in  
676 conservation presence/absence models, *Environ Conserv*, 24, 38–49, 1997.
- 677 Friedman, J. H. and Stuetzle, W.: Projection Pursuit Regression, *J Am Stat Assoc*, 76, 817,  
678 <https://doi.org/10.2307/2287576>, 1981.
- 679 Fryda T; Erin LeDell, N. G. S. A. A. F. A. C. C. T. K. T. N. P. A. M. K. M. M.: Package  
680 “h2o” Title R Interface for the “H2O” Scalable Machine Learning Platform,  
681 <https://doi.org/10.32614/CRAN.package.h2o>, 2024.
- 682 Gavilán-Acuña, G., Olmedo, G.F., Mena-Quijada, P., Guevara, M., Barría-Knopf, B., Watt,  
683 M.S.: Reducing the Uncertainty of Radiata Pine Site Index Maps Using an Spatial Ensemble of  
684 Machine Learning Models. *Forests* 12, 77. <https://doi.org/10.3390/f12010077>, 2021.
- 685 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B.:  
686 Bayesian Data Analysis, Chapman and Hall/CRC, <https://doi.org/10.1201/b16018>, 2013.
- 687 Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., and Moore, R.: Google  
688 Earth Engine: Planetary-scale geospatial analysis for everyone, *Remote Sens Environ*, 202, 18–27,  
689 <https://doi.org/10.1016/j.rse.2017.06.031>, 2017.
- 690 Gray, J. M., Bishop, T. F. A., and Yang, X.: Pragmatic models for the prediction and digital  
691 mapping of soil properties in eastern Australia, *Soil Research*, 53, 24,  
692 <https://doi.org/10.1071/SR13306>, 2015.



- 693 Greenwell, B. M. and Boehmke, B. C.: Variable Importance Plots - An Introduction to the  
694 vip Package, *R J*, 12, 343–366, <https://doi.org/10.32614/RJ-2020-013>, 2020.
- 695 Gregorutti, B., Michel, B., and Saint-Pierre, P.: Correlation and variable importance in  
696 random forests, *Stat Comput*, 27, 659–678, <https://doi.org/10.1007/s11222-016-9646-1>, 2017.
- 697 Guevara, M., Olmedo, G. F., Stell, E., Yigini, Y., Aguilar Duarte, Y., Arellano Hernández, C.,  
698 Arévalo, G. E., Arroyo-Cruz, C. E., Bolivar, A., Bunning, S., Bustamante Cañas, N., Cruz-  
699 Gaistardo, C. O., Davila, F., Dell Acqua, M., Encina, A., Figueredo Tacona, H., Fontes, F.,  
700 Hernández Herrera, J. A., Ibelle Navarro, A. R., Loayza, V., Manueles, A. M., Mendoza Jara, F.,  
701 Olivera, C., Osorio Hermosilla, R., Pereira, G., Prieto, P., Ramos, I. A., Rey Brina, J. C., Rivera, R.,  
702 Rodríguez-Rodríguez, J., Roopnarine, R., Rosales Ibarra, A., Rosales Riveiro, K. A., Schulz, G. A.,  
703 Spence, A., Vasques, G. M., Vargas, R. R., and Vargas, R.: No silver bullet for digital soil mapping:  
704 country-specific soil organic carbon estimates across Latin America, *Soil*, 4, 173–193,  
705 <https://doi.org/10.5194/soil-4-173-2018>, 2018.
