

GRADIENT BOOSTING TREES WITH BAYESIAN OPTIMIZATION TO PREDICT ACTIVITY FROM OTHER GEOTECHNICAL PARAMETERS

Esteban Díaz¹, Giovanni Spagnoli^{2*}

¹ Departamento de Ingeniería Civil. Escuela Politécnica Superior, Universidad de Alicante, P.O.

Box 99, E-03080 Alicante, Spain, esteban.diaz@ua.es

² DMT GmbH & Co. KG, Am TÜV 1, 45307 Essen, Germany, giovanni.spagnoli@dm-t-

group.com ORCID: 0000-0002-1866-4345

* Corresponding author.

Abstract

Clay swell potential can be classified based on the value of activity and it is defined as the ratio of plasticity index to clay content as a percentage. This paper outlines the investigation into how activity correlates with other key properties of clayey soils. Specifically, four approaches were evaluated for predicting activity using: a) liquid limit (LL), specific surface area (SSA), cation exchange capacity (CEC) and clay content; b) LL, SSA and CEC; c) LL; and d) SSA and CEC. For this purpose, a database of 104 samples was collected from which 35 machine learning algorithms were trained. Gradient Boosting Trees showed the highest prediction accuracy in the four approaches and, to improve its prediction performance, a Bayesian Optimization was applied to tune their hyperparameters, resulting in the final models. The performance of the developed models was evaluated, showing prominent results with exceptionally good metrics, except in the approach from SSA and CEC where the trained algorithm was not capable of predicting activity with confidence ($R^2=0.46$). This algorithm can predict activity using only LL with high accuracy ($R^2=0.94$), and when combined with SSA and CEC, the precision is further enhanced ($R^2=0.96$). Finally, a variable importance analysis was performed, indicating LL is the variable with the greatest influence in predicting activity.

Keywords: machine learning; liquid limit; specific surface area; cation exchange capacity; clay content; activity.

29 **1. Introduction**

30 Key properties of clays such as Atterberg limits, Specific Surface Area (SSA) and Cation
31 Exchange Capacity (CEC), are important in geotechnical engineering and in particular for
32 characterizing expansive soils. Expansive soils are a very problematic matter in several fields of
33 civil engineering, posing significant amount of damage (Jones Jr and Holtz 1973). Expansive soils
34 are those in which the variation of water content results in a large volume change (Kar 2021). The
35 construction on these soils often generates serious issues, especially when lightweight structures
36 are built. There are several approaches to estimate the expansive potential of a soil. Expansive
37 soils can be identified by the Atterberg limits, clay content, or a combination of both. Skempton
38 (1953) proposed the concept of activity, i.e. the ratio of plasticity index (PI) to clay fraction
39 content, which can be utilized as an index property to establish the swelling potential of expansive
40 soil. Peck et al. (1974) correlated PI with the expansion potential. Zapata et al. (2006) suggested
41 that considering the % passing 75 μ m improved the prediction. Different authors developed other
42 types of tests to identify expansive soils (e.g. Lambe 1960; Yao et al. 2004). SSA (Chittoori and
43 Puppala 2011) and CEC (Mitchell and Soga 2005; Nelson et al. 2015) are also used to indirectly
44 identify expansive soils. According to Low (1987), the surface's level of hydration has a
45 significant impact on how clays behave. This makes the SSA of a clayey soil extremely
46 significant, and it has been proposed that SSA can be used to forecast the engineering behaviour
47 of fine-grained soils (e.g. Warkentin 1972). The principal clay minerals, such as montmorillonite,
48 illite, and kaolinite, have a significant impact on the SSA of soils, which reflects the consistency
49 traits as well as the clay concentration (Spagnoli and Shimobe 2019). Muhunthan (1991)
50 attempted a rheological correlation between SSA and LL. CEC and the soil's swelling
51 characteristics are strongly connected. With an increase in CEC, there is also an increase in soil
52 swelling (Christidis 1998). The higher the presence of smectitic clay minerals (i.e. expansive
53 clays), the more the clay swells. According to Al-Rawas (1999), the cations are what regulate
54 how expansive soils are. Therefore, both LL and CEC are variables that regulate the soils'
55 propensity to swell (Spagnoli and Shimobe 2019). Chittoori and Puppala (2011) suggested that

56 since SSA of smectitic soils is higher than those of kaolitic soils, SSA can be used to predict the
57 expansion potential, as the water holding capacity is higher.

58 Although activity is not a complex or difficult parameter to obtain, it is a relevant property of
59 clayey soils correlated with swelling potential. Several authors have proposed a correlation
60 between soil activity and other geotechnical parameters (e.g. Polidori 2009; Spagnoli and
61 Shimobe 2019). However, in general, the proposed geotechnical correlations consider only two
62 parameters at the same time (input and output). Machine learning (ML) techniques have
63 progressively emerged as an alternative approach to address numerous geotechnical challenges
64 (e.g. Díaz et al. 2018; Díaz et al. 2021; Díaz and Tomás 2021; Phoon and Zhang 2023; Salvatore
65 et al. 2022; Wang et al. 2020; Zhang et al. 2022a; Zhang et al. 2023; Zhang et al. 2022b).
66 Motivated by these advantages, this paper employs machine learning techniques to explore
67 potential correlations between soil activity and other properties of clayey soils (LL, SSA, CEC
68 and clay content). To achieve this, the study focused on four approaches based on various
69 properties associated with activity, comparing their outcomes to evaluate the performance of each.
70 The objective is to understand clearly which properties can be utilized to determine activity with
71 an appropriate degree of accuracy. To this end, four prediction models to forecast activity were
72 developed using a dataset of experimental results. With this dataset, a comparative study of
73 various ML algorithms was carried out. The algorithms that performed best underwent Bayesian
74 optimization to determine the appropriate model hyperparameters. The results of the final tuned
75 models in predicting activity are then evaluated and discussed. Finally, an importance analysis of
76 the variables is conducted to identify the most important parameters in the prediction of activity.
77

78 **2. Database**

79 104 data points regarding clay content $<2\mu\text{m}$, SSA (m^2/g), CEC ($\text{meq}/100\text{g}$), and LL obtained by
80 means of the Casagrande cup were acquired from 12 different publications and a single datapoint
81 belonging to the authors. The activity value was obtained from the data. Pure clays and natural
82 clays were selected, in order to have a relatively high heterogeneity in the data. Specifically data
83 from (Arifin 2008), Cerato (2001), Cerato and Lutenegeger (2002), Cerato and Lutenegeger (2004),

84 Cerato and Lutenege (2005), Marcial (2013), Mishra et al. (2012), Santhoshkumar et al. (2016),
85 Schwing et al. (2013), Sivapullaiah et al. (2008), Spagnoli et al. (2013) and Zhang et al. (2003)
86 were used.

87

88 **2.1. Statistical description**

89 Some descriptive statistical information is provided in Table 1. Skewness and kurtosis are two
90 measures of the shape of a distribution in statistics. Skewness measures the degree of asymmetry
91 in a distribution, while kurtosis measures the degree of peakedness or flatness. Skewness is
92 defined as the third standardized moment of a distribution and can be positive (skewed right),
93 negative (skewed left), or zero (symmetric). Positive skewness indicates that the tail of the
94 distribution is longer on the right-hand side, while negative skewness signifies that the tail is
95 longer on the left-hand side. Kurtosis is defined as the fourth standardized moment of a
96 distribution and can be high (leptokurtic) indicating that the distribution has a sharp peak and fat
97 tails, or low (platykurtic) kurtosis signifying a flatter distribution with fewer outliers (Wackerly
98 et al. 2014).

