

1                   **A SUPER-LEARNER MACHINE LEARNING MODEL FOR A GLOBAL**  
2                   **PREDICTION OF COMPRESSION INDEX IN CLAYS**

3                   Esteban Díaz<sup>1\*</sup>, Giovanni Spagnoli <sup>2</sup>

4                   <sup>1</sup> Departamento de Ingeniería Civil, Escuela Politécnica Superior, Universidad de Alicante, P.O.

5                   Box 99, E-03080 Alicante, Spain, [esteban.diaz@ua.es](mailto:esteban.diaz@ua.es)

6                   <sup>2</sup> DMT GmbH & Co. KG, Am TÜV 1, 45307 Essen, Germany, [spagnoli\\_giovanni@yahoo.de](mailto:spagnoli_giovanni@yahoo.de),

7                   [giovanni.spagnoli@dm-group.com](mailto:giovanni.spagnoli@dm-group.com); ORCID: 0000-0002-1866-4345

8                   \* Corresponding author.

9  
10                   **Abstract**

11                   Settlement of structures is determined by the stiffness of the soil where they are built.  
12                   Compression index ( $c_c$ ) quantifies the compressibility of the soil and is a key parameter in the  
13                   design of geotechnical structures. To predict the value of  $c_c$  in clay soils, a global database of  
14                   more than 1000 data points was collected and analysed. Liquid limit, plasticity index, natural  
15                   water content, and initial void ratio were considered as predictor variables. A super-learner  
16                   machine learning model was developed to predict  $c_c$  from these variables. The model  
17                   demonstrated a reasonable predictive performance and was subsequently integrated into an online  
18                   tool. Additionally, four symbolic regression expressions were obtained to relate  $c_c$  with some of  
19                   the input variables when not all data are available, providing simple and practical alternatives for  
20                    $c_c$ , estimation. This study provided two major contributions: (1) the non-local nature of the data  
21                   expands the scope and generalizability of the findings, and (2) the availability of the proposed  
22                   algorithm through an online application ensures its accessibility for geotechnical engineers,  
23                   enhancing the work's practical applicability and intrinsic value.

24  
25                   **Keywords:** machine learning; compression index; liquid limit; plasticity index; natural water  
26                   content; initial void ratio; clay.

## 29 **1. Introduction**

30 Compression index ( $c_c$ ), which is the slope of the linear section of the  $e$ - $\log\sigma'$  plot (and is  
31 dimensionless), and the coefficient of consolidation ( $c_v$ ), primarily define the compressibility  
32 properties of fine-grained soils (Craig 2004).  $c_c$  is a crucial parameter for characterizing soil  
33 compressibility and deformation behaviour under load. It is commonly used in geotechnical  
34 engineering to estimate the settlement and deformation of soil structures, such as foundations and  
35 embankments (Das 2021). Soils with a higher  $c_c$  tend to be more compressible and deformable  
36 under load, whereas soils with a lower  $c_c$  are less compressible and have a greater capacity to  
37 withstand deformation when they are loaded. The  $c_c$  value in soils is obtained through the  
38 oedometer test, which is relatively time-consuming and results in higher costs when compared to  
39 standard index tests. Since Atterberg limits initially characterize soils, these state parameters have  
40 been used to establish correlations with various other engineering properties of soils. Furthermore,  
41 plasticity is influenced by the electrochemical behaviour of clay minerals (Carter and Bentley  
42 1991) in the same way as  $c_c$  (Onyejekwe et al. 2015). For this reason, several attempts were made  
43 in the past to correlate basic geotechnical properties with  $c_c$ . Numerous authors provided linear  
44 equations relating  $c_c$  to the liquid limit (LL) of soils (e.g. Azzouz et al. 1976; Bowles 1979; Park  
45 and Lee 2011; Sridharan and Nagaraj 2000; Terzaghi et al. 1967; Tsuchida 1991). On the other  
46 hand, plasticity index (PI), was also correlated with  $c_c$  (Sridharan and Nagaraj 2000; Wroth and  
47 Wood 1978). Additionally, many correlations based on a linear relationship with natural water  
48 content ( $w$ ) were proposed (e.g. Azzouz et al. 1976; Koppula 1981; Rendon-Herrero 1980). The  
49 relationship with initial void ratio ( $e_0$ ) was examined by Nishida (1956), Hough (1957) and  
50 (Bowles 1979), among others. Other studies included more than one index property in the  
51 estimation of  $c_c$ , such as  $w$  and LL (Koppula 1981) or  $e_0$  and LL (Al-Khafaji and Andersland  
52 1992). However, when these correlations were tested with new data, they exhibited significant  
53 scatter, with deviations reaching up to 30% (Spagnoli and Shimobe 2020), suggesting a lack of  
54 universally applicable validity. They are applicable within specific limits and should be restricted  
55 to the soil type or location where they were validated (Verbrugge and Schroeder 2018). Using

56 these correlations in different conditions may lead to unsatisfactory outcomes (Onyejekwe et al.  
57 2016). To address the limits of classic regression approaches in geotechnical engineering, the  
58 application of machine learning (ML) algorithms have been extensively developed,  
59 demonstrating improved performance for predicting several soil engineering properties compared  
60 to traditional statistical methods (e.g. Bardhan et al. 2023; Bardhan et al. 2021; Dam Nguyen et  
61 al. 2022; Díaz and Tomás 2021; Salvatore et al. 2022; Singh et al. 2023; Trong et al. 2021).  
62 However, it is important to be aware of the limitations and uncertainties associated with ML  
63 approaches before applying them to real-world geotechnical engineering projects. Numerous  
64 studies (Baghbani et al. 2022; Zhang et al. 2023; Zhang et al. 2022) have exposed these  
65 limitations, which are primarily: a) the scarcity of high-quality data, b) the difficulty in  
66 interpreting the models, and c) the lack of generalization. Regarding data availability,  
67 geotechnical data can be costly and often incomplete or uncertain. This can lead to ML models  
68 that are not as accurate or reliable as desired. Another limitation of ML is the interpretability of  
69 the models (i.e., black box). ML models are often complex and nonlinear, making it challenging  
70 to understand the relationship between input data and output predictions. Finally, due to the  
71 inherent heterogeneity and spatial variability of soil deposits, it is difficult for empirical models  
72 trained on limited datasets to reliably extrapolate beyond the geographic scope represented by  
73 those data. Therefore, while ML shows great potential to complement traditional approaches in  
74 geotechnical engineering, it is crucial to address these limitations before the implementation of  
75 any ML model.

76 Recently, numerous studies have employed ML algorithms to predict  $c_c$  from some parameters  
77 related to this property showing promising results (e.g. Desai et al. 2009; Kalantary and Kordnaeij  
78 2012; Kumar and Rani 2011; Nesamatha and Arumairaj 2015; Samui et al. 2012). Alam et al.  
79 (2014) compiled a database of 125 clay samples, which included  $w$ , LL,  $e_0$ , and PI as input  
80 variables and created an Artificial Neural Network (ANN) model to predict  $c_c$ . Kumar and Rani  
81 (2011) also used an ANN to predict  $c_c$ , considering 41 samples with the following input variables:  
82 fine contents, LL, PI, maximum dry density, and optimal moisture content. Park and Lee (2011)  
83 developed an artificial neural network (ANN) model by utilizing 947 consolidation tests

84 performed on soil samples gathered from 67 construction sites in the Republic of Korea. They  
85 considered as input variables,  $w$ , LL, PI,  $e_0$ , specific gravity of soil particles ( $G_s$ ), and weight  
86 percentage of sand, silt and clay. Benbouras et al. (2019) developed an ANN model with 373  
87 oedometer test samples to correlate  $c_c$  with wet density,  $w$ ,  $e_0$ , fine content, LL, PI, and soil type.  
88 The dataset utilized in this research comprised samples gathered from various projects executed  
89 in the city of Algiers (Algeria). Zhang et al. (2021) used a Random Forest algorithm that utilized  
90 a database with 311 samples, encompassing three input variables (LL, PI,  $e_0$ ). Asteris et al. (2022)  
91 introduced extreme learning machine models that applied Manta ray foraging optimization to  
92 predict  $c_c$  from void ratio at liquid limit, LL and PI. It should be noted that the void ratio at LL is  
93 a parameter that is not usually available in the design phases of geotechnical projects, whereas  
94  $e_0$ , is more frequently encountered. Long et al. (2023) established a relationship between  $c_c$  and  
95  $w$ , LL, PI,  $e_0$ , and  $G_s$  using Tree-Based Techniques from 391 samples from Northern Iran.  
96 However, all these approaches have, to a greater or lesser extent, a local character, or a relatively  
97 small number of samples. All the studies presented either rely on a limited dataset, or the collected  
98 samples have a local nature (i.e., they come from the same area or country), or they do not have  
99 a direct application because they are based on ML algorithms. Many engineers lack knowledge  
100 of the programming language or specific software required to utilize the proposed algorithms.  
101 Therefore, the main objective of this work is threefold: 1) to create an algorithm based on a dataset  
102 with a large number of samples, 2) to ensure these samples are collected worldwide to eliminate  
103 any local bias, and 3) to design the proposed algorithm in a way that any geotechnical engineer  
104 can easily use it. The present paper took advantage of ensemble learners, to develop and validate  
105 an accurate prediction model for forecasting the compression index of clay soils for a large  
106 number of data (1008) collected worldwide. Subsequently, this model was deployed to a user-  
107 friendly web application. This paper is structured as follows. A brief summary of the experimental  
108 database is given in Section 2. The machine learning process performed is presented in Section  
109 3. Later, in Section 4, the main results obtained are studied and interpreted by evaluating the  
110 super-learner machine learning model. Section 5 describes the online tool deployed. An analysis

111 of the dataset using Symbolic Regression is included in Section 6. The final section (Section 7)  
 112 ends the paper and presents the main conclusions obtained.

