

1 **Machine learning regression and classification algorithms utilised for strength**
2 **prediction of OPC/by-product materials improved soils.**

3
4 Eyo E. U.^{a*}, Abbey S. J.^a

5
6 (*Corresponding author Email: eyo.eyo@uwe.ac.uk)

7
8 ^a*Faculty of Environment and Technology, Department of Geography and Environmental Management, Civil*
9 *Engineering Cluster, University of the West of England, Bristol, United Kingdom.*

10
11 **Abstract**

12 In this study, stand-alone machine (ML) models (Bayesian regressor (BLR), least square
13 linear regressor (REG), artificial neural networks (ANN), and logistic regression (LR)), tree-
14 ensemble ML models (boosted decision tree (BDT), random decision forest (RDF) decision
15 jungle (DJ)) and meta-ensemble ML models (voting (VE) and stacking (SE)) are applied to
16 predict the strength of different soils improved by part-substitution of OPC with PFA and
17 GGBS in various combinations and proportions. Multiclass elements of these proposed ML
18 models are also deployed to provide analysis across multiple cross-validation methods.
19 Results of regression analysis indicated higher statistical variance of OPC-substituted
20 predictor variables compared to soils improved by OPC alone when using both stand-alone
21 and tree-based algorithms. On average, the REG model produced strength predictions with
22 higher accuracy (RMSE of 0.39 and R^2 of 0.86) compared to ANN (RMSE of 0.44 and R^2 of
23 0.82), but with comparatively lower accuracy compared to tree-based models (average RMSE
24 of 0.33 and R^2 of 0.90) and meta-ensemble models (average RMSE of 0.06 and R^2 of 0.91).
25 For ML classification, multiclass neural network algorithm (*mANN*) produced higher
26 accuracy (0.78), precision (0.67) and rate of recall (0.67) compared to tree-based models but
27 fell short to meta-ensemble models (average accuracy of 0.80, precision of 0.70 and recall of
28 0.71). Diagnostic tests across different validation methods indicated better performance of the
29 VE model compared to its SE ML counterpart when adopting the train-validation split
30 technique. Overall, the ensemble methods were more versatile on regression and multiclass
31 classification problems because they aggregated multiple learners to provide robust
32 predictions.

33
34 **Keywords:** Machine learning; cement; PFA; GGBS; OPC; Bayesian regressor; linear
35 regression; artificial neural networks; logistic regression; ensembles

List of abbreviations

ANN	Artificial neural network
ASTM	American system of testing for materials
AUC	Area under curve
BDT	Boosted decision tree.
BLR	Bayesian linear regression
BP	Back propagation
CEM I	Cement
CRISP-DM	CRoss industry standard process data mining
CV	Cross validation
DAG	Directed acyclic graph.
DJ	Decision jungle
EML	Extreme machine learning
FN	Functional networks
FPR	False positive rate
GA	Genetic algorithm
GGBS	Ground granulated blast furnace slag.
K-FCV	k-fold cross validation
KNN	k-Nearest neighbours
LR	Logistic regression
MAE	Mean absolute error.
MARS	Multivariate adaptive regression splines
MCCV	Monte Carlo cross validation
MGP	Multi-genetic programming
ML	Machine learning
OPC	Ordinary Portland Cement
PFA	Pulverised fuel ash
PI	Plasticity index
RDF	Random decision forest
REG	Linear regression
RMSE	Root mean square error.
ROC	Receiver operating characteristic
SE	Stacking ensemble
SVM	Support vector machine
TPR	True positive rate
TVS	Train-validation split.
UCS	Unconfined compressive strength
VE	Voting ensemble

43

44

45

46

47

48

49

50

51

List of symbols

α_m	neural network activities
Y	predictor variable
Y_n	normalised UCS influential factor
Y_0	raw UCS influential factor
m_y	mean of distribution
σ_y	standard deviation
N	dataset points
n_v	size of validation dataset
n_t	size of training dataset
X_n	independent variable
β_0	regression constant
t_m	sample target value
ε_m	additive noise
w	weight vector
β	Bayesian precision parameter
α	hyper-parameter controlling distribution
Σ	posterior variance
μ	mean of weights
w_{ij}	weight between two neurons
y_n	neural network output signal
x	activation of n th neuron
σ	neural network activation function
t	hypothesis test value
SE	standard error
X_m	mean of actual observations.
SS_{xx}	explained variation.