99 Figure 1 presents the box plots for the data. A box plot is a visual representation of the distribution
100 of a dataset that shows the median, quartiles, and outliers. It consists of a rectangular box, which
101 extends from the first quartile (Q1) to the third quartile (Q3), with a vertical line inside that
102 represents the median. The distance between the upper and lower edges of the box, known as the
103 interquartile range (IQR), contains the middle 50% of the data. The whiskers extend from the box
104 to the minimum and maximum values within a certain range of the dataset. The range is typically
105 set at 1.5 times the IQR (Gelman and Hill 2006). Figure 1 shows several outliers, which are data
106 points outside of this range and are plotted individually as dots outside the whiskers. Figure 2
107 shows selected histograms of the data. The goodness of fit for the probability distributions is
108 conducted using the Anderson-Darling (AD) test. The AD is a statistical test based on the idea of
109 comparing the cumulative distribution function (CDF) of the sample data to the CDF of the
110 theoretical distribution being tested (Hollander, et al., 2015). Notably, with the exception of clay

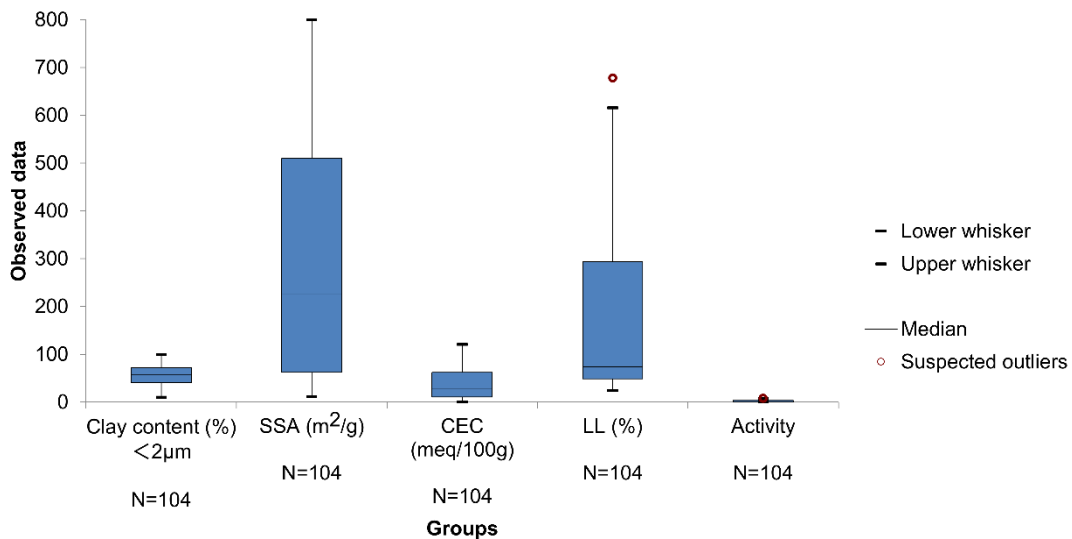
111 content, none of the data adheres to the normal distribution. AD gives more weight to the tails
 112 than does the Kolmogorov-Smirnov (KS) test (Stephens 1974).

113

Statistics	Clay content (%) <2 μ m	SSA (m ² /g)	CEC (meq/100g)	LL (%)	Activity
Count	104	104	104	104	104
Mean	56.039	285.459	37.834	181.361	2.189
Median	57	225.5	28.25	74	0.99099
Mode	36.2	15	2	42	0.44
Minimum	10.2	11	0.8	24	0.2
Maximum	100	800	120.9	678	8.0132
Skewness	-0.18	0.49	0.81	1.14	1.15
Kurtosis	-0.61	-1.12	-0.135	-0.132	-0.053
Standard Deviation	20.23	235.25	31.07	180.77	2.24

114

Table 1. Descriptive statistics of the data analyzed.



115

116

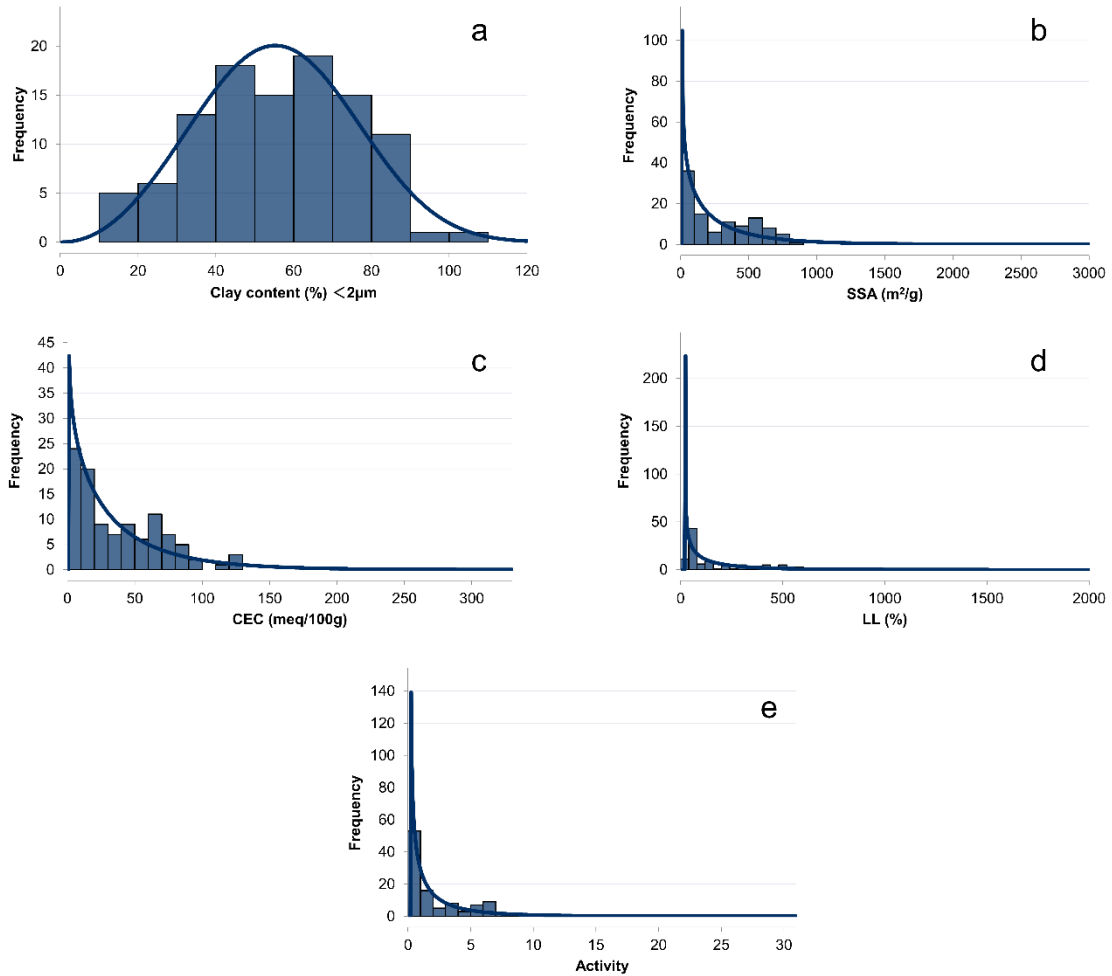
Figure 1. Box plots of the data analyzed.

117

118 For the data presented in Table 1, the clay content values follow a 2-Param Weibull distribution,
 119 while the remaining parameters follow a 3-Param Weibull distributions (see histograms in Figure
 120 2). Additionally, a normality test has been performed for all data sets. A normality test is a
 121 statistical test used to determine whether a given set of data comes from a normally distributed
 122 population. Normality tests are used to check the assumption of normality, which is often made
 123 in statistical inference procedures such as hypothesis testing and confidence interval estimation

124 (Devore et al. 2013). The normality test used was the Shapiro-Wilk test, which is based on the
125 deviation between the observed data and the expected normal distribution.