- 706 Harden, J. W., Hugelius, G., Ahlström, A., Blankinship, J. C., Bond-Lamberty, B.,  
707 Lawrence, C. R., Loisel, J., Malhotra, A., Jackson, R. B., Ogle, S., Phillips, C., Ryals, R., Todd-  
708 Brown, K., Vargas, R., Vergara, S. E., Cotrufo, M. F., Keiluweit, M., Heckman, K. A., Crow, S. E.,  
709 Silver, W. L., DeLonge, M., Nave, L. E., Bond-Lamberty, B., Lawrence, C. R., Loisel, J., Malhotra,  
710 A., Jackson, R. B., Ogle, S., Phillips, C., Ryals, R., Todd-Brown, K., Vargas, R., Vergara, S. E.,  
711 Cotrufo, M. F., Keiluweit, M., Heckman, K. A., Crow, S. E., Silver, W. L., DeLonge, M., and Nave,  
712 L. E.: Networking our science to characterize the state, vulnerabilities, and management  
713 opportunities of soil organic matter, *Glob Chang Biol*, 24, e705–e718,  
714 <https://doi.org/10.1111/gcb.13896>, 2018.
- 715 Hengl, T., Jesus, J.M. de, Heuvelink, G.B.M., Gonzalez, M.R., Kilibarda, M., Blagotić, A.,  
716 Shangquan, W., Wright, M.N., Geng, X., Bauer-Marschallinger, B., Guevara, M.A., Vargas, R.,  
717 MacMillan, R.A., Batjes, N.H., Leenaars, J.G.B., Ribeiro, E., Wheeler, I., Mantel, S., and Kempen,  
718 B.: SoilGrids250m: Global gridded soil information based on machine learning. *PLOS ONE* 12,  
719 e0169748. <https://doi.org/10.1371/journal.pone.0169748>, 2017.
- 720 Hengl, T., Kempen, B., Heuvelink, G., and Malone, B.: Package “GSIF” - Global Soil  
721 Information Facilities, 2020.
- 722 Hijmans, R. J.: Geographic Data Analysis and Modeling [R package raster version 3.6-28],  
723 2024.
- 724 IGME, 1995: Mapa Geológico de España a escala 1:1.000.000  
725 [http://info.igme.es/cartografiadigital/geologica/Geologicos1MMapa.aspx?](http://info.igme.es/cartografiadigital/geologica/Geologicos1MMapa.aspx?Id=Geologico1000_(1994))  
726 [Id=Geologico1000\\_\(1994\)](http://info.igme.es/cartografiadigital/geologica/Geologicos1MMapa.aspx?Id=Geologico1000_(1994)).
- 727 IGN, 2012: Mapa de suelos de España: Escala 1:1.000.000  
728 [http://info.igme.es/cartografiadigital/geologica/Geologicos1MMapa.aspx?](http://info.igme.es/cartografiadigital/geologica/Geologicos1MMapa.aspx?Id=Geologico1000_(1994))  
729 [Id=Geologico1000\\_\(1994\)](http://info.igme.es/cartografiadigital/geologica/Geologicos1MMapa.aspx?Id=Geologico1000_(1994)).
- 730 Jenny, H.: Factors of soil formation; a sytem of quantitative pedology, 1941.



- 731 Jobbágy, E. G., Jackson, R. B., Processes, B., and Change, G.: The vertical distribution of  
732 soil organic carbon and its relation to climate and vegetation, *Ecological Applications*, 10, 423–436,  
733 [https://doi.org/https://doi.org/10.1890/1051-0761\(2000\)010\[0423:TVDOS0\]2.0.CO;2](https://doi.org/https://doi.org/10.1890/1051-0761(2000)010[0423:TVDOS0]2.0.CO;2), 2000.
- 734 Keitt, T., Bivand, R., Pebesma, E., Rowlingson, B., and ‘Rgdal,’ P.: Bindings for the  
735 geospatial data abstraction library, 2012.
- 736 Kottek, M., Grieser, J., Beck, C., Rudolf, B., and Rubel, F.: World Map of the Köppen-  
737 Geiger climate classification updated, *Meteorologische Zeitschrift*, 15, 259–263,  
738 <https://doi.org/10.1127/0941-2948/2006/0130>, 2006.