113

114 **2. Database description and analysis.**

115 The collection of the dataset is the first step in the building of a machine learning model. In the  
 116 present work, the experimental database of  $c_c$  (1008 samples) comprised 913 samples from  
 117 research papers (Alhaji et al. 2017; Benbouras et al. 2019; Kalantary and Kordnaeij 2012; LCPC  
 118 1977; McCabe et al. 2014; Mitachi and Ono 1985; Widodo and Ibrahim 2012), and 95 samples  
 119 from authors' own data. This database has samples of different countries as Nigeria, Ireland,  
 120 Spain, Iran, Indonesia, France, Algeria, Bangladesh, among others. In the same way, the database  
 121 has soils with low plasticity to very high plasticity. In fact, LL ranges from 17.1% to 199.0% and  
 122 PI from 2.0% to 82.0%. On the other hand, the database included soils with very high  
 123 compressibility and soils with low compressibility according to  $e_0$  values (ranging from 0.279 to  
 124 7.114). This fact is also corroborated, by examining the values of  $c_c$  varying these between 0.013  
 125 and 2.2. Finally, the  $w$  values, also varies widely ranging from 8% to 244.1%. Based on all the  
 126 above, the database includes a significant amount of data and a wide range of them, making it  
 127 suitable for a reliable study. The final dataset can be found in Table S1 (accessible online), which  
 128 includes details such as references, soil type, mineralogy, and origin when available. Furthermore,  
 129 the dataset includes both the actual  $c_c$  values and those obtained by the super-learner ML model.  
 130 The main statistics of the compiled database are shown in Table 1.

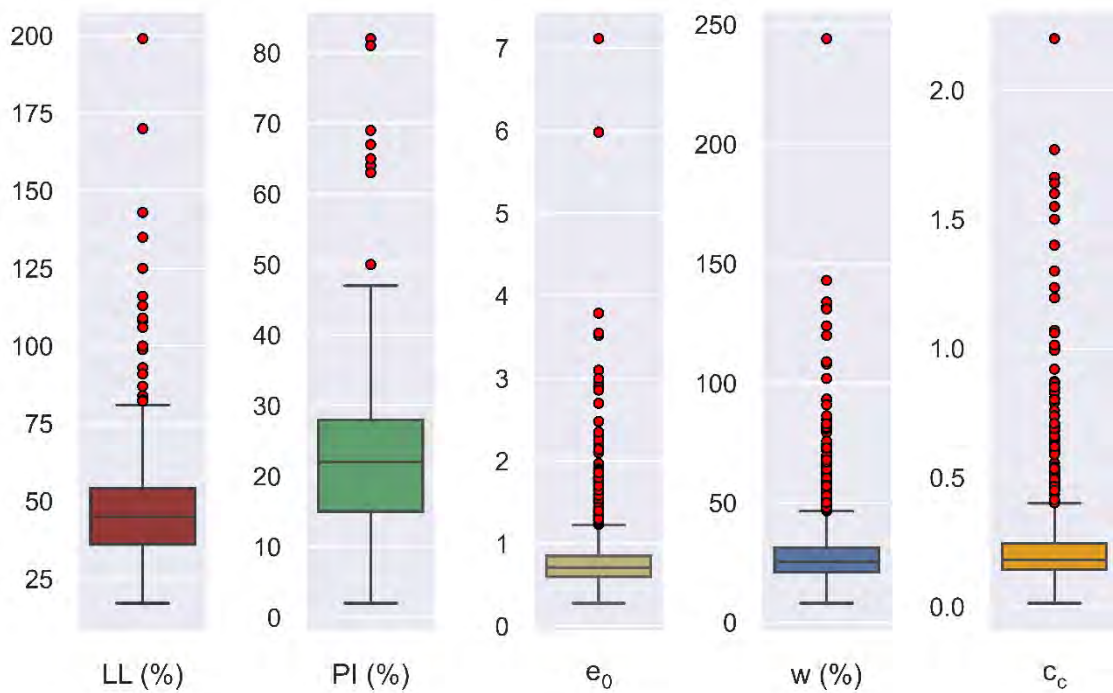
131

	LL (%)	PI (%)	$e_0$	$w$ (%)	$c_c$
Samples	1008	947	1008	1008	1008
Mean	46.49	21.78	0.817	29.73	0.236
Standard deviation	15.90	9.83	0.467	17.06	0.210
Minimum	17.10	2.00	0.279	8.00	0.013
25 <sup>th</sup> percentile	36.00	15.07	0.598	21.00	0.143
Median	44.85	22.03	0.712	25.40	0.180
75 <sup>th</sup> percentile	54.14	28.05	0.852	31.20	0.246
Maximum	199.00	82.00	7.114	244.10	2.200

132 *Table 1. Descriptive statistics of the data analysed.*

133

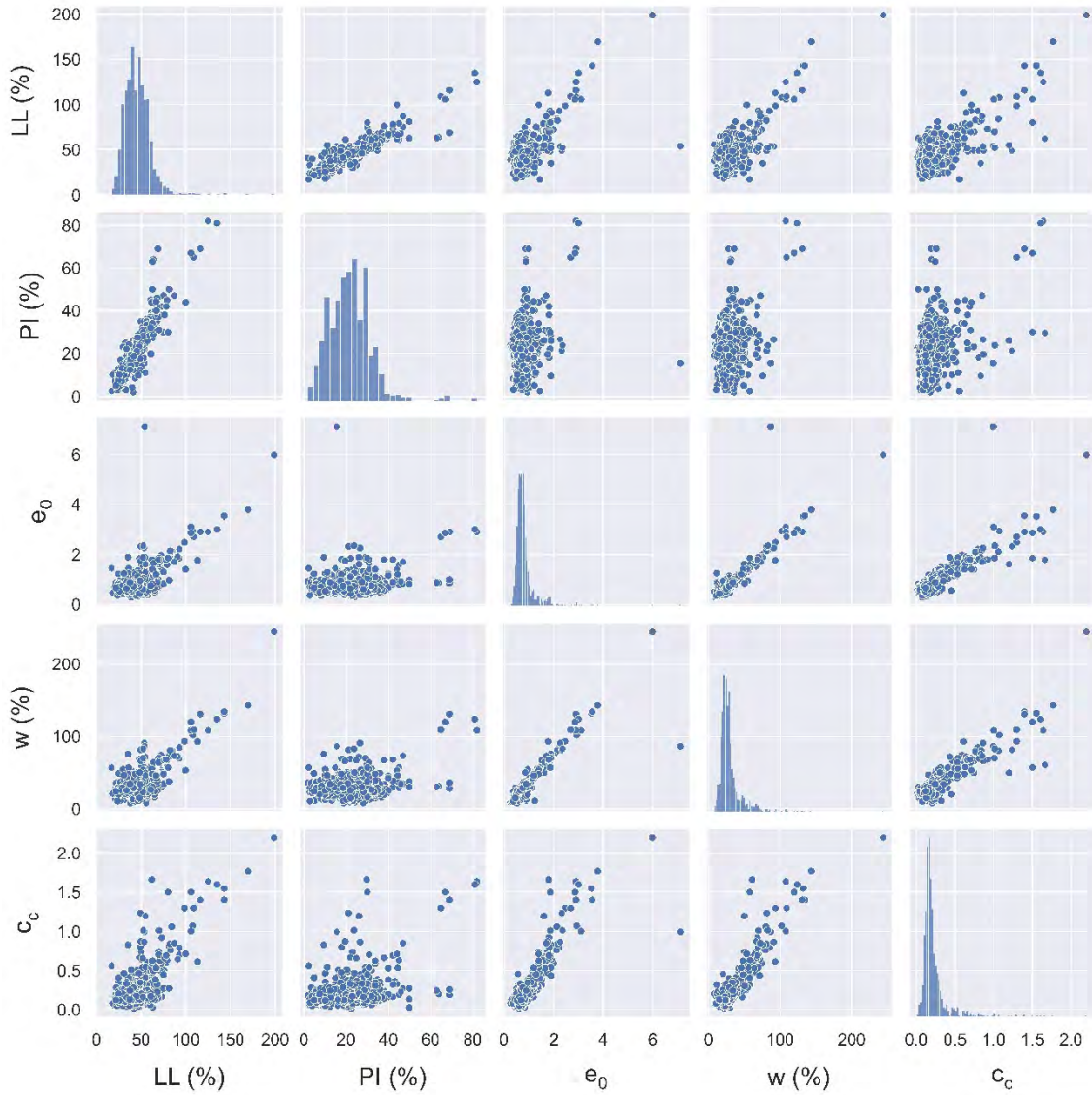
134 Figure 1 shows box plots for LL, PI,  $e_0$ ,  $w$  and  $c_c$  of the database. Box plots are used to visualize  
 135 data dispersion, which was split into quartiles. The method is used to detect outliers (if any), data  
 136 symmetry, dispersion, and skewness (Reagan and Kiemele 2008). The box in a box plot shows  
 137 the interquartile range (IQR), with the bottom and top of the box representing the 25<sup>th</sup> and 75<sup>th</sup>  
 138 percentiles, respectively. The whiskers extend to the final data value inside the inner fence, which  
 139 is 1.5 times the IQR from the box's edge. The height of the box represents the interquartile range.  
 140 Outliers are defined as data points extending to  $3 \times \text{IQR}$  (Reagan and Kiemele 2008). Some points  
 141 in Figure 1 are identified as outliers, mainly corresponding to high values of the variables. The  
 142 analysis of these box plots reveals the wide range of variation and high dispersion of the variables  
 143 studied clearly related to the worldwide character of the compiled database. For a better  
 144 interpretation of Figure 1, the data included in Table 1 can be consulted.  
 145