126



127

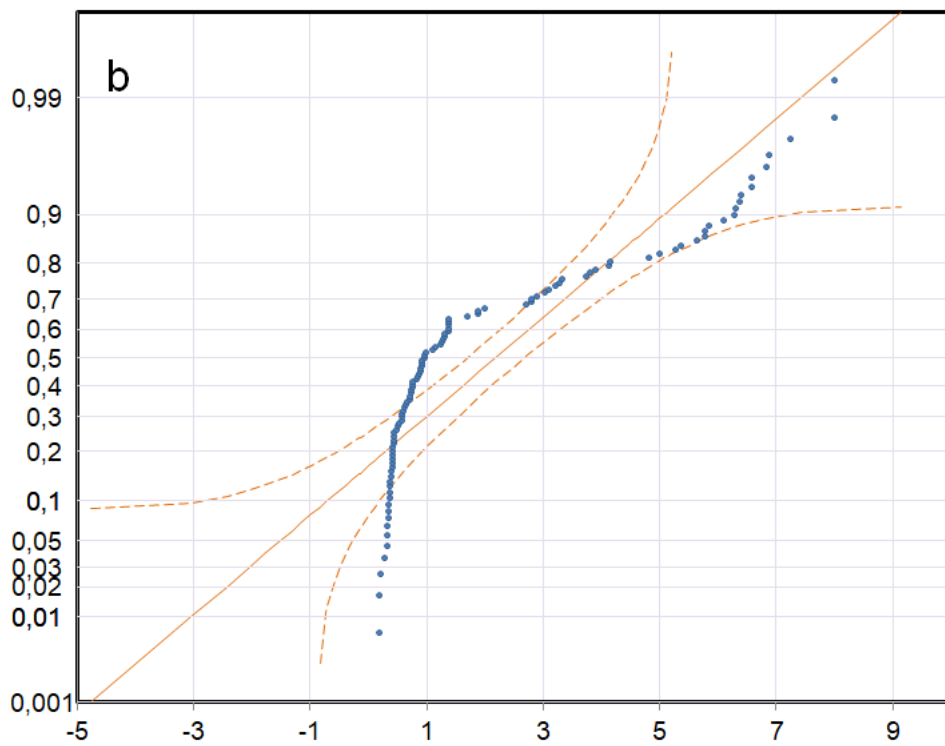
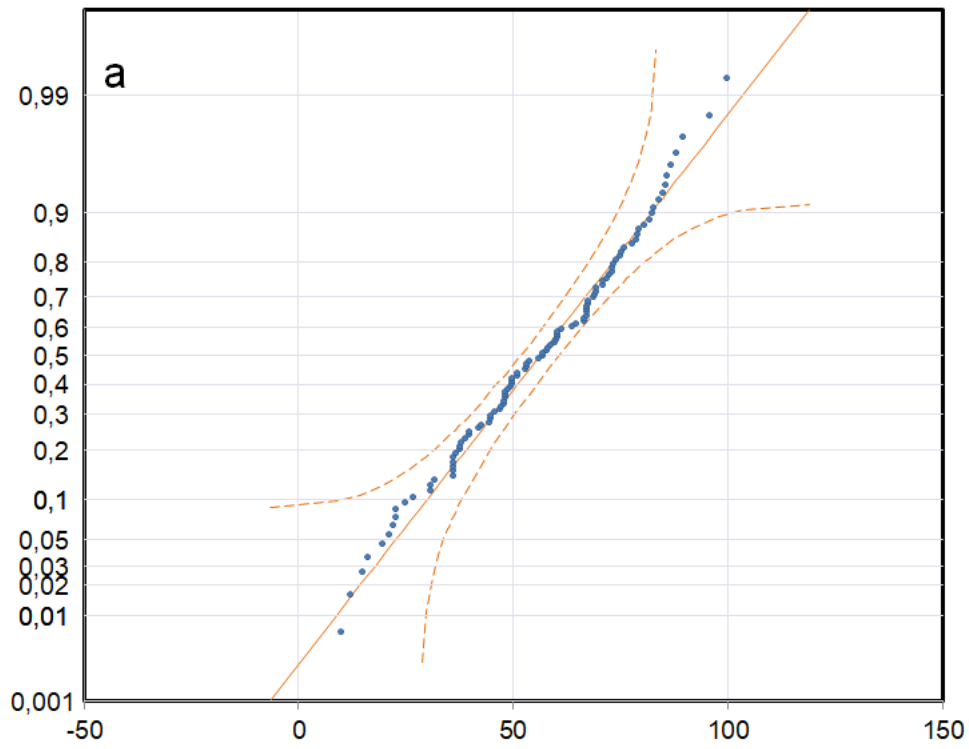
128 *Figure 2. Histograms for the data from Table 1 for a) clay content, b) SSA, c) CEC, d) LL, and*
129 *e) activity*

130

131 Figure 3 shows the probability plots for clay content and activity, as an example. A probability
132 plot is a graphical technique used in statistics to assess whether a set of data follows a particular
133 distribution. It compares the ordered values of a dataset to the expected values of a theoretical
134 distribution. If the data follows the distribution, the plot should show a roughly straight line. If
135 the data does not follow the distribution, the plot will show deviations from the straight line
136 (Moore and McCabe 1989). In a probability plot, the x-axis represents the expected values of the

137 theoretical distribution being tested, while the y-axis represents the ordered values of the dataset
138 being tested. The points on the plot are plotted based on their rank order, such that the smallest
139 observation is plotted at the far left and the largest observation is plotted at the far right. The
140 closer the points are to the straight line, the better the fit of the data to the theoretical normal
141 distribution being tested. The probability plots for CEC, SSA and LL are similar to Figure 3b. It
142 is possible to observe that while the data follow a straight line for clay content (Figure 3a), this is
143 not the case for activity (Figure 3b).

144



145

146 *Figure 3. Probability plots for a) clay content and b) activity. To note the difference considering*

147

their normality path.

148

149 **3. Methodology.**

150 With the 104 samples described in the previous section, the four variables previously presented
151 were used to predict activity of clayey soils. Specifically, four approaches have been studied
152 taking into account different variables (Table 2) and these predict activity from: 1) LL, clay
153 content, SSA and CEC, 2) LL, SSA and CEC, 3) LL, and 4) SSA and CEC.

154

Approach 1	Approach 2	Approach 3	Approach 4
LL	LL	LL	SSA
Clay content	SSA	-	CEC
SSA	CEC	-	-
CEC	-	-	-

155 *Table 2. Variables considered in each of the approaches considered in the prediction of activity.*

156 **3.2. Machine learning methods**

157 **3.2.1. Gradient boosting trees**

158 Gradient boosting trees are a powerful family of machine learning algorithms for performing
159 gradient descent on decision trees using the boosting ensemble learning method. The main idea
160 behind them is to combine iteratively several simple models (i.e. weak learners) to obtain a model
161 with enhanced prediction accuracy (i.e. strong learner). Boosting algorithms were initially
162 proposed for classification tasks (Freund 1995; Freund and Schapire 1996; Schapire 1990).
163 Friedman (2001) expanded the boosting to regression tasks by creating the gradient boosting
164 machines method (GBM). The boosting method adjusts the weights of the training sample
165 according to the last iteration and assigns more weight to observations that are difficult to predict
166 and less weight to those that have already been well managed. It can be understood as a numerical
167 optimization algorithm aiming to find an additive model that maximally reduces the loss function.
168 The GBM algorithm builds successive decision trees to fit one training example at a time (the tree
169 that best reduces the loss function). As it fits each new sample, it updates its knowledge of which
170 features are important for the prediction of future samples. It starts with an initial estimation for
171 model parameters and iteratively enhances these estimations until a required level of accuracy has
172 been achieved or some other stopping criterion has been satisfied. Specifically, in regression
173 problems, the algorithm begins by initializing the model with a first prediction, which is a decision

174 tree that maximally minimizes the loss function (mean squared error in regression), then at each
175 stage a new decision tree is fitted to the existing residual and added to the prior model to update
176 the previously obtained residuals. The algorithm keeps on iterating until the prefixed maximum
177 number of iterations is reached. This process is called stage wise, meaning that at each new stage,
178 the decision trees included into the model at previous steps remain unchanged. With the process
179 of fitting decision trees to the residuals, the algorithm is enhanced in the zones where it does not
180 perform well. Four hyperparameters mainly govern behaviour of the GBM: 1) the learning rate,
181 2) the number of boosting stages to perform, 3) the number of features to consider when looking
182 for the best split, and 4) the maximum depth of the individual regression estimators.

183

184 **3.2.2. Bayesian hyperparameter optimization**

185 The performance of the machine learning algorithms is strongly determined by the model
186 parameters (i.e. hyperparameters) which need to be set before training. In this study, the optimal
187 hyperparameters were established applying Bayesian optimization (Shahriari et al. 2015; Snoek
188 et al. 2012) which is a general technique for function optimization. Bayesian optimization builds
189 a probability model based on the previous evaluation results of the target for finding the value
190 that minimizes the objective function. Bayesian hyperparameter optimization was performed
191 using BayesSearchCV from the Scikit-optimize package on Python (Head et al. 2020). This
192 method uses stepwise Bayesian Optimization to discover the most promising hyperparameters in
193 the problem-space. This optimization method was chosen for hyperparameter tuning due to its
194 efficiency (Eriksson et al. 2019) and because it has been proven to be superior to other
195 optimization algorithms on many optimization benchmark functions (Jones 2001). Finally, it must
196 be remarked that this technique has been extensively used for machine learning hyperparameter
197 tuning in geotechnics (e.g. Li et al. 2022; Zhang et al. 2021).