- 739 De la Rosa, D., Mayol, F., Moreno, D., Rosales, A., Castillo, V., Moreno, F., Cabrera, F.,  
740 Colomer, J. C., Antoine, J., Masui, S., Brinkman, R., Horn, R., and Prange, N.: SEIS.NET: Sistema  
741 Español de Información de Suelos en Internet, *Edafología*, 8, 45–56, 2001.
- 742 Llorente, M., Rovira, P., Merino, A., Rubio, A., Turrión, M., Badalía, D., Romanya, J., and  
743 González, J. C. J. A.: The CARBOSOL Database: a georeferenced soil profile analytical database  
744 for Spain, <https://doi.org/10.1594/PANGAEA.884517>, 2018.
- 745 MAPA: Anuario dde Estadística de España 2020, Ministerio de Agricultura, Pesca y  
746 Alimentación, 2021.
- 747 McBratney, A. B. B., Mendonça Santos, M. L. L., and Minasny, B.: On digital soil mapping,  
748 *Geoderma*, 117, 3–52, [https://doi.org/10.1016/S0016-7061\(03\)00223-4](https://doi.org/10.1016/S0016-7061(03)00223-4), 2003.
- 749 McElreath, R.: *Statistical Rethinking*, Chapman and Hall/CRC,  
750 <https://doi.org/10.1201/9781315372495>, 2018.
- 751 McElreath, R.: *rethinking: Statistical Rethinking book package*, 2016, R package version, 1,  
752 2020.
- 753 Meinshausen, N.: Quantile Regression Forests, *Journal of Machine Learning Research*, 7,  
754 983–999, 2006.
- 755 Mulder, V. L., Lacoste, M., Richer-de-Forges, A. C., Martin, M. P., and Arrouays, D.:  
756 National versus global modelling the 3D distribution of soil organic carbon in mainland France,  
757 *Geoderma*, 263, 16–34, <https://doi.org/10.1016/j.geoderma.2015.08.035>, 2016.
- 758 Muñoz-Rojas, M., Jordán, A., Zavala, L. M., De la Rosa, D., Abd-Elmabod, S. K., and  
759 Anaya-Romero, M.: Organic carbon stocks in Mediterranean soil types under different land uses  
760 (Southern Spain), *Solid Earth*, 3, 375–386, <https://doi.org/10.5194/se-3-375-2012>, 2012.
- 761 Ninyerola, M., Pons, X., and Roure, J.: *Atlas climático digital de la Península Ibérica:*  
762 *metodología y aplicaciones en bioclimatología y geobotánica*, edited by: Universitat Autònoma de  
763 Barcelona, Barcelona, 2005.
- 764 Padarian, J. and McBratney, A. B.: A new model for intra- and inter-institutional soil data  
765 sharing, *SOIL*, 6, 89–94, <https://doi.org/10.5194/soil-6-89-2020>, 2020.
- 766 Pebesma, E. and Bivand, R. S.: *Classes and Methods for Spatial Data: the sp Package*, 2005.
- 767 Pebesma, E. J.: Multivariable geostatistics in S: the gstat package, *Comput Geosci*, 30, 683–  
768 691, 2004.





- 769 Peterson, B. G., Carl, P., Boudt, K., Bennett, R., Ulrich, J., Zivot, E., Lestel, M., Balkissoon,  
770 K., and Wuertz, D.: PerformanceAnalytics: Econometric tools for performance and risk analysis, R  
771 package version, 1, 2014.
- 772 Phillips, S. J., Anderson, R. P., and Schapire, R. E.: Maximum entropy modeling of species  
773 geographic distributions, *Ecol Modell*, 190, 231–259,  
774 <https://doi.org/10.1016/j.ecolmodel.2005.03.026>, 2006.
- 775 Poeplau, C., Vos, C., and Don, A.: Soil organic carbon stocks are systematically  
776 overestimated by misuse of the parameters bulk density and rock fragment content, *SOIL*, 3, 61–66,  
777 <https://doi.org/10.5194/soil-3-61-2017>, 2017.
- 778 Poggio, L., de Sousa, L. M., Batjes, N. H., Heuvelink, G. B. M., Kempen, B., Ribeiro, E.,  
779 and Rossiter, D.: SoilGrids 2.0: producing soil information for the globe with quantified spatial  
780 uncertainty, *SOIL*, 7, 217–240, <https://doi.org/10.5194/soil-7-217-2021>, 2021.