146  
 147 *Figure 1. Box plots of the variables considered.*  
 148

149 The scatter plots of the considered parameters are depicted in Figure 2, to show a descriptive  
 150 overview of the data distribution. These plots indicate a positive relationship between  $c_c$  and the  
 151 rest of the variables. This fact implies, to a greater or lesser extent, that an increase in the

152 considered input variables tends to proportionally increase  $c_c$ . This relationship is clearer in  
 153 variables as  $e_0$  and  $w$  and more diffuse in LL and PI. Finally, a strong relationship between  $e_0$   
 154 and  $w$  must also be noted.  
 155



156

157

*Figure 2. Scatter plots and distribution histograms of the variables.*

158

159 These linear trends can be quantified numerically using the Pearson correlation coefficient ( $r$ ) and  
 160 a correlation matrix. A correlation matrix is a table showing the correlation values, which measure  
 161 the degree of linear relationship between each pair of variables. Correlation values can be between  
 162 -1 and +1. If the two variables tend to increase or decrease at the same time, the correlation value

163 is positive. In Table 2 is shown the correlation matrix providing an overview of the Pearson  
 164 correlation coefficient. From the analysis of the table, it can be verified that the same relationships  
 165 established previously in a visual manner are the ones that obtain the highest correlation values.  
 166 Indeed, a strong and positive correlation is shown between  $c_c$  and  $e_0$  ( $r=0.87$ ) and  $w$  ( $r=0.89$ ), and  
 167 a somewhat less strong and also positive correlation between  $c_c$  and LL ( $r=0.66$ ). Between  $c_c$  and  
 168 PI, the correlation is low and positive ( $r=0.34$ ). On the other hand, there is a strong and positive  
 169 correlation between  $e_0$  and  $w$  ( $r=0.92$ ).  
 170

	<b>LL</b>	<b>PI</b>	$e_0$	$w$	$c_c$
<b>LL</b>	1.00				
<b>PI</b>	0.92	1.00			
$e_0$	0.61	0.23	1.00		
$w$	0.64	0.28	0.92	1.00	
$c_c$	<b>0.66</b>	<b>0.34</b>	<b>0.87</b>	<b>0.89</b>	<b>1.00</b>

171 *Table 2. Correlation matrix of the variables considered. Note that the values indicated*  
 172 *correspond to the Pearson correlation coefficient (r).*

173

174 The equations obtained by linear regression for the cases with the highest value of Pearson  
 175 correlation coefficient are (Equations 1 to 3):

176

177  $c_c = 0.393e_0 - 0.0842$  (1)

178  $c_c = 0.011w(\%) - 0.0905$  (2)

179  $e_0 = 0.0252w(\%) + 0.0687$  (3)

180



181 The coefficient of determination ( $R^2$ ) of each of the three equations above presented is 0.76, 0.80,  
182 0.85, respectively. Equation 3 is very interesting, as the  $e_0$  value can be estimated by  $w$ ,  
183 considering that  $e_0$  relates the soil structure with its geologic history (Onyejekwe et al. 2015).

184

### 185 **3. Machine learning procedure.**

#### 186 **3.1. Model selection process.**

187 According to the analysis of the dataset carried out in the previous section, PI has 61 data less  
188 than the rest of the variables. To match the number with the rest of the variables, a data imputation  
189 technique was used. This is a technique widely used in ML algorithms for dealing with missing  
190 values and it has been used in geotechnical issues satisfactorily (e.g. Aydın et al. 2023; Díaz et al.  
191 2023). For the imputation of values, a multivariate feature imputation algorithm has been chosen  
192 (Little and Rubin 2019; Van Buuren and Oudshoorn 2000). This technique uses the information  
193 of all of the available features in order to estimate the missing value of one variable by considering  
194 samples which have a similar situation in terms of all of the features in the dataset. After imputing  
195 the missing data, an outlier detection analysis was conducted using the one-class SVM algorithm  
196 developed by Schölkopf et al. (1999) and successfully employed in similar works (e.g. Díaz et al.  
197 2023). This algorithm is a method designed to identify outliers and anomalies within a dataset,  
198 utilizing the principles of traditional Support Vector Machines (SVM). The core idea behind SVM  
199 is to identify a hyperplane that maximizes the separation between different classes within the  
200 dataset. Once this hyperplane is established, new data points can be classified based on their  
201 position relative to the hyperplane. In the case of the one-class SVM, there is only one class, and  
202 it defines the hyperplane for normal data points while classifying data points located outside the  
203 hyperplane as outliers. Upon applying this algorithm, 9 outliers were identified. These 9 data  
204 points were thoroughly examined, but no anomalous values in their properties were observed.  
205 Therefore, due to their limited number, it was decided not to remove them from the dataset. With  
206 the final dataset including 1008 records of each variable, a machine learning algorithm selection  
207 process was carried out using the k-fold cross-validation technique (with  $k=10$ ). The results of  
208 this selection process are shown in Table 3, ordering the selected algorithms, from best to worst

209 performance considering four statistical indicators, mean absolute percentage error (MAPE),  
 210 coefficient of determination ( $R^2$ ), root mean squared error (RMSE), and mean absolute error  
 211 (MAE). Three algorithms, all of them based on decision trees, outperform over the rest: Extra  
 212 Trees Regressor (Geurts et al. 2006), Random Forest Regressor (Ho 1995), and Gradient Boosting  
 213 Regressor (Friedman 2001).  
 214

Model	RMSE	MAPE	MAE	$R^2$
Extra Trees Regressor	0.0767	0.2521	0.0466	0.87
Random Forest Regressor	0.0776	0.2534	0.0473	0.86
Gradient Boosting Regressor	0.0781	0.2529	0.0478	0.86
Huber Regressor	0.1050	0.2766	0.0559	0.75
Light Gradient Boosting Machine	0.1059	0.2831	0.0576	0.74
K Neighbors Regressor	0.1000	0.2919	0.0578	0.75
AdaBoost Regressor	0.0876	0.3427	0.0581	0.81
Ridge Regression	0.1084	0.3106	0.0610	0.74
Linear Regression	0.1117	0.3106	0.0613	0.72
Least Angle Regression	0.1117	0.3106	0.0613	0.72
Bayesian Ridge	0.1113	0.3124	0.0616	0.72
Decision Tree Regressor	0.0966	0.3296	0.0622	0.79
Elastic Net	0.1040	0.3307	0.0625	0.74
Orthogonal Matching Pursuit	0.1034	0.3471	0.0654	0.74
Lasso Regression	0.1170	0.3479	0.0668	0.68
Lasso Least Angle Regression	0.1170	0.3479	0.0668	0.68
Passive Aggressive Regressor	0.1232	0.3562	0.0728	0.63

215 *Table 3. Performance of base machine learning models obtained by k-fold cross validation.*

216  
 217 This process was carried out with the variables without normalization. But it should be noted that  
 218 the same algorithm selection process was also performed normalizing the variables using the min-  
 219 max method, which scales each variable individually between zero and one. The results of this  
 220 process normalizing the variables, were exactly the same in the three models with the best  
 221 performance.

222

### 223 **3.2. Model development.**

224 To ensure a proper generalization of the algorithms, it is good practice to assess their performance  
 225 on unknown data. For this purpose, the dataset was divided into two groups (training and test)  
 226 with an 80/20 partition. Subsequently, the three algorithms with the best performance were

227 subjected to a tuning process of their hyperparameters to maximize their performance. For this  
 228 purpose, the particle swarm optimisation (Kennedy and Eberhart 1995) was used. This is a  
 229 computational technique which optimises problems by iteratively improving candidate solutions  
 230 (aka particles). The results of this optimization are shown in Table 4, in terms of  $R^2$ , RMSE, and  
 231 MAE and with the result in both in the training and in the test set.

232

<b>Model</b>	<b>Set</b>	<b><math>R^2</math></b>	<b>MAE</b>	<b>RMSE</b>
Extra Trees	Training	0.92	0.034	0.059
Regressor	Test	0.92	0.040	0.055
Random Forest	Training	0.92	0.039	0.061
Regressor	Test	0.92	0.041	0.055
Gradient Boosting	Training	0.93	0.038	0.056
Regressor	Test	0.91	0.042	0.058

233 *Table 4. Summary of performance metrics of the trained machine learning algorithms.*

234

235 With the tuned algorithms, they were ensembled to improve their predictive capacity. Among the  
 236 existing approaches, the super-learner algorithm (Laan et al. 2007) was chosen. The super-learner  
 237 algorithm is a type of ensemble method that applies stacked generalization to k-fold cross-  
 238 validation. It combines multiple prediction models (base learners) by assigning different weights  
 239 to these models to find their optimal combination and produce a single best prediction function.  
 240 Thus, the predictions of the base learners are used to train a regression model (meta learner) that  
 241 assigns relative weights to the predictions of each base-model. In this case, and after a process of  
 242 trial and error, the best results were obtained using the tuned Extra Trees Regressor and Gradient  
 243 Boosting Regressor algorithms as base learners, and the tuned Random Forest Regressor  
 244 algorithm as meta learner. It must be indicated that the meta learner was trained on the base  
 245 models' predictions as well as the original training data. The performance metrics resulting in the  
 246 super-learner machine learning model, are shown in Table 5.

247

<b>Model</b>	<b>Set</b>	<b><math>R^2</math></b>	<b>MAE</b>	<b>RMSE</b>
Super-Learner	Training	0.93	0.034	0.057
	Test	0.93	0.039	0.053

248 *Table 5. Summary of performance metrics of the super-learner machine learning model.*

249 As can be seen the super-learner machine learning model improved slightly the predictions of the  
250 best of the three selected algorithms (Table 4).