198

199 **3.2.3. Verification and evaluation of the machine learning models**

200 In this paper, a training set was employed to choose and build the predictive models and a test set
201 was used to examine the trained models in each of the four approaches analysed. The coefficient

202 of determination (R^2), the mean square error (MSE) and the mean absolute error (MAE), were
203 applied with the aim of evaluating the reliability of the algorithms considered and to interpret the
204 correspondence between predictions and observed values. The definition of R^2 , MSE and MAE
205 is expressed by (Equations 1 to 3):

206

$$207 \quad R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (1)$$

208

$$209 \quad MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (2)$$

210

$$211 \quad MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (3)$$

212

213 where y is the measured value, \hat{y} is the model predicted value, \bar{y} is the average of the measured
214 values, and n is the number of samples in the training or testing sets.

215

216 On the other hand, k-fold cross-validation (Stone 1974), was used in the Bayesian optimization.

217 In the k-fold cross-validation approach, the dataset is randomly shuffled and then divided into k

218 folds. k-1 folds are used to train the model and the remaining fold (the test set) is employed for

219 the evaluation. The process is repeated k times, and performance of the model is evaluated by the

220 mean prediction error of k sub-datasets. In this study, the optimization was done with a 5-fold

221 cross-validation.

222

223 **3.2.4. Data preparation. Feature rescaling**

224 To improve the performance of the machine learning algorithms, the dataset was pre-processed.

225 In particular, the input parameters of the dataset were standardized using the min-max scaling,

226 which involves rescaling the range of features to scale the range in [0, 1]. Using this data rescaling

227 method, the impact of parameters with different scales on the algorithm performance can be

228 minimised.

229 **3.2.5. Machine learning algorithm selection**

230 35 ML algorithms were built and trained in each of the four approaches considered, using the
231 Scikit-learn package (Pedregosa et al. 2011) which is the most useful and robust library for
232 machine learning in Python. The initial selection included: SVR, Random Forest Regressor, Extra
233 Trees Regressor, AdaBoost Regressor, NuSVR, Gradient Boosting Regressor, K-Neighbors
234 Regressor, Histogram-based Gradient Boosting Regressor, Bagging Regressor, MLP Regressor,
235 Huber Regressor, Linear SVR, Ridge CV, Bayesian Ridge, Ridge, Linear Regression,
236 Transformed Target Regressor, Lasso CV, Elastic Net CV, Lasso Lars CV, Lasso Lars IC, Lars
237 CV, Lars, SGD Regressor, RANSAC Regressor, Elastic Net, Lasso, Orthogonal Matching Pursuit
238 CV, Passive Aggressive Regressor, Gaussian Process Regressor, Orthogonal Matching Pursuit,
239 Decision Tree Regressor, Dummy Regressor, Lasso Lars and Kernel Ridge.

240

241 **3.3. Model conception**

242 To choose the optimal model to predict activity in each of the four approaches considered, the
243 next phases were followed:

- 244 1. Building a database, collecting data from different research papers. In this phase, 104
245 samples of clayey soils containing values of LL, clay content, SSA, CEC and activity
246 were gathered.
- 247 2. Rescaling of the input variables.
- 248 3. Application of 35 Machine Learning Algorithms: Using the chosen inputs for each of the
249 considered approaches.
- 250 4. Identification of the best models, considering R^2 as the main statistical performance
251 indicator.
- 252 5. Optimizing the best models for each of the four approaches using Bayesian optimization.
- 253 6. Assessing the predictive capability of the four selected models considering the test set.
- 254 7. Performing a feature importance analysis of the selected models to find the inputs with a
255 higher influence on the predictions.

256

257 **4. Application and results**

258 **4.1. Predictive comparisons among different algorithms**

259 In this section, from the 35 machine learning algorithms considered, a selection process is carried
 260 out, choosing R^2 as the reference metric. Table 3 displays the results of this analysis, but only
 261 highlights the top ten results for each of the considered approaches. From the analysis of the
 262 comparative study, it can be deduced that for the four approaches, the algorithm with the best
 263 performance is the Gradient Boosting Regressor Trees (GBRT).

Ranking	Approach 1		Approach 2		Approach 3		Approach 4	
	Algorithm	R^2	Algorithm	R^2	Algorithm	R^2	Algorithm	R^2
1	Gradient Boosting Regressor	0.98	Gradient Boosting Regressor	0.96	Gradient Boosting Regressor	0.95	Gradient Boosting Regressor	0.45
2	Extra Tree Regressor	0.96	Extra Trees Regressor	0.94	Bagging Regressor	0.94	Passive Aggressive Regressor	0.44
3	Gaussian Process Regressor	0.95	Decision Tree Regressor	0.93	AdaBoost Regressor	0.94	Extra Trees Regressor	0.43
4	Random Forest Regressor	0.93	Random Forest Regressor	0.93	SVR	0.94	MLP Regressor	0.42
5	Linear SVR	0.93	RANSAC Regressor	0.93	Random Forest Regressor	0.94	Linear SVR	0.42
6	Huber Regressor	0.93	Linear SVR	0.93	Decision Tree Regressor	0.93	AdaBoost Regressor	0.41
7	Ridge	0.93	Huber Regressor	0.93	XGB Regressor	0.93	Poisson Regressor	0.40
8	Passive Aggressive Regressor	0.92	Bagging Regressor	0.93	Nu SVR	0.93	Gamma Regressor	0.40
9	Elastic Net CV	0.92	AdaBoost Regressor	0.92	Linear SVR	0.93	SGD Regressor	0.36
10	RANSAC Regressor	0.92	Lars	0.92	RANSAC Regressor	0.93	Random Forest Regressor	0.36

264
 265 *Table 3. Ranking of the first ten algorithms with the best coefficient of determination according*
 266 *to the analysis of the 35 initially selected.*

267
 268 Since the GBRT selected models are decision tree-based algorithms and these are fairly
 269 insensitive to the scale of the features, the developed models were compared with and without

270 rescaling. The results without scaling of features suggested that it had virtually no impact on the
 271 GBRT results.

272

273 4.2. Bayesian optimization

274 With the four GBRT models previously selected, a Bayesian optimization process was performed
 275 to find the hyperparameters that give the models better predictive accuracy. In this Bayesian
 276 searching process, the R^2 was chosen as the reference metric. The cross-validation folds were
 277 created by a stratified fold with a splitting number of 5 and an iteration number of 70. With this
 278 process, the hyperparameters were established using cross-validation on the training set and the
 279 predictions were performed on the test set. The values of the hyperparameters obtained in this
 280 optimization process for each of the approaches taken into consideration are shown in Table 4.

281

Hyperparameter	Approach 1	Approach 2	Approach 3	Approach 4
Learning rate	0.07427	0.01259	0.00194	0.00118
Maximum depth	2	4	2	1
Number of estimators	1013	894	1881	1830
Number of features	None	None	None	2

282 *Table 4. Values of the hyperparameters in the final GBRT models obtained by Bayesian*

283 *optimization for each of the approaches considered.*

284 4.3. Performance analysis

285 Following the fine-tuning of the hyperparameters, the selected models with the best performance
 286 in each one of the approaches considered were analysed with the test dataset. In Table 5 is
 287 included a summary of the performance metrics obtained in the training and test sets of the
 288 definitive tuned models.