- 781 R Foundation for Statistical Computing: R: A Language and Environment for Statistical  
782 Computing., 2022.
- 783 Rodríguez Martín, J. A., Álvaro-Fuentes, J., Gonzalo, J., Gil, C., Ramos-Miras, J. J., Grau  
784 Corbí, J. M., and Boluda, R.: Assessment of the soil organic carbon stock in Spain, *Geoderma*, 264,  
785 117–125, <https://doi.org/10.1016/j.geoderma.2015.10.010>, 2016.
- 786 Rosell, R. A., Gasparoni, J. C., and Galantini, J. A.: Soil organic matter evaluation, in:  
787 Assessment methods for soil carbon, Lewis Publishers Boca Raton, 311–322, 2001.
- 788 Saatchi, S. S., Houghton, R. A., Dos Santos Alvalá, R. C., Soares, J. V., and Yu, Y.:  
789 Distribution of aboveground live biomass in the Amazon basin, *Glob Chang Biol*, 13, 816–837,  
790 <https://doi.org/10.1111/j.1365-2486.2007.01323.x>, 2007.
- 791 Scharlemann, J. P., Tanner, E. V., Hiederer, R., and Kapos, V.: Global soil carbon:  
792 understanding and managing the largest terrestrial carbon pool, *Carbon Manag*, 5, 81–91,  
793 <https://doi.org/10.4155/cmt.13.77>, 2014.
- 794 Searle, R., McBratney, A., Grundy, M., Kidd, D., Malone, B., Arrouays, D., Stockman, U.,  
795 Zund, P., Wilson, P., Wilford, J., Van Gool, D., Triantafilis, J., Thomas, M., Stower, L., Slater, B.,  
796 Robinson, N., Ringrose-Voase, A., Padarian, J., Payne, J., Orton, T., Odgers, N., O’Brien, L.,  
797 Minasny, B., Bennett, J. M., Liddicoat, C., Jones, E., Holmes, K., Harms, B., Gray, J., Bui, E., and  
798 Andrews, K.: Digital soil mapping and assessment for Australia and beyond: A propitious future,  
799 *Geoderma Regional*, 24, e00359, <https://doi.org/10.1016/j.geodrs.2021.e00359>, 2021.
- 800 Serrano, F. S.: Deformaciones Recientes en el Centro de la Península Ibérica, 2000.
- 801 Shangguan, W., Dai, Y., Duan, Q., Liu, B., and Yuan, H.: A global soil data set for earth  
802 system modeling, *J Adv Model Earth Syst*, 6, 249–263, <https://doi.org/10.1002/2013MS000293>,  
803 2014.
- 804 Shangguan, W., Hengl, T., Mendes de Jesus, J., Yuan, H., and Dai, Y.: Mapping the global  
805 depth to bedrock for land surface modeling, *J Adv Model Earth Syst*, 9, 65–88,  
806 <https://doi.org/10.1002/2016MS000686>, 2017.



- 807 Smith, P., Soussana, J. F., Angers, D., Schipper, L., Chenu, C., Rasse, D. P., Batjes, N. H.,  
808 Egmond, F., McNeill, S., Kuhnert, M., Arias-Navarro, C., Olesen, J. E., Chirinda, N., Fornara, D.,  
809 Wollenberg, E., Álvaro-Fuentes, J., Sanz-Cobena, A., Klumpp, K., van Egmond, F., McNeill, S.,  
810 Kuhnert, M., Arias-Navarro, C., Olesen, J. E., Chirinda, N., Fornara, D., Wollenberg, E., Álvaro-  
811 Fuentes, J., Sanz-Cobena, A., and Klumpp, K.: How to measure, report and verify soil carbon  
812 change to realize the potential of soil carbon sequestration for atmospheric greenhouse gas removal,  
813 *Glob Chang Biol*, 26, 219–241, <https://doi.org/10.1111/gcb.14815>, 2020.