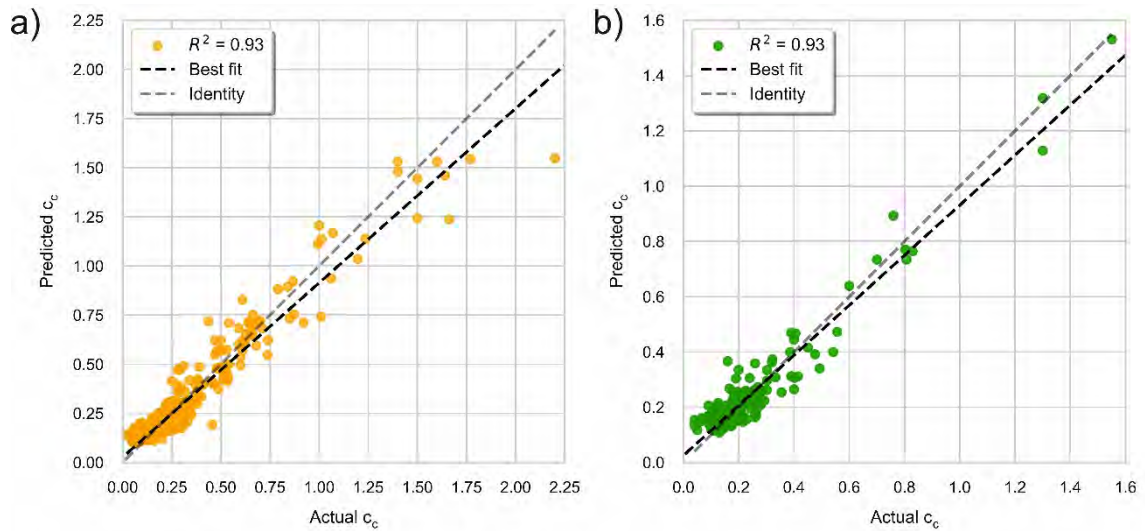
251

## 252 **4. Results.**

### 253 **4.1. Analysis of the super-learner ML model for predicting $c_c$ .**

254 The prediction accuracy of the super-learner machine learning model is evaluated using training  
255 and test datasets. Figure 3 shows the scatter plots for the actual values of  $c_c$  (x-axis) versus the  
256 predicted values of  $c_c$  (y-axis). This figure shows that the vast majority of data are located close  
257 to the no error line, indicating a good agreement between the predicted and measured values.  
258 Additionally, in Figure 3 are included the values of  $R^2$  on both the training and test datasets. These  
259 values (0.93 in both sets) are indicative of a high predictive capacity, and since they are the same,  
260 it is assured a correct behaviour of the model in unseen data, discarding overfitting issues. The  
261 predictive performance of the selected model was also evaluated based on the anteriorly defined  
262 performance metrics (MAE and RMSE). In the training dataset, MAE and RMSE values of 0.034  
263 and 0.057 were obtained, respectively. In the test set, these values were similar (MAE of 0.039  
264 and RMSE of 0.053). These metrics are summarized in Table 5 (previous section) for both the  
265 train and test datasets. Additionally, the a20-index (Apostolopoulou et al. 2020; Asteris et al.  
266 2021a; Asteris et al. 2021b) was calculated. The a20-index offers the advantage of having a clear  
267 engineering interpretation, indicating the percentage of samples that meet the predicted values  
268 within a  $\pm 20\%$  deviation from the experimental ones. The values obtained in the training and test  
269 sets were 73.45% and 71.27%, respectively, showing a similar and reasonable performance for  
270 the prediction of  $c_c$ .

271



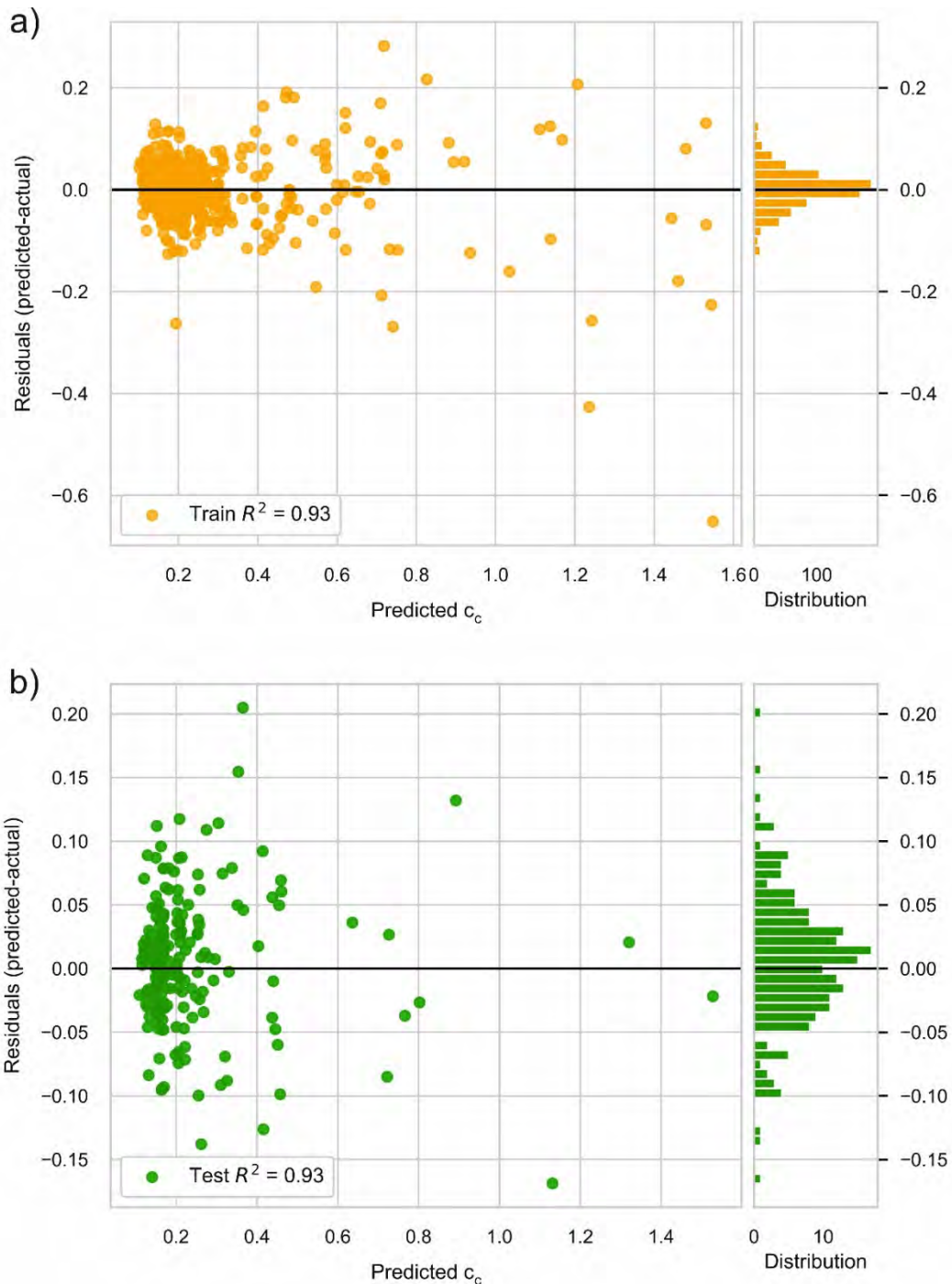
272

273 *Figure 3. Scatter plot showing graphical performance of the super-learner machine learning*  
 274 *model in (a) training set, and (b) test set. The identity line (i.e., no error line) and the best fit*  
 275 *line were included in both figures.*

276

277 In order to evaluate the results of the adopted model, a residuals study was performed. The  
 278 residuals were defined as the difference between the predicted and actual value obtained by the  
 279 super-learner model. In Figure 4 are shown the residuals for both the training and test datasets  
 280 with two different visualizations, a scatter plot, and a histogram distribution. An inspection of this  
 281 figure shows that for the two datasets considered, the residuals are concentrated around zero and  
 282 have Gaussian distributions (i.e., residuals are normally distributed) and, thus, there are not many  
 283 outliers. These facts are indicative of a robust algorithm with a great generalization ability and do  
 284 not reveal significant issues in the predictions.

285



286

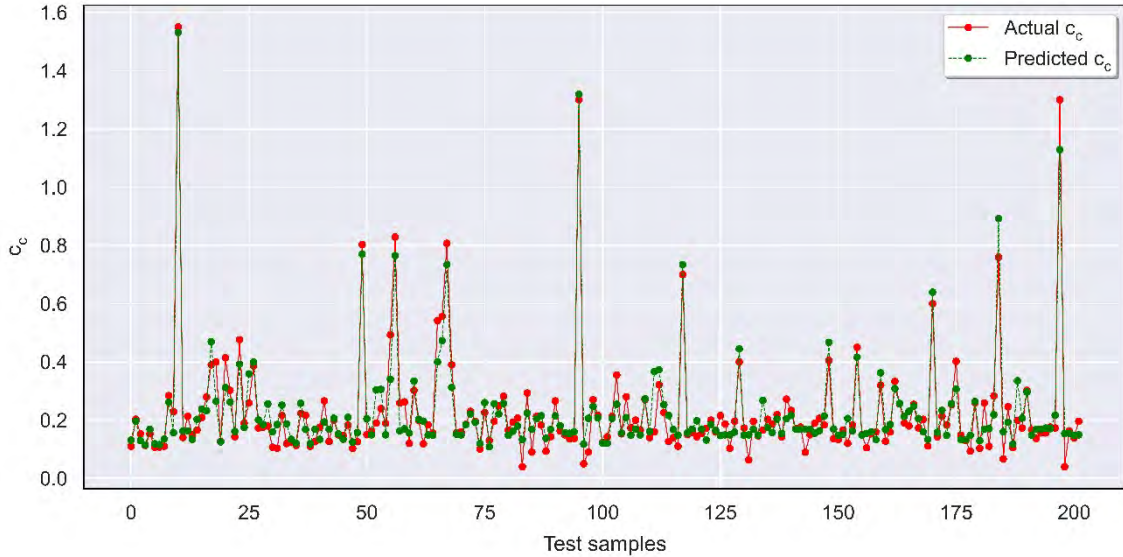
287 *Figure 4. Residuals scatter plot and histogram distribution in (a) training set and (b) test set of*  
 288 *the super-learner machine learning model.*

289

290 Figure 5 compares the experimental and predicted  $c_c$  values on the test set (unseen data). The  
 291 prediction performance of the testing set, according to what was previously discussed, was very  
 292 good. Figure 5 shows that the super-learner machine learning model, can accurately capture the

293 evolution of the actual values, the prediction results are consistent, and the difference between the  
 294 minimum and maximum predicted values is relatively small. No large deviations are observed in  
 295 a general manner, affecting only isolated samples.