289

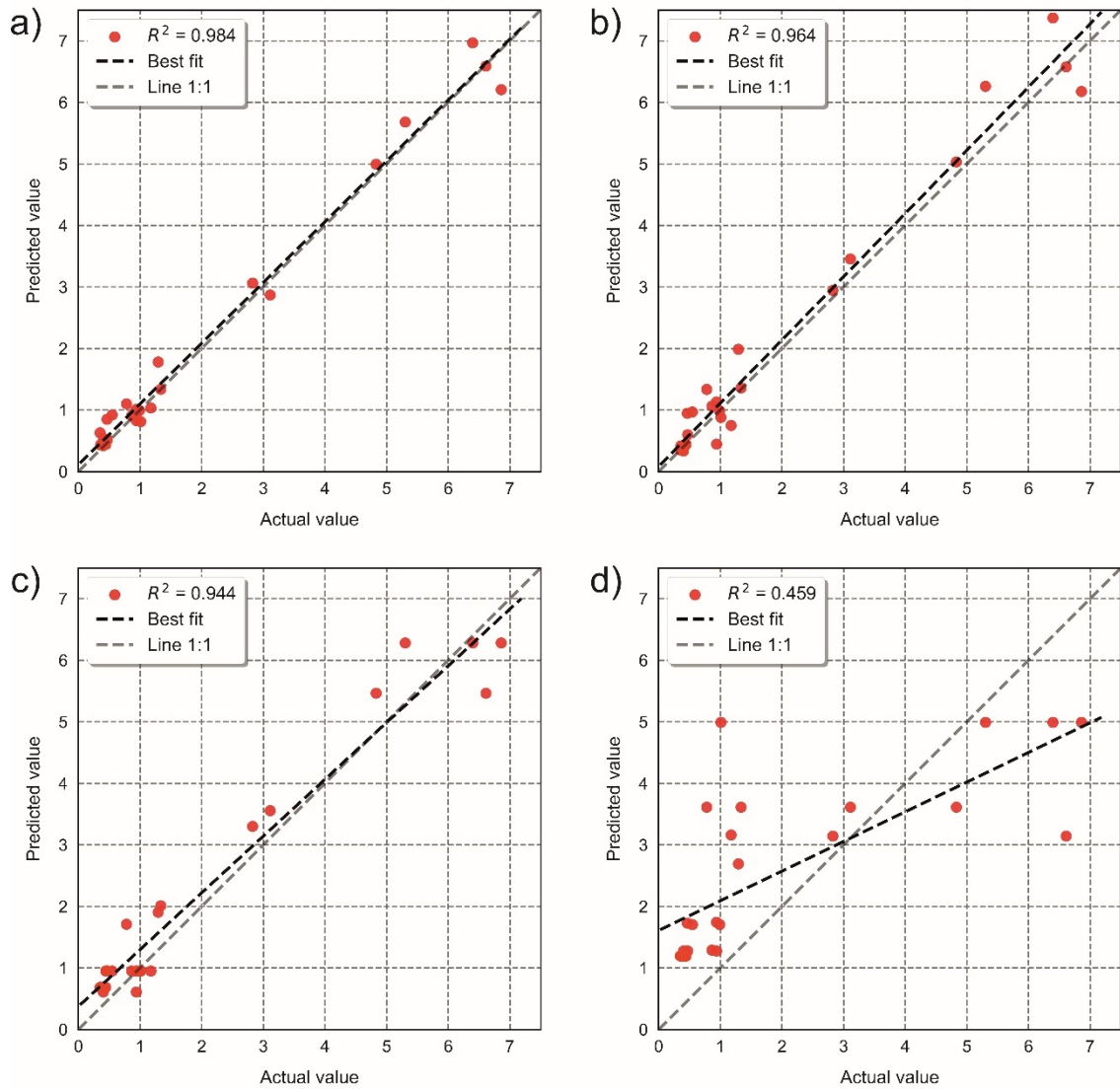
Approach	R^2		MAE		MSE	
	Training	Test	Training	Test	Training	Test
1	0.999	0.984	0.02	0.19	0.01	0.08
2	0.999	0.964	0.02	0.28	0.01	0.16
3	0.934	0.944	0.40	0.42	0.34	0.26
4	0.667	0.459	1.01	1.26	1.71	2.45

290 *Table 5. Summary of accuracy parameters in the train and test sets of the models for each*
291 *approach.*

292

293 Figure 4 examines the relationships between the models' predictions and the actual activity with
294 a linear regression analysis. For a perfect fit, all data should fall along a 1:1 line, as the model
295 outputs would be equal to the measured values. An excellent fit was obtained for the first two
296 approaches with values of $R^2 > 0.96$. A slightly lower level of accuracy is obtained in the approach
297 3 although it still gets good metrics ($R^2 > 0.94$). Approach 4 does not offer good results and it can
298 be concluded that there is no reliable way to predict activity from CEC and SSA because less than
299 half of the variance in the outcome variable is explained by the model ($R^2 = 0.46$).

300 On the other hand, from the analysis of the MAE values in the test set, the proposed GBRT
301 deviates, on average, from the predictions by ± 0.19 , ± 0.28 , ± 0.42 and ± 1.26 respectively for each
302 of the respective approaches. It should be noted when analyzing these values that approximately
303 40% of the test set has activity values exceeding 2.5, reaching up to values of nearly 7 (Figure 4).
304 The values of the first three approaches show accurate models, but in the fourth approach the
305 value is too high, showing a notable dispersion in the prediction of activity. Approach 1 yielded
306 the best predictive performance and includes as predictor variables clay content (directly related
307 to activity) and LL (associated with activity through PI). Nevertheless, Approach 2, which no
308 longer includes clay content, and Approach 3 (using only LL) also presented outstanding
309 performance metrics. This is an important observation, as this algorithm enables the prediction of
310 soil activity solely based on LL without significant errors, or even minimizing these errors if SSA
311 and CEC are additionally available alongside LL, rendering further soil properties unnecessary.

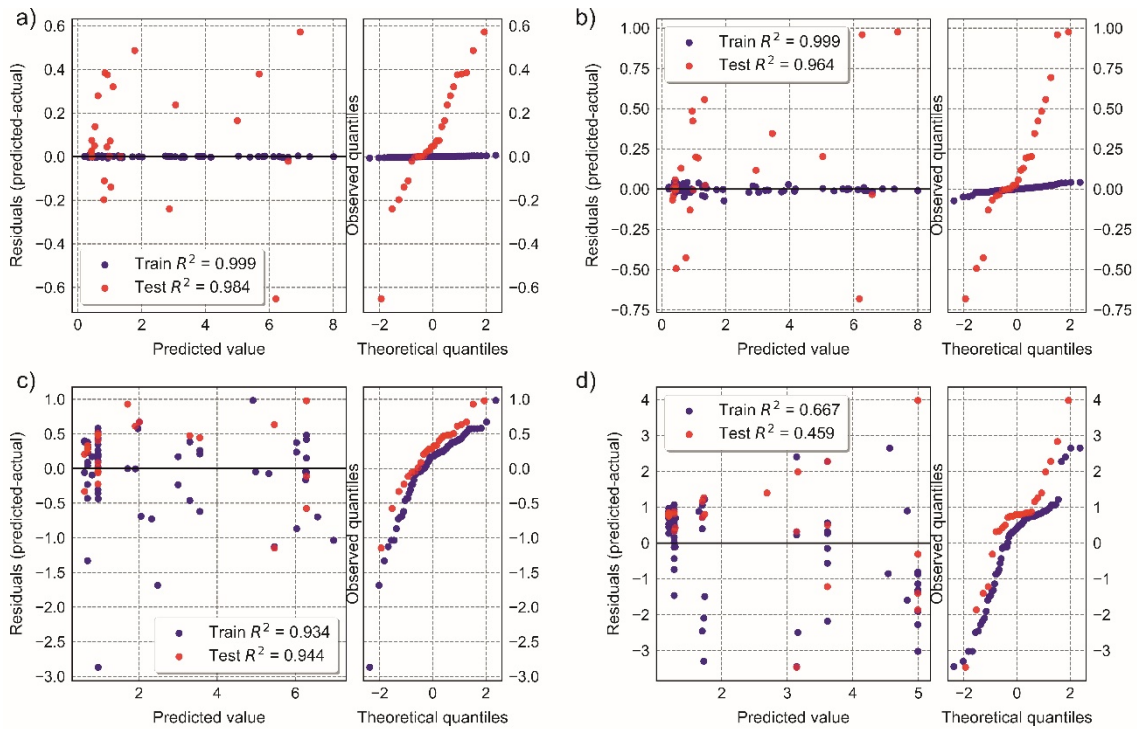


312

313 *Figure 4. Performance of the selected models for the training and test datasets in each of the*
 314 *approaches considered, a) activity from LL, clay content, SSA and CEC, b) activity from LL,*
 315 *SSA and CEC, c) activity from LL, and d) activity from SSA and CEC.*

316

317 Figure 5 presents the residuals of the GBRT model predictions for each of the four approaches,
 318 which are calculated as the difference between predicted and observed values. Additionally, the
 319 figure showcases the quantile-quantile (Q-Q) plots for these residuals. A Q-Q plot contrasts the
 320 quantiles of a given dataset with the quantiles of a theoretical probability distribution, in this case,
 321 the Gaussian distribution.