- 814 Stockmann, U., Padarian, J., McBratney, A., Minasny, B., de Brogniez, D., Montanarella, L.,  
815 Hong, S. Y., Rawlins, B. G., and Field, D. J.: Global soil organic carbon assessment, *Glob Food*  
816 *Sec*, 6, 9–16, <https://doi.org/10.1016/j.gfs.2015.07.001>, 2015.
- 817 USDA: Module 3. USDA Textural Classification Study Guide, in: Soil mechanics level 1,  
818 National Employee Development Staff, Soil Conservation Service, U.S. Department of Agriculture.,  
819 Washington DC, 1987.
- 820 Vargas, R., Alcaraz-Segura, D., Birdsey, R., Brunzell, N. A., Cruz-Gaistardo, C. O., de Jong,  
821 B., Etchevers, J., Guevara, M., Hayes, D. J., Johnson, K., Loescher, H. W., Paz, F., Ryu, Y.,  
822 Sanchez-Mejia, Z., and Toledo-Gutierrez, K. P.: Enhancing interoperability to facilitate  
823 implementation of REDD+: case study of Mexico, *Carbon Manag*, 8, 57–65,  
824 <https://doi.org/10.1080/17583004.2017.1285177>, 2017.
- 825 Varón-Ramírez, V.M., Araujo-Carrillo, G.A., Guevara Santamaría, M.A.: Colombian soil  
826 texture: building a spatial ensemble model. *Earth System Science Data* 14, 4719–4741,  
827 <https://doi.org/10.5194/essd-14-4719-2022>, 2022.
- 828 Villarreal, S., Guevara, M., Alcaraz-Segura, D., Brunzell, N. A., Hayes, D., Loescher, H. W.,  
829 and Vargas, R.: Ecosystem functional diversity and the representativeness of environmental  
830 networks across the conterminous United States, *Agric For Meteorol*, 262, 423–433,  
831 <https://doi.org/10.1016/j.agrformet.2018.07.016>, 2018.
- 832 Wang, B., Waters, C., Orgill, S., Cowie, A., Clark, A., Li Liu, D., Simpson, M., McGowen,  
833 I., and Sides, T.: Estimating soil organic carbon stocks using different modelling techniques in the  
834 semi-arid rangelands of eastern Australia, *Ecol Indic*, 88, 425–438,  
835 <https://doi.org/10.1016/j.ecolind.2018.01.049>, 2018a.
- 836 Wang, B., Waters, C., Orgill, S., Gray, J., Cowie, A., Clark, A., and Liu, D. L.: High  
837 resolution mapping of soil organic carbon stocks using remote sensing variables in the semi-arid  
838 rangelands of eastern Australia, *Science of the Total Environment*, 630, 367–378,  
839 <https://doi.org/10.1016/j.scitotenv.2018.02.204>, 2018b.
- 840 Wilks, D. S.: Statistical Methods in the Atmospheric Sciences. 4th Edition. (International  
841 Geophysics), Elsevier, <https://doi.org/10.1016/C2017-0-03921-6>, 2019.
- 842 WRB-IUSS: World Reference Base for Soil Resources. 2014, FAO, Rome, Italy, 1–191 pp.,  
843 2014.
- 844 Yan, X. and Su, X.: Linear Regression Analysis: Theory and Computing, World Scientific,  
845 2009.



- 846 Yigini, Y., Olmedo, G. F., Reiter, S., Baritz, R., Viatkin, K., and Vargas, R.: Soil organic  
847 carbon mapping: Cookbook 2nd Edition, FAO, Rome, 220 pp., 2018.
- 848 Zhang, C. and Ma, Y.: Ensemble Machine Learning: Methods and applications, edited by:  
849 Zhang, C. and Ma, Y., Springer US, Boston, MA, VIII, 332 pp., [https://doi.org/10.1007/978-1-4419-](https://doi.org/10.1007/978-1-4419-9326-7)  
850 9326-7, 2012.