296



297

298 *Figure 5. Performance plot of super-learner machine learning model prediction results vs*  
 299 *actual values of  $c_c$  in the test set.*

300

301 **4.2. Sensitivity analysis.**

302 Subsequently, a sensitivity analysis of the parameters involved in the super-learner ML model  
 303 was conducted. This analysis employed Sobol's method (Sobol 1990), a technique for assessing  
 304 the significance of input parameters in computational models. Sobol's method evaluates their  
 305 impact on the output by quantifying their contribution to output variance. Within this method, the  
 306 First-order sensitivity indices (S1) measure the influence of each variable in isolation while  
 307 holding all other variables constant. Total-order sensitivity indices (ST) encompass not only the  
 308 individual effects of each variable but also their interactions with other variables on the output.  
 309 The results of this analysis are collected in Table 6, showing that  $e_0$  was the most influential  
 310 variable, with a ST of 0.839, implying that it had the greatest effect on the variability of the  
 311 model's output. Following  $e_0$ ,  $w$  emerged as the second most important factor in this analysis (ST  
 312 = 0.122). Next, LL possessed a moderate ST, with a value of 0.066. Finally, PI exhibited the

313 lowest ST, with a value of 0.018. The values of ST and S1 for each variable were similar,  
 314 suggesting that the majority of the influence of the variables on the model's output was attributed  
 315 to their individual effects rather than complex interactions among them.

316

Variable	ST	S1
LL	0.066	0.064
PI	0.018	0.016
$e_0$	0.839	0.802
$w$	0.122	0.080

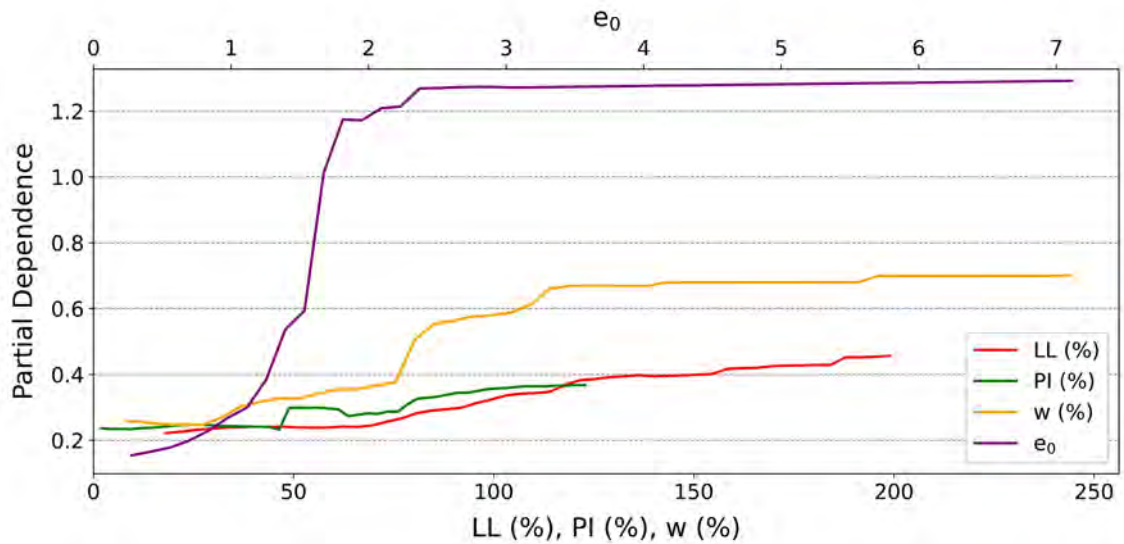
317 *Table 6. Sensitivity analysis of variables using Sobol's method comparing the Total (ST) and*  
 318 *First-Order (S1) Sensitivity Indices for input variables.*

319

320 After establishing the significance of the variables, an analysis was conducted to understand the  
 321 impact of variations in these inputs on the overall results. Initially, partial dependence plots for  
 322 single input features were analysed to visualize the relationship between individual input variables  
 323 and the output variable in a regression model. These plots are instrumental in understanding how  
 324 a single input variable influences the output while holding other variables constant. The results of  
 325 this analysis are displayed in Figure 6 in terms of partial dependence (average predicted effect).  
 326 For the variables LL and PI, an increasing trend was observed, with some fluctuations, indicating  
 327 a positive and near linear relationship with  $c_c$ . In the case of variable  $e_0$ , a very strong growth  
 328 trend was noted until  $e_0=2.5$ , after which the impact stabilizes. This suggests that within this range  
 329 of values (up to  $e_0=2.5$ ), small variations in  $e_0$  will cause a significant impact on the value of  $c_c$ ,  
 330 clearly indicating that the relationship is not linear. As for  $w$ , a similar pattern to  $e_0$  was observed,  
 331 but with a lesser impact on the value of  $c_c$ , with the effect becoming stable around values of 110%.

332



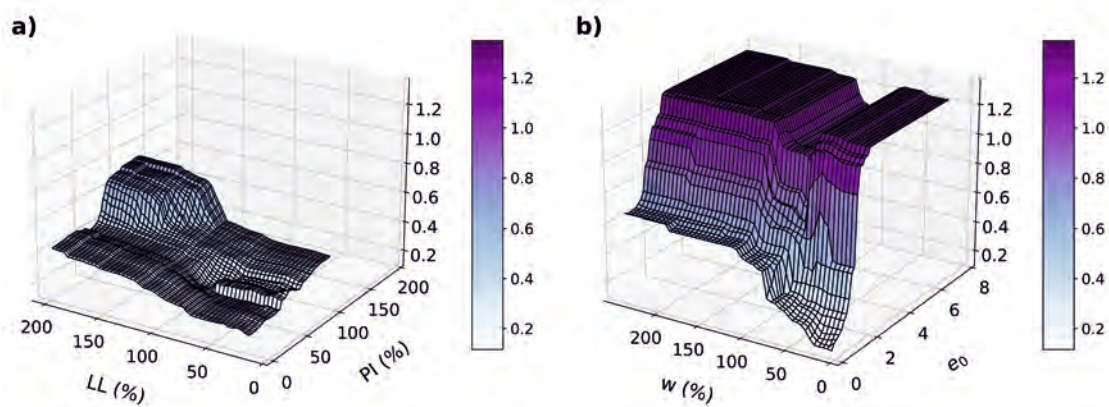


333

334 *Figure 6. Partial dependence plots for single input features in terms of the average predicted*  
 335 *effect. Note that the figure has a dual x-axis, with the upper one for  $e_0$  and the lower one for the*  
 336 *variables LL, PI and w.*

337

338 Finally, partial dependence plots for pairs of input features were computed. These plots are  
 339 designed to illustrate the joint influence of input feature pairs on the model's predictions, offering  
 340 insights into their collective impact on the output variable. They assist in elucidating complex  
 341 interactions between input features and the model's response. These graphs depicting the pairs of  
 342 variables LL and PI, as well as  $w$  and  $e_0$ , are presented in Figure 7. The graphs showing the  
 343 remaining variable pairs are included in Figure S1. From this analysis, it was established that most  
 344 of the impact on the model's output can be attributed to the individual effects of the variables,  
 345 rather than complex interactions among them. In the LL and PI pair (Figure 7a), as well as in the  
 346  $w$  and  $e_0$  pair (Figure 7b), some interaction was observed, although not excessively significant.  
 347 These conclusions were consistent with the Sobol analysis.



348

349 *Figure 7. Partial dependence plots for pairs of input features. a) LL - PI and b) w -  $e_0$ .*

350

351 In conclusion, the analyses conducted highlighted that the most significant variable in the model  
 352 was  $e_0$ , which also caused the most substantial changes in the predicted value of  $c_c$ , both  
 353 individually and in terms of interactions between pairs of variables.