322

323 *Figure 5. Residuals and Q-Q plots of each of the approaches considered, a) activity from LL,*

324 *clay content, SSA and CEC, b) activity from LL, SSA and CEC, c) activity from LL, and d)*

325 *activity from SSA and CEC.*

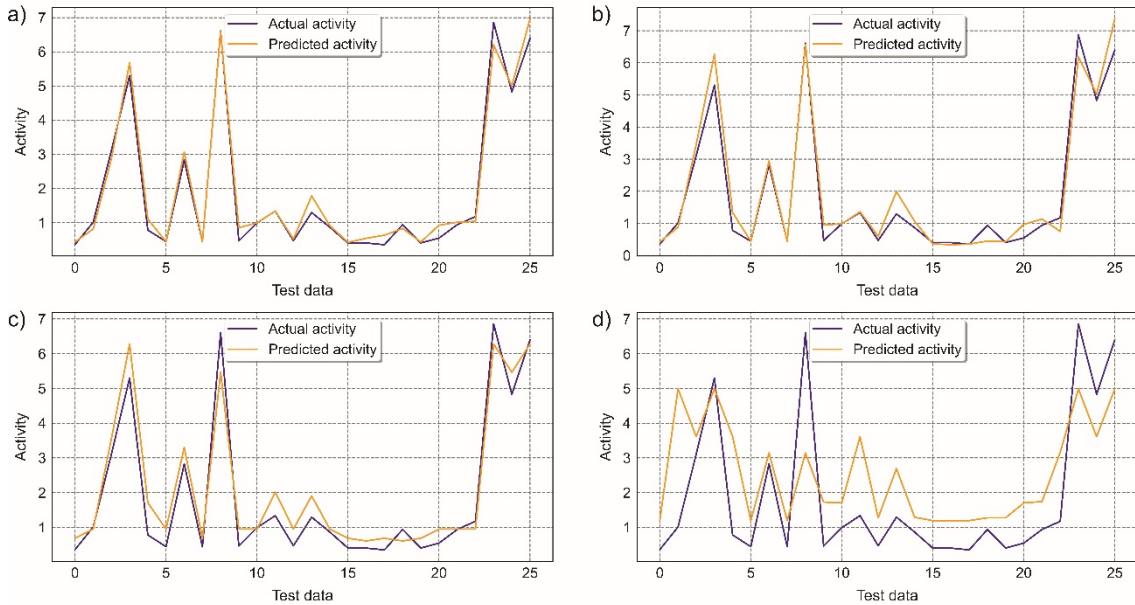
326

327 From examining the Q-Q plots, it can be inferred that in the test dataset, the residuals follow a
 328 distribution close to Gaussian although with a slight tendency towards positive residuals in the
 329 first two approaches. In approaches 3 and 4, the residuals move away from the Gaussian
 330 distribution, showing a light negative skewed distribution. These trends to Gaussian distributions
 331 are desirable since the algorithms, on average, predict the values with low error, and there are no
 332 extreme deviations in the predictions.

333 The evolution of the predicted activity in the test dataset for the four approaches considered is
 334 shown in Figure 6, compared with the measured values. Approaches 1, 2, and 3 accurately capture
 335 the evolution of actual values. In these 3 approaches, there are no areas where the algorithm tends
 336 to have less reliable predictions, both in the maximum and in the minimum values of activity, the
 337 difference between the predicted and the actual values is relatively small. Instead, in approach 4,

338 the predictive power is less, and the algorithm seems not to capture well the prediction of the high
 339 values and it can fail the prediction to identify large values of activity.

340



341

342 *Figure 6. Activity prediction results for the test dataset in each of the approaches considered, a)*
 343 *activity from LL, clay content, SSA and CEC, b) activity from LL, SSA and CEC, c) activity from*
 344 *LL, and d) activity from SSA and CEC.*

345

346 4.4. Variable importance in the prediction of activity

347 The feature importance is an important tool for the model interpretability, providing an evaluation
 348 of the predictive capacity of the input variables which can help to know the contributions of these
 349 to the output of the model. The trained GBRT models can automatically calculate feature
 350 importance, which can be obtained through the Gini Importance or mean decrease in impurity
 351 (Breiman et al. 1984). This method determines each feature importance as the sum over the
 352 number of splits in all the trees that include a specific feature, proportionally to the number of
 353 samples it splits.

354 Figure 7 shows the acquired importance order of the input variables in the approaches considered,
 355 taking into consideration that the approach 3 only considers one variable, so this analysis is not
 356 necessary. A higher value compared to another means greater importance of a feature for making

357 a prediction. In the developed GBRT models, LL is clearly the most important feature variable in
358 the approaches 1 and 2, and the rest of the variables have a less impact on the prediction of
359 activity. This is an interesting observation, as in approach 1 our model gives greater importance
360 in the prediction to LL than to clay content, which is a variable that by definition is proportional
361 to the value of activity. However, LL is also related to plasticity index and therefore to activity,
362 and this strong relationship has already been shown in numerous works (e.g. Spagnoli and
363 Shimobe 2019). In the approach 4, SSA has a significant impact in the prediction of activity,
364 while the contribution of CEC is small.

365

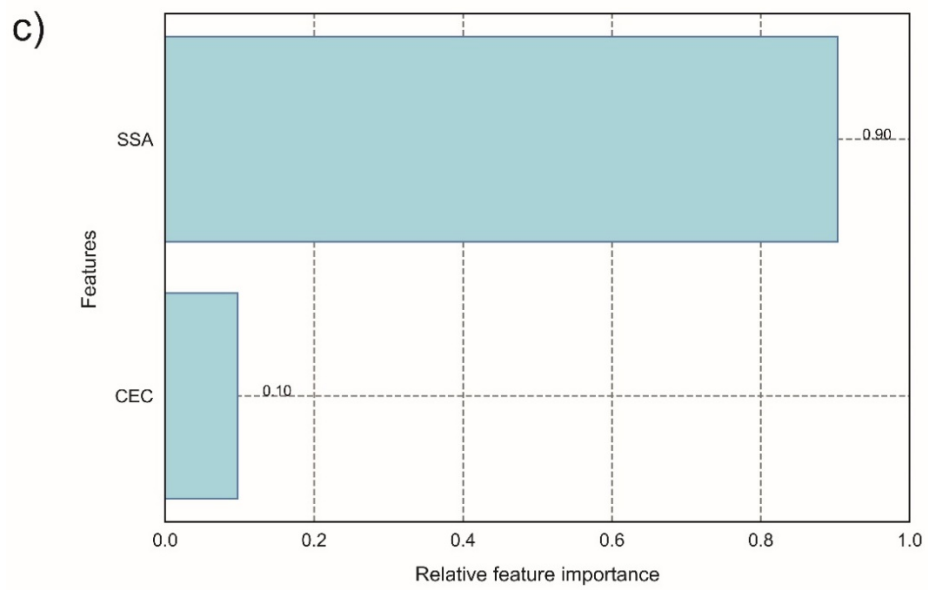
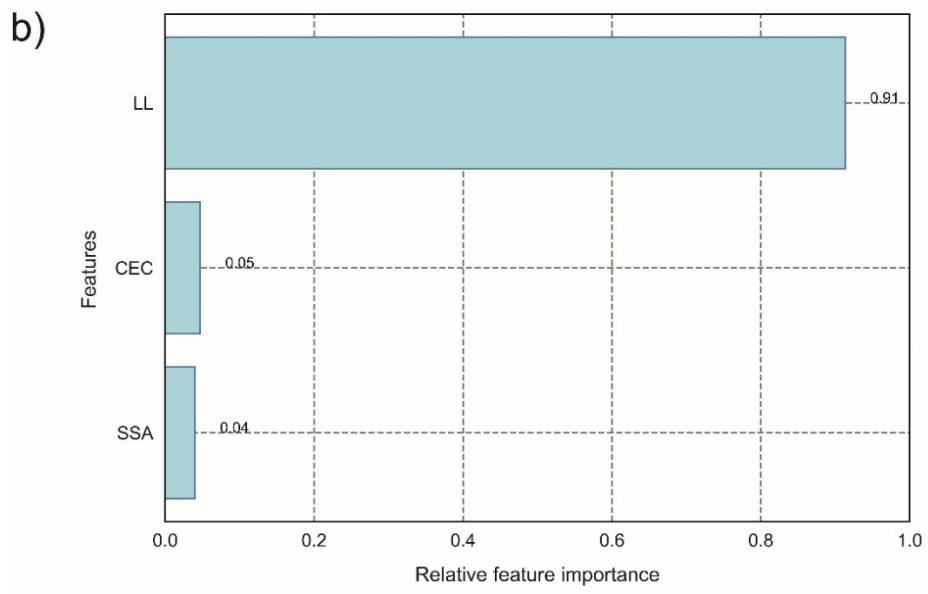
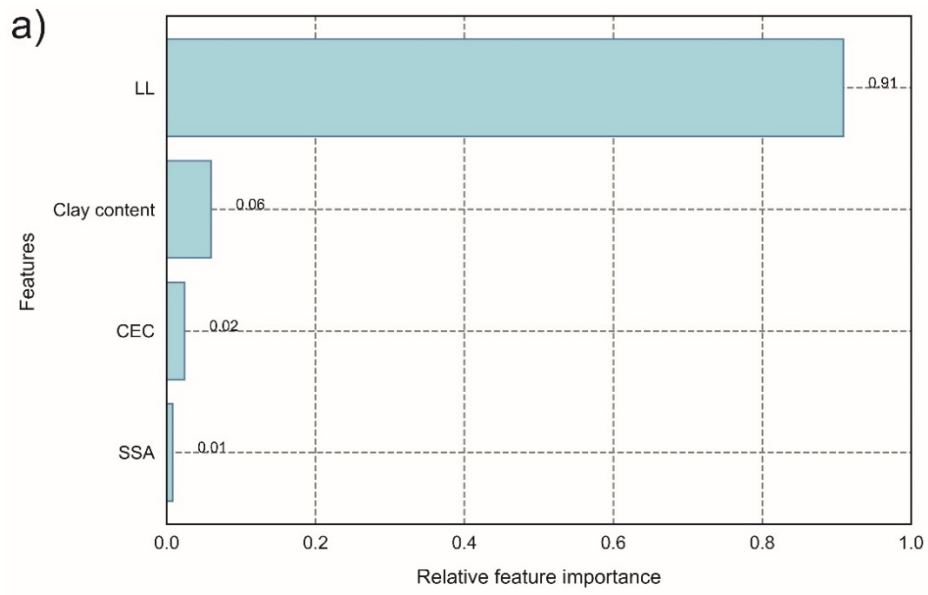
366