354

### 355 **4.3. Overfitting analysis**

356 One of the common challenges encountered in ML models is overfitting. This occurs when a  
 357 model performs exceptionally well in replicating the data used for its development and training.  
 358 However, when applied to input parameter values outside those used during development and  
 359 training, the model may generate highly unusual predictions. To analyse the overfitting of a ML  
 360 model, there are several techniques, highlighting the works of Asteris et al. (2019) or Armaghani  
 361 and Asteris (2021), where the final model is checked with datasets in which the inputs gradually  
 362 increase in value as they are tests made for a subsequent analysis. In this study, it is not like this  
 363 since the inputs do not increase gradually as they were tests conducted on soils where the values  
 364 of the input variables do not follow a gradual pattern. Thus, alternative techniques were chosen  
 365 to enable the identification of overfitting in a model. To this end, the difference in model  
 366 performance between the training set and the test set was initially assessed. In this case, a very  
 367 small difference was observed between the training and test errors (Table 5), indicating that there  
 368 was no evidence of overfitting. Alternatively, analysing the residuals can provide insights into  
 369 potential overfitting in the model. Specifically, the residuals of the train and test sets were

370 analysed based on the following criteria: a) size: if the test residuals are significantly larger than  
 371 the train residuals, it is a sign of overfitting; b) variance: if similar variances exist in train and test,  
 372 it is indicative that the model generalizes well, not evidencing overfitting; and c) correlation: if  
 373 the residuals are correlated, it is a sign that the model is learning patterns that are not relevant for  
 374 predicting the dependent variable, which is a sign of overfitting. In this case, the size of the  
 375 residuals and their variance were given by the mean and standard deviation, respectively. The  
 376 correlation of the residuals was carried out through the Durbin-Watson test, which is a method to  
 377 detect the presence of autocorrelation among the residuals in regression analysis. The value of the  
 378 Durbin-Watson statistic varies between 0 and 4. A value close to 2 suggests that there is no  
 379 autocorrelation; values less than 2 indicate positive autocorrelation; and values greater than 2  
 380 indicate negative autocorrelation. The summary of this analysis is shown in Table 7.

381

<b>Set</b>	<b>Mean</b>	<b>Standard deviation</b>	<b>Durbin-Watson statistic</b>
Training	0.0003	0.05730	1.95
Test	-0.0038	0.05301	1.87

382

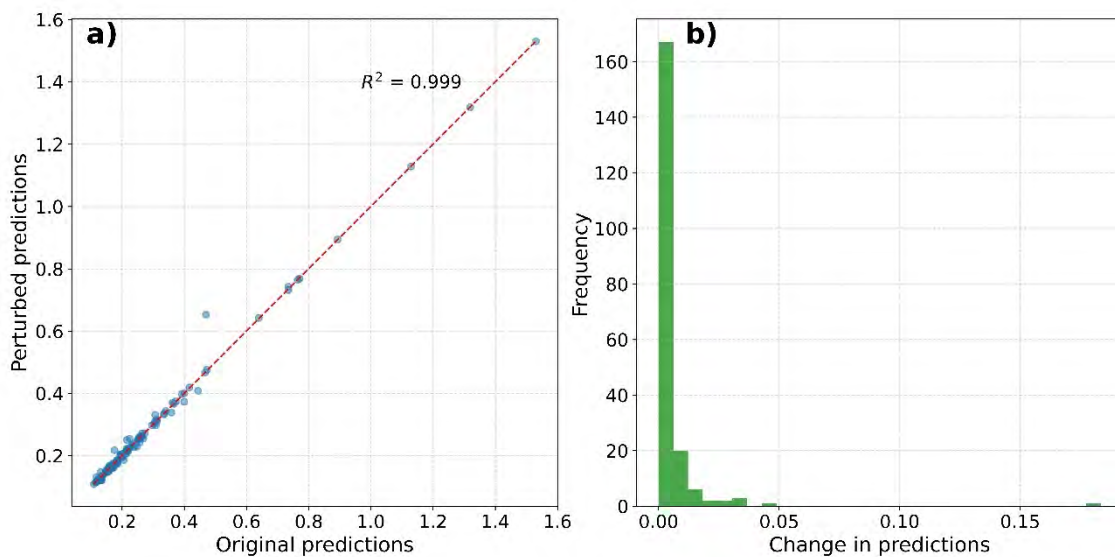
*Table 7. Summary of residual metrics in the training and test sets.*

383

384 From the analysis of the results, it was evident that the mean of the residuals is very close to zero  
 385 for both sets, which is good as it suggested that there is no systematic bias in the predictions. The  
 386 standard deviations were also comparable, indicating that the dispersion of errors was similar in  
 387 both sets. The Durbin-Watson values were close to 2, implying that the residuals exhibited  
 388 independence. According to the results obtained in the residual analysis, there is no evidence of  
 389 overfitting from this perspective.

390 Although none of the conducted analyses indicated signs of overfitting in the model, an additional  
 391 analysis was undertaken by perturbing the input values of the test set. In this analysis, a random  
 392 perturbation of  $\pm 1\%$  was applied to the test set data (including all input variables) to simulate  
 393 natural variability and assess the robustness of the model and its propensity for overfitting.  
 394 Subsequently, the super-learner ML model was tasked with making predictions on both the  
 395 original test set and the perturbed test set to evaluate how the model reacted to disturbances in the

396 input data. A model that is excessively sensitive to small perturbations in the data might be  
 397 overfitting, meaning it has learned to adjust to the particularities and noise of the training set,  
 398 rather than capturing general trends. The results of this analysis are presented in Figure 8, where  
 399 the scatter plot revealed an almost perfect correlation ( $R^2$  of 0.999) between the original and  
 400 perturbed predictions. The histogram showed that the average changes in predictions due to  
 401 perturbations were minimal, with a mean of 0.0036, suggesting significant stability of the model  
 402 against minor variations in the data. As a conclusion of this analysis, the model was highly robust  
 403 and showed low sensitivity to disturbances in the input data. Therefore, the three conducted  
 404 analyses concluded that the super-learner ML model exhibited strong generalization abilities,  
 405 capturing underlying trends rather than fitting to specific noises or peculiarities of the training  
 406 data.  
 407



408  
 409 *Figure 8. Analysis of model robustness to input perturbation. a) Scatter plot comparing original*  
 410 *and perturbed model predictions, and b) Histogram of absolute changes in model predictions*  
 411 *due to perturbation.*

#### 413 **4.4. Comparison of the super-learner ML model with geotechnical correlations.**

414 In order to assess the performance of the proposed model, a comparative analysis was conducted  
 415 using the database compiled for this study. This analysis involved nine empirical correlations for

416 estimating  $c_c$  with a global scope. The results are depicted in Figure S2 as a scatter plot,  
417 showcasing the comparison between each correlation's predictions and the actual values from the  
418 comprehensive dataset. Additionally, in this figure, the formulas of the tested correlations are  
419 included. The super-learner ML model proposed in this work clearly achieved the best  
420 performance overall, with the highest  $R^2$  score, lowest RMSE, and lowest MAE (Figure S2a). In  
421 contrast, the Sridharan and Nagaraj (2000) correlation (Figure S2b) performed very poorly, with  
422 a negative  $R^2$  score (-0.79), high RMSE (0.281), and high MAE (0.257), clearly showing it did  
423 not work well for this dataset. Other correlations such as Azzouz et al. (1976) (Figure S2e),  
424 Rendon Rendon-Herrero (1980) (Figure S2g), and Koppula (1981) (Figure S2f) achieved  
425 moderate results, with  $R^2$  scores around 0.7, RMSE values around 0.10, and MAE values around  
426 0.06. The Wroth and Wood (1978) (Figure S2c), Bowles (1979) (Figure S2h), and Al-Khafaji and  
427 Andersland (1992) (Figure S2i) correlations demonstrated low performance, characterized by low  
428  $R^2$  scores around 0.20, RMSE values around 0.19, and high MAE values. Additionally, the  
429 Koppula (1981) (Figure S2j) correlation also had poor outcomes, with negative  $R^2$  scores, high  
430 RMSE, and high MAE. It is noteworthy that some correlations, such as Sridharan and Nagaraj  
431 (2000) (Figure S2b and d), Koppula (1981) (Figure S2f and j), and in less extent, Wroth and Wood  
432 (1978) (Figure S2c), consistently predicted values significantly above the actual ones.  
433 Conversely, correlations like Al-Khafaji and Andersland (1992) (Figure S2i) and Bowles (1979)  
434 (Figure S2h) tended to predict consistently below the actual values. Furthermore, the four  
435 expressions obtained through Symbolic Regression in this study (Equations 4 to 7) also exhibited  
436 superior performance metrics compared to the analysed correlations. In summary, the model  
437 proposed in this paper was unequivocally the most accurate and effective for this dataset.

438

## 439 **5. Online tool developed**

440 Machine learning models have seen their use increase exponentially in recent years due to their  
441 high performance on the predictions and for their ability of discovering robust patterns in complex  
442 datasets. However, their application is usually difficult because it is necessary to know the trained  
443 algorithm fully, with all its hyperparameters and to know the tools with which it has been

444 developed (e.g., programming language or specific software). This means that currently, many  
445 machine learning researches do not have a great practical application. In this study, an online  
446 application has been developed to facilitate the non-expert user, the usage of the machine learning  
447 model presented here, becoming one of the few studies that offers a practical implementation of  
448 the trained machine learning algorithm. Figure 9 shows the graphical user interface of the  
449 developed prediction tool, which is located at  
450 [https://huggingface.co/spaces/EstebanDC/Compression\\_Index](https://huggingface.co/spaces/EstebanDC/Compression_Index). In this tool, the value of  $c_c$  of a  
451 clay soil can be easily obtained by the user, defining the input values (i.e., LL(%), (PI %),  $w$ (%)  
452 and  $e_0$ ). Once these values are introduced, the compression index is calculated using the optimized  
453 super-learner ML model. This online tool can be used on any device with an internet connection  
454 (computers, tablets, or smartphones).  
455

A SUPER-LEARNER MACHINE LEARNING  
MODEL FOR A GLOBAL PREDICTION OF  
COMPRESSION INDEX IN CLAYS

Liquid limit (%)  
30

Plasticity index (%)  
15

Initial void ratio  
0.8

Natural water content (%)  
25

Clear Submit

Compression index  
[0.183]

456  
457 *Figure 9. Application developed with the trained super-learner machine learning model.*  
458

459  
460

## 461 **6. Symbolic regression**

462 In addition to the analysis carried out using the machine learning models presented previously,  
463 the dataset was analysed using Symbolic Regression. The aim of this analysis is to propose  
464 mathematical expressions, which relate  $c_c$  with some of the input variables. These expressions  
465 can be applied alternatively to the super-learner machine learning model when data for all the  
466 input variables are not available. Symbolic Regression-based approaches provide alternative  
467 machine learning methods that are recently gaining popularity. These approaches, as other ML  
468 algorithms, learn patterns from observed data and they have a great advantage over other ML  
469 algorithms due to their interpretability and explanation capabilities. Symbolic Regression  
470 attempts to explain a target variable by multiple input variables using a mathematical expression  
471 involving of a predefined set of basic computation functions. In this study, a symbolic regressor  
472 based on Feynman's path integral formulation from quantum field theory (Broløs et al. 2021) was  
473 applied. With this approach, four expressions (Equations 4 to 7) were obtained, and their  
474 performance metrics are presented in Table 8.