367

368

369

370



372 *Figure 7. Importance ranking of the input parameters in the activity prediction for each of the*
373 *approaches considered, a) activity from LL, clay content, SSA and CEC, b) activity from LL,*
374 *SSA and CEC and c) activity from SSA and CEC.*

375

376 **5. Conclusions.**

377 In geotechnical engineering, dealing with expansive soils is crucial due to their challenging
378 swelling and shrinking behaviors. This study introduces a novel approach, utilizing GBRT to
379 predict activity in clayey soils, a property related to swelling potential. Such predictive
380 capabilities have direct applications in pre-anticipating soil behavior, leading to safer and more
381 efficient construction and infrastructure planning. To achieve this aim four approaches to predict
382 activity were considered using different input variables: 1) LL, clay content, SSA and CEC, 2)
383 LL, SSA and CEC, 3) LL, and 4) SSA and CEC. This research paper has offered several key
384 insights and contributions:

385 1. GBRT outperformed the 35 algorithms that were initially selected and evaluated. These
386 algorithms were subsequently tuned using a Bayesian optimization process, obtaining the
387 definitive algorithms for the four approaches considered.

388 2. In the first three approaches, models yield high prediction accuracy. Approach 1, with its
389 inclusion of clay content (linked to activity) and LL (related via PI), had the highest predictive
390 accuracy ($R^2 = 0.98$).

391 3. The notable performance of Approaches 2 and 3 suggests that soil activity can be reliably
392 predicted using only LL ($R^2 = 0.94$), with even greater precision when combined with SSA and
393 CEC ($R^2 = 0.96$).

394 4. The fourth approach, i.e. predicting activity from SSA and CEC, did not show satisfactory
395 results in any of the models analysed, obtaining an R^2 value of 0.46, being able to conclude that
396 there is no reliable way to predict activity from SSA and CEC.

397 5. The conducted feature importance analysis indicated that LL is the most influential variable in
398 predicting activity for Approaches 1 and 2, with other variables having a lesser impact. In
399 Approach 4, SSA is the primary contributor to the prediction of activity.

400 Finally, the potential for scalability and adaptability of the proposed algorithms might be subject
401 to further detailed investigation. By incorporating a larger dataset, the proposed algorithms could
402 be refined, potentially extending the relevance of the current work. This could also create the way
403 for accounting for the effects of additional variables, marking a direction for future research.

404

405 **Availability statement**

406 The data supporting the research results can be accessed at
407 <https://huggingface.co/datasets/EstebanDC/activity-clays> or can be obtained by contacting the
408 corresponding author.

409

410

411

412 **References**

- 413 Al-Rawas, A.A. 1999. The factors controlling the expansive nature of the soils and rocks of
414 northern Oman. *Engineering Geology*, **53**, 327-350.
- 415 Arifin, Y.F. 2008. *Thermo-hydro-mechanical behavior of compacted bentonite-sand mixtures: an*
416 *experimental study*. <https://nbn-resolving.org/urn:nbn:de:gbv:wim2-20081006-14288>. doi:
417 <https://doi.org/10.25643/bauhaus-universitaet.852>.
- 418 Breiman, L.I., Friedman, J.H., Olshen, R.A. & Stone, C.J. 1984. Classification and regression trees
419 (cart). *Encyclopedia of Ecology*, **40**, 582-588.
- 420 Cerato, A. 2001. Influence of specific surface area on geotechnical characteristics of fine-grained
421 soils. Unpublished MSc Thesis, Department of Civil and Environmental Engineering, University
422 of Massachusetts.
- 423 Cerato, A.B. & Lutenegeger, A.J. 2002. Determination of surface area of fine-grained soils by the
424 ethylene glycol monoethyl ether (EGME) method. *Geotechnical testing journal*, **25**, 315-321.
- 425 Cerato, A.B. & Lutenegeger, A.J. 2004. Determining intrinsic compressibility of fine-grained soils.
426 *Journal of Geotechnical and Geoenvironmental Engineering*, **130**, 872-877.
- 427 Cerato, A.B. & Lutenegeger, A.J. 2005. Activity, relative activity and specific surface area of fine-
428 grained soils. *Proceedings of the International Conference on Soil Mechanics and Geotechnical*
429 *Engineering*. AA Balkema Publishers, 325.
- 430 Chittoori, B. & Puppala, A.J. 2011. Quantitative estimation of clay mineralogy in fine-grained
431 soils. *Journal of Geotechnical and Geoenvironmental Engineering*, **137**, 997-1008.
- 432 Christidis, G.E. 1998. Physical and chemical properties of some bentonite deposits of Kimolos
433 Island, Greece. *Applied Clay Science*, **13**, 79-98.
- 434 Devore, J.L., Farnum, N.R. & Doi, J.A. 2013. *Applied statistics for engineers and scientists*.
435 Cengage Learning.
- 436 Díaz, E., Brotons, V. & Tomás, R. 2018. Use of artificial neural networks to predict 3-D elastic
437 settlement of foundations on soils with inclined bedrock. *Soils and Foundations*, **58**, 1414-1422,
438 doi: <https://doi.org/10.1016/j.sandf.2018.08.001>.
- 439 Díaz, E., Pastor, J., Rabat, Á. & Tomás, R. 2021. Machine learning techniques for relating liquid
440 limit obtained by Casagrande cup and fall cone test in low-medium plasticity fine grained soils.
441 *Engineering Geology*, **294**, 106381.
- 442 Díaz, E. & Tomás, R. 2021. Upgrading the prediction of jet grouting column diameter using deep
443 learning with an emphasis on high energies. *Acta Geotechnica*, **16**, 1627-1633, doi:
444 <https://doi.org/10.1007/s11440-020-01091-8>.
- 445 Eriksson, D., Pearce, M., Gardner, J., Turner, R.D. & Poloczek, M. 2019. Scalable global
446 optimization via local bayesian optimization. *Advances in neural information processing*
447 *systems*, **32**.
- 448 Freund, Y. 1995. Boosting a weak learning algorithm by majority. *Information and computation*,
449 **121**, 256-285.
- 450 Freund, Y. & Schapire, R.E. 1996. Experiments with a new boosting algorithm. *icml*. Citeseer,
451 148-156.
- 452 Friedman, J.H.J.A.o.s. 2001. Greedy function approximation: a gradient boosting machine. 1189-
453 1232.
- 454 Gelman, A. & Hill, J. 2006. *Data analysis using regression and multilevel/hierarchical models*.
455 Cambridge university press.
- 456 Head, T., Kumar, M., Nahrstaedt, H., Louppe, G. & Shcherbatyi, I. 2020. *scikit-optimize/scikit-*
457 *optimize: v0. 8.1*. Zenodo.
- 458 Jones, D.R. 2001. A taxonomy of global optimization methods based on response surfaces.
459 *Journal of global optimization*, **21**, 345-383.
- 460 Jones Jr, D.E. & Holtz, W.G. 1973. Expansive soils-the hidden disaster. *Civil engineering*, **43**.
- 461 Kar, K.K. 2021. *Handbook of Fly Ash*. Butterworth-Heinemann.

462 Lambe, T.W. 1960. The Character and Identification of Expansive Soils: A Report Completed for
463 the Technical Studies Program of the Federal Housing Administration. Federal Housing
464 Administration.

465 Li, D., Liu, Z., Xiao, P., Zhou, J. & Jahed Armaghani, D. 2022. Intelligent rockburst prediction model
466 with sample category balance using feedforward neural network and Bayesian optimization.
467 *Underground Space*, **7**, 833-846, doi: <https://doi.org/10.1016/j.undsp.2021.12.009>.

468 Low, P.F. 1987. Structural component of the swelling pressure of clays. *Langmuir*, **3**, 18-25.

469 Marcial, D. 2013. Measuring water retention properties of a series of bentonite clays a wide
470 range of suctions. *Advances in unsaturated soils*. Taylor Francis Group, London, 135-140.

471 Mishra, A.K., Ohtsubo, M., Li, L.Y. & Higashi, T. 2012. Influence of various factors on the
472 difference in the liquid limit values determined by Casagrande's and fall cone method.
473 *Environmental Earth Sciences*, **65**, 21-27.

474 Mitchell, J.K. & Soga, K. 2005. *Fundamentals of soil behavior*. John Wiley & Sons New York.

475 Moore, D.S. & McCabe, G.P. 1989. *Introduction to the Practice of Statistics*. WH Freeman/Times
476 Books/Henry Holt & Co.

477 Muhunthan, B. 1991. Liquid limit and surface area of clays. *Geotechnique*, **41**, 135-138.

478 Nelson, J.D., Chao, K.C., Overton, D.D. & Nelson, E.J. 2015. *Foundation engineering for expansive*
479 *soils*. John Wiley & Sons.

480 Peck, R.B., Hanson, W.E. & Thornburn, T.H. 1974. *Foundation engineering*. John Wiley & Sons.

481 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
482 Prettenhofer, P., Weiss, R. & Dubourg, V. 2011. Scikit-learn: Machine learning in Python. the
483 *Journal of machine Learning research*, **12**, 2825-2830.

484 Phoon, K.-K. & Zhang, W. 2023. Future of machine learning in geotechnics. *Georisk: Assessment*
485 *and Management of Risk for Engineered Systems and Geohazards*, **17**, 7-22.

486 Polidori, E. 2009. Reappraisal of the activity of clays. *Activity chart*. *Soils and Foundations*, **49**,
487 431-441.

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

491 Santhoshkumar, T., Abraham, B., Sridharan, A. & Jose, B. 2016. Role of bentonite in improving
492 the efficiency of cement grouting in coarse sand. *Geotechnical Engineering Journal of the SEAGS*
493 *& AGSSEA*, **47**, 1-8.

494 Schapire, R.E. 1990. The strength of weak learnability. *Machine learning*, **5**, 197-227.

495 Schwing, M., Chen, Z., Scheuermann, A. & Wagner, N. 2013. Dielectric properties of a clay soil
496 determined in the frequency range from 1 MHz to 40 GHz. *Proc 10th Int. Conf. Electromagn.*
497 *Wave Interact. Water and Moist Substances*, 242-250.

498 Shahriari, B., Swersky, K., Wang, Z., Adams, R.P. & De Freitas, N. 2015. Taking the human out of
499 the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, **104**, 148-175.

500 Sivapullaiah, P., Prasad, B.G. & Allam, M. 2008. Methylene blue surface area method to correlate
501 with specific soil properties. *Geotechnical testing journal*, **31**, 503-512.

502 Skempton, A. 1953. The colloidal activity of clays. *Selected papers on soil mechanics*, **1**, 57-61.

503 Snoek, J., Larochelle, H. & Adams, R.P. 2012. Practical bayesian optimization of machine learning
504 algorithms. *Advances in neural information processing systems*, **25**.

505 Spagnoli, G., Feinendegen, M. & Rubinos, D. 2013. Modification of clay adhesion to improve
506 tunnelling excavation. *Proceedings of the Institution of Civil Engineers-Ground Improvement*,
507 **166**, 21-31.

508 Spagnoli, G. & Shimobe, S. 2019. A statistical reappraisal of the relationship between liquid limit
509 and specific surface area, cation exchange capacity and activity of clays. *Journal of Rock*
510 *Mechanics and Geotechnical Engineering*, **11**, 874-881.

511 Stephens, M.A. 1974. EDF statistics for goodness of fit and some comparisons. *Journal of the*
512 *American statistical Association*, **69**, 730-737.

513 Stone, M. 1974. Cross-validatory choice and assessment of statistical predictions. Journal of the
514 royal statistical society: Series B (Methodological), **36**, 111-133.

515 Wackerly, D., Mendenhall, W. & Scheaffer, R.L. 2014. Mathematical statistics with applications.
516 Cengage Learning.

517 Wang, L., Wu, C., Tang, L., Zhang, W., Lacasse, S., Liu, H. & Gao, L. 2020. Efficient reliability
518 analysis of earth dam slope stability using extreme gradient boosting method. Acta Geotechnica,
519 **15**, 3135-3150.

520 Warkentin, B. 1972. Use of the liquid limit in characterizing clay soils. Canadian Journal of Soil
521 Science, **52**, 457-464.

522 Yao, H.-l., Yang, Y. & Cheng, P. 2004. Standard moisture absorption water content of soil and its
523 testing standard. ROCK AND SOIL MECHANICS-WUHAN-, **25**, 856-859.

524 Zapata, C., Houston, S., Houston, W. & Dye, H. 2006. Expansion index and its relationship with
525 other index properties. *Unsaturated Soils 2006*, 2133-2137.

526 Zhang, W., Gu, X., Tang, L., Yin, Y., Liu, D. & Zhang, Y. 2022a. Application of machine learning,
527 deep learning and optimization algorithms in geoenvironment and geoscience: Comprehensive
528 review and future challenge. Gondwana Research, **109**, 1-17.

529 Zhang, W., He, Y., Wang, L., Liu, S. & Meng, X. 2023. Landslide Susceptibility mapping using
530 random forest and extreme gradient boosting: A case study of Fengjie, Chongqing. Geological
531 Journal.

532 Zhang, W., Wu, C., Zhong, H., Li, Y. & Wang, L. 2021. Prediction of undrained shear strength using
533 extreme gradient boosting and random forest based on Bayesian optimization. Geoscience
534 Frontiers, **12**, 469-477, doi: <https://doi.org/10.1016/j.gsf.2020.03.007>.

535 Zhang, W., Zhang, R., Wu, C., Goh, A.T. & Wang, L. 2022b. Assessment of basal heave stability
536 for braced excavations in anisotropic clay using extreme gradient boosting and random forest
537 regression. Underground Space, **7**, 233-241.

538 Zhang, Y., Qu, Y., Liu, G. & Wu, S. 2003. Engineering geological properties of Miocene hard clays
539 along the middle line of the North–South Diversion Water Project in China. Bulletin of
540 Engineering Geology and the Environment, **62**, 213-219.

541

542

543

544

545

546

547

548

549

550

551