475

$$476 \quad c_c = 1.95 - 2.04 \tanh(1.82(0.204e_0 - 1)^2) \quad (4)$$

477

$$478 \quad c_c = 1.72e^{-4.02(1-0.209e_0)^4 - 0.599(0.00363LL-1)^2} + 0.0513 \quad (5)$$

479

$$480 \quad c_c = 3.32e^{-3.82(1-0.219e_0)^2 - 2.0e^{-9.78(1-0.244e_0)^2 - 0.237(-0.0135PI-1)^2}} - 0.0409 \quad (6)$$

481

$$482 \quad c_c = 2.67e^{-1.99(1-0.209e_0)^2 - 0.406(0.00286w-1)^2} - 0.26 \quad (7)$$

483

484 With  $w$ ,  $LL$ , and  $PI$  expressed in %.

485

486

487

Equation	Set	R <sup>2</sup>	MAE	RMSE
4	Training	0.85	0.049	0.083
	Test	0.87	0.042	0.065
5	Training	0.86	0.047	0.072
	Test	0.86	0.049	0.088
6	Training	0.86	0.047	0.078
	Test	0.91	0.043	0.063
7	Training	0.85	0.048	0.082
	Test	0.92	0.047	0.068

Table 8. Summary of performance metrics of the proposed equations.

488

489

490 Symbolic Regression approach gets good prediction performance, albeit being less accurate than  
491 the trained super-learner machine learning model, although presenting moderately long training  
492 times. Alternatively to the Equations 4 to 7, those obtained by linear regression and included in  
493 section 2 (Equations 1 to 3) can be used, although they offer worse metrics than those obtained  
494 by Symbolic Regression. These expressions must be fed with values of the input variables within  
495 the ranges with which the algorithm has been trained, since the behaviour outside these ranges is  
496 unknown.

497

## 498 7. Conclusions.

499 This study proposes a novel super-learner ML algorithm that provides a reliable and accurate  
500 model for predicting  $c_c$ , a key parameter in engineering applications. This algorithm is reasonably  
501 capable of predicting the value of  $c_c$  based on variables previously related to this parameter (i.e.,  
502 LL, PI,  $e_0$  and  $w$ ). To this aim, a database was built including more than 1000 samples collected  
503 worldwide, in order to reduce the local character of the proposals presented. The high and  
504 satisfactory error metrics, including R<sup>2</sup>, MAE, RMSE, and a20-index, demonstrate the model's  
505 robust performance, particularly considering the global nature of the dataset and the dispersion in  
506 the input variables. To facilitate the use of the model, a user-friendly prediction tool has been  
507 deployed online ([https://huggingface.co/spaces/EstebanDC/Compression\\_Index](https://huggingface.co/spaces/EstebanDC/Compression_Index)). Additionally,  
508 four symbolic regression expressions have been proposed for cases where all input parameters  
509 are not available. A comparative analysis was carried out using the compiled database against



510 global empirical correlations, demonstrating that super-learner ML model was the most accurate  
511 and effective for this dataset. The global scope of this study and the development of an online tool  
512 enhance the proposed algorithm's applicability and value beyond other methods in the field,  
513 thereby making it more useful for geotechnical engineers. Future work will focus on expanding  
514 the dataset and updating the model to increase its generalization ability and application range.  
515 Finally, it is important to exercise caution when applying the model or expressions beyond the  
516 specified input variable ranges, as this can lead to predictions of questionable reliability, a  
517 common issue with predictive algorithms.

518

519

520 **Availability statement**

521 The data that support the findings are available from the corresponding author upon reasonable  
522 request.

523

524 **References**

525

- 526 Al-Khafaji, A. & Andersland, O. 1992. Equations for compression index approximation. Journal  
527 of Geotechnical Engineering, **118**, 148-153, doi: [https://doi.org/10.1061/\(ASCE\)0733-](https://doi.org/10.1061/(ASCE)0733-9410(1992)118:1(148))  
528 [9410\(1992\)118:1\(148\)](https://doi.org/10.1061/(ASCE)0733-9410(1992)118:1(148)).
- 529 Alam, S., Khuntia, S. & Patra, C. 2014. Prediction of compression index of clay using artificial  
530 neural network. *International conference on industrial engineering science and applications-*  
531 *NIT, Durgapur*.
- 532 Alhaji, M.M., Alhassan, M., Tsado, T.Y. & Mohammed, Y.A. 2017. Compression Index Prediction  
533 Models for Fine-grained Soil Deposits in Nigeria. Proceeding of the 2nd International Engineering  
534 Conference, Federal ...
- 535 Apostolopoulou, M., Asteris, P.G., Armaghani, D.J., Douvika, M.G., Lourenço, P.B., Cavaleri, L.,  
536 Bakolas, A. & Moropoulou, A. 2020. Mapping and holistic design of natural hydraulic lime  
537 mortars. *Cement and Concrete Research*, **136**, 106167, doi:  
538 <https://doi.org/10.1016/j.cemconres.2020.106167>.
- 539 Armaghani, D.J. & Asteris, P.G. 2021. A comparative study of ANN and ANFIS models for the  
540 prediction of cement-based mortar materials compressive strength. *Neural Computing and*  
541 *Applications*, **33**, 4501-4532.
- 542 Asteris, P.G., Armaghani, D.J., Hatzigeorgiou, G.D., Karayannis, C.G. & Pilakoutas, K. 2019.  
543 Predicting the shear strength of reinforced concrete beams using Artificial Neural Networks.  
544 *Computers and Concrete, An International Journal*, **24**, 469-488.
- 545 Asteris, P.G., Koopialipour, M., Armaghani, D.J., Kotsonis, E.A. & Lourenço, P.B. 2021a. Prediction  
546 of cement-based mortars compressive strength using machine learning techniques. *Neural*  
547 *Computing and Applications*, **33**, 13089-13121, doi: 10.1007/s00521-021-06004-8.
- 548 Asteris, P.G., Lourenço, P.B., Hajihassani, M., Adami, C.-E.N., Lemonis, M.E., Skentou, A.D.,  
549 Marques, R., Nguyen, H., Rodrigues, H. & Varum, H. 2021b. Soft computing-based models for  
550 the prediction of masonry compressive strength. *Engineering Structures*, **248**, 113276, doi:  
551 <https://doi.org/10.1016/j.engstruct.2021.113276>.
- 552 Asteris, P.G., Mamou, A., Ferentinou, M., Tran, T.-T. & Zhou, J. 2022. Predicting clay  
553 compressibility using a novel Manta ray foraging optimization-based extreme learning machine  
554 model. *Transportation Geotechnics*, **37**, 100861, doi:  
555 <https://doi.org/10.1016/j.trgeo.2022.100861>.
- 556 Aydın, Y., Işıkdağ, Ü., Bekdaş, G., Nigdeli, S.M. & Geem, Z.W. 2023. Use of Machine Learning  
557 Techniques in Soil Classification. *Sustainability*, **15**, 2374.
- 558 Azzouz, A.S., Krizek, R.J. & Corotis, R.B. 1976. Regression analysis of soil compressibility. *Soils*  
559 *and Foundations*, **16**, 19-29, doi: [https://doi.org/10.3208/sandf1972.16.2\\_19](https://doi.org/10.3208/sandf1972.16.2_19).
- 560 Baghbani, A., Choudhury, T., Costa, S. & Reiner, J. 2022. Application of artificial intelligence in  
561 geotechnical engineering: A state-of-the-art review. *Earth-Science Reviews*, **228**, 103991.
- 562 Bardhan, A., Alzo'ubi, A.K., Palanivelu, S., Hamidian, P., GuhaRay, A., Kumar, G., Tsoukalas, M.Z.  
563 & Asteris, P.G. 2023. A hybrid approach of ANN and improved PSO for estimating soaked CBR of  
564 subgrade soils of heavy-haul railway corridor. *International Journal of Pavement Engineering*,  
565 **24**, 2176494, doi: 10.1080/10298436.2023.2176494.

566 Bardhan, A., Gokceoglu, C., Burman, A., Samui, P. & Asteris, P.G. 2021. Efficient computational  
567 techniques for predicting the California bearing ratio of soil in soaked conditions. *Engineering*  
568 *Geology*, **291**, 106239, doi: <https://doi.org/10.1016/j.enggeo.2021.106239>.

569 Benbouras, M.A., Kettab Mitiche, R., Zedira, H., Petrisor, A.-I., Mezouar, N. & Debiche, F. 2019.  
570 A new approach to predict the compression index using artificial intelligence methods. *Marine*  
571 *Georesources & Geotechnology*, **37**, 704-720.

572 Bowles, J.E. 1979. *Physical and geotechnical properties of soils*.

573 Broløfs, K.R., Machado, M.V., Cave, C., Kasak, J., Stentoft-Hansen, V., Batanero, V.G., Jelen, T. &  
574 Wilstrup, C. 2021. An approach to symbolic regression using feyn. arXiv preprint  
575 arXiv:2104.05417.

576 Carter, M. & Bentley, S.P. 1991. *Correlations of soil properties*. Pentech press publishers.

577 Craig, R.F. 2004. *Craig's soil mechanics*. CRC press.

578 Dam Nguyen, D., Roussis, P.C., Thai Pham, B., Ferentinou, M., Mamou, A., Quang Vu, D., Thi Bui,  
579 Q.-A., Kien Trong, D. & Asteris, P.G. 2022. Bagging and Multilayer Perceptron Hybrid Intelligence  
580 Models Predicting the Swelling Potential of Soil. *Transportation Geotechnics*, **36**, 100797, doi:  
581 <https://doi.org/10.1016/j.trgeo.2022.100797>.

582 Das, B.M. 2021. *Principles of geotechnical engineering*. Cengage learning.

583 Desai, V.M., Desai, V. & Rao, D. 2009. Prediction of compression index using artificial neural  
584 networks. *Indian geotechnical conference (IGC-2009), Guntur, India*, 614-617.

585 Díaz, E., Salamanca-Medina, E.L. & Tomás, R. 2023. Assessment of compressive strength of jet  
586 grouting by machine learning. *Journal of Rock Mechanics and Geotechnical Engineering*, doi:  
587 <https://doi.org/10.1016/j.jrmge.2023.03.008>.

588 Díaz, E. & Tomás, R. 2021. Upgrading the prediction of jet grouting column diameter using deep  
589 learning with an emphasis on high energies. *Acta Geotechnica*, **16**, 1627-1633, doi:  
590 <https://doi.org/10.1007/s11440-020-01091-8>.

591 Friedman, J.H.J.A.o.s. 2001. Greedy function approximation: a gradient boosting machine. 1189-  
592 1232.

593 Geurts, P., Ernst, D. & Wehenkel, L. 2006. Extremely randomized trees. *Machine learning*, **63**, 3-  
594 42.

595 Ho, T.K. 1995. Random decision forests. *Proceedings of 3rd international conference on*  
596 *document analysis and recognition*. IEEE, 278-282.

597 Hough, B.K. 1957. *Basic Soils Engineering*.

598 Kalantary, F. & Kordnaeij, A. 2012. Prediction of compression index using artificial neural  
599 network. *Scientific Research and Essays*, **7**, 2835-2848.

600 Kennedy, J. & Eberhart, R. 1995. Particle swarm optimization. *Proceedings of ICNN'95-*  
601 *international conference on neural networks*, **4**, 1942-1948, doi: 10.1109/ICNN.1995.488968.

602 Koppula, S. 1981. Statistical estimation of compression index. *ASTM Geotechnical Testing*  
603 *Journal*, **4**.

604 Kumar, V.P. & Rani, C.S. 2011. Prediction of compression index of soils using artificial neural  
605 networks (ANNs). *Int J Eng Res Appl*, **1**, 1554-1558.

606 Laan, M.J.v.d., Polley, E.C. & Hubbard, A.E. 2007. Super Learner. *Statistical Applications in*  
607 *Genetics and Molecular Biology*, **6**, doi: doi:10.2202/1544-6115.1309.

608 LCPC. 1977. Remblais sur sols compressibles. *Bulletin des Laboratoires des Ponts et Chaussees*,  
609 *Special T*, 58. Paris (France). 361p.

610 Little, R.J. & Rubin, D.B. 2019. *Statistical analysis with missing data*. John Wiley & Sons.

611 Long, T., He, B., Ghorbani, A. & Khatami, S.M.H. 2023. Tree-Based Techniques for Predicting the  
612 Compression Index of Clayey Soils. *Journal of Soft Computing in Civil Engineering*, **7**, 52-67.

613 McCabe, B.A., Sheil, B.B., Long, M.M., Buggy, F.J. & Farrell, E.R. 2014. Empirical correlations for  
614 the compression index of Irish soft soils. *Proceedings of the Institution of Civil Engineers-*  
615 *Geotechnical Engineering*, **167**, 510-517.

616 Mitachi, T. & Ono, T. 1985. Prediction of undrained shear strength of overconsolidated clay.  
617 *Tsuchi to Kiso, JSSMFE*, **33**, 21-26.

618 Nesamatha, R. & Arumairaj, P. 2015. Numerical modeling for prediction of compression index  
619 from soil index properties. *Electron J Geotech Eng*, **20**, 4369-4378.

620 Nishida, Y. 1956. A brief note on compression index of soil. *Journal of the Soil Mechanics and*  
621 *Foundations Division*, **82**, 1027-1021-1027-1014.

622 Onyejekwe, S., Kang, X. & Ge, L. 2015. Assessment of empirical equations for the compression  
623 index of fine-grained soils in Missouri. *Bulletin of Engineering Geology and the Environment*, **74**,  
624 705-716.

625 Onyejekwe, S., Kang, X. & Ge, L. 2016. Evaluation of the scale of fluctuation of geotechnical  
626 parameters by autocorrelation function and semivariogram function. *Engineering Geology*, **214**,  
627 43-49.

628 Park, H.I. & Lee, S.R. 2011. Evaluation of the compression index of soils using an artificial neural  
629 network. *Computers and Geotechnics*, **38**, 472-481.

630 Reagan, L.A. & Kiemele, M.J. 2008. *Design for six sigma: The tool guide for practitioners*. CTQ  
631 Media.

632 Rendon-Herrero, O. 1980. Universal compression index equation. *Journal of the Geotechnical*  
633 *Engineering Division*, **106**, 1179-1200.

634 Salvatore, E., Modoni, G., Spagnoli, G., Arciero, M., Mascolo, M.C. & Ochmański, M. 2022.  
635 Conditioning clayey soils with a dispersant agent for Deep Soil Mixing application: laboratory  
636 experiments and artificial neural network interpretation. *Acta Geotechnica*, **17**, 5073-5087.

637 Samui, P., Kim, D., Das, S. & Yoon, G.L. 2012. Determination of Compression Index for Marine  
638 Clay: A Relevance Vector Machine Approach. *Marine Georesources & Geotechnology*, **30**, 263-  
639 273, doi: 10.1080/1064119X.2011.614323.

640 Schölkopf, B., Williamson, R.C., Smola, A., Shawe-Taylor, J. & Platt, J. 1999. Support vector  
641 method for novelty detection. *Advances in neural information processing systems*, **12**.

642 Singh, M.J., Kaushik, A., Patnaik, G., Xu, D.-S., Feng, W.-Q., Rajput, A., Prakash, G. & Borana, L.  
643 2023. Machine learning-based approach for predicting the consolidation characteristics of soft  
644 soil. *Marine Georesources & Geotechnology*, 1-15, doi: 10.1080/1064119X.2023.2193174.

645 Sobol, I.Y.M. 1990. On sensitivity estimation for nonlinear mathematical models.  
646 *Matematicheskoe modelirovanie*, **2**, 112-118.

647 Spagnoli, G. & Shimobe, S. 2020. Statistical analysis of some correlations between compression  
648 index and Atterberg limits. *Environmental Earth Sciences*, **79**, 532.

649 Sridharan, A. & Nagaraj, H. 2000. Compressibility behaviour of remoulded, fine-grained soils and  
650 correlation with index properties. *Canadian Geotechnical Journal*, **37**, 712-722.

651 Terzaghi, K., Peck, R.B. & Mesri, G. 1967. *Soil mechanics in engineering practice*. 2nd edition.  
652 John Wiley and Sons, New York.

653 Trong, D.K., Pham, B.T., Jalal, F.E., Iqbal, M., Roussis, P.C., Mamou, A., Ferentinou, M., Vu, D.Q.,  
654 Duc Dam, N. & Tran, Q.A. 2021. On random subspace optimization-based hybrid computing  
655 models predicting the california bearing ratio of soils. *Materials*, **14**, 6516.

656 Tsuchida, T. 1991. A new concept of e-logp relationship for clays. *Proce, 9th Asian regional*  
657 *Conference on soil mechanics and foundation Engineering*, 87-90.

658 Van Buuren, S. & Oudshoorn, C.G. 2000. *Multivariate imputation by chained equations*. Leiden:  
659 TNO.

660 Verbrugge, J.-C. & Schroeder, C. 2018. *Geotechnical correlations for soils and rocks*. John Wiley  
661 & Sons.

662 Widodo, S. & Ibrahim, A. 2012. Estimation of primary compression index (Cc) using physical  
663 properties of Pontianak soft clay. *International Journal of Engineering Research and*  
664 *Applications*, **2**, 2231-2235.

665 Wroth, C. & Wood, D. 1978. The correlation of index properties with some basic engineering  
666 properties of soils. *Canadian Geotechnical Journal*, **15**, 137-145.

667 Zhang, P., Yin, Z.-Y., Jin, Y.-F., Chan, T.H.T. & Gao, F.-P. 2021. Intelligent modelling of clay  
668 compressibility using hybrid meta-heuristic and machine learning algorithms. *Geoscience*  
669 *Frontiers*, **12**, 441-452, doi: <https://doi.org/10.1016/j.gsf.2020.02.014>.

670 Zhang, W., Gu, X., Hong, L., Han, L. & Wang, L. 2023. Comprehensive review of machine learning  
671 in geotechnical reliability analysis: Algorithms, applications and further challenges. Applied Soft  
672 Computing, **136**, 110066, doi: <https://doi.org/10.1016/j.asoc.2023.110066>.  
673 Zhang, W., Gu, X., Tang, L., Yin, Y., Liu, D. & Zhang, Y. 2022. Application of machine learning,  
674 deep learning and optimization algorithms in geoengineering and geoscience: Comprehensive  
675 review and future challenge. Gondwana Research, **109**, 1-17.

676

677