1. Introduction

The method of soil stabilisation that involves treatment with pozzolanic binders, continues to remain one of the most effective and economical means of ground improvement. The temptation to utilise Ordinary Portland Cement (CEM-I or OPC) as a traditional binder to stabilise weak soils seems nearly unavoidable in certain projects given their good hydraulic and binding qualities. However, the negative effects on the environment due to the continuous production and usage of OPC cannot be over-emphasized [1]. Hence, a complete replacement or part-substitution of OPC with relatively low-carbon secondary alternatives or by-products such as ground granulated blast furnace slag (GGBS) and pulverised fuel ash (PFA) has become inevitable in soil stabilisation [2–10].

Determination of the mechanical properties of stabilised soils based on some composite binder mixture parameters is undoubtedly an important first step towards the development of design mix guidelines for subsequent field application [11–13]. For a stabilised soil with multiple binder combinations, the challenges of establishing a property such as compressive strength may involve some time-consuming and laborious laboratory trial batching (soil-binder type quantities and optimum combinations), choice of curing duration, selection and testing of other related properties which can have the potential of affecting the target variable. Meanwhile, on the basis of theory, conventional models of forecasting the compressive strength of stabilised soils consist essentially of relationships that are developed empirically from statistical methods whereby, linear, and sometimes nonlinear regression techniques are applied [14,15]. The analytical equations generated through these models tend to determine unknown coefficients that affect the relationship of other variables and the compressive strength. These models, though effective in some cases, are riddled with shortcomings such as those associated with the complexities of the stabilised soil mentioned above.

In recent times, machine learning (ML) techniques have been introduced to compensate for the limitations of traditional methods of compressive strength prediction of soils [16,17]. However, the adoption of ML models for performance evaluation of improved ground properties has been very slow and only reported in few studies as follows: strength, dry density, moisture content additive content, resilient modulus modelling and prediction using artificial neural networks (ANN), support vector machines and regression (SVM & R), meta-ensembles (voting, stacking, tiering & bagging), functional networks (FN), multivariate adaptive regression splines (MARS), Logistic regression (LR), k-nearest neighbours (KNN), Genetic algorithm (GA), multi-genetic programming (MGP) [16–30]. These authors have used materials such as cement, lime, fly ash, fibres and geopolymers to strengthen the weak soils.

It is obvious that the application of multiple ML algorithms and a critical analysis that compares the relative performances of each one is not plentiful in literature. Moreover, an application of ML models to predict the unconfined compressive strength (UCS) of soils stabilised by partial substitution of OPC with cementitious by-products (PFA and GGBS) has not been done.

In this study, stand-alone ML models (Bayesian regression, linear regression, artificial neural networks, and logistic regression), tree-ensemble ML models (boosted decision tree, random decision forest and decision jungle) and meta-ensemble ML models (voting and stacking) are applied to investigate and predict the strength properties of five different soils stabilised by partial replacement of OPC with PFA and GGBS in various combinations and proportions.

Until recently, most ML predictions reported in literature have relied on conventional statistical metrics such as coefficient of determination and other standard error analyses for performance assessment [31–39]. However, this study aims to extend the scope of investigation, evaluation and prediction to include other diagnostic tests to support, confirm

101 and validate the commonly used statistical measures. This ensures that adequate sensitivity
102 analysis of performance is accounted for while also emphasising the effect of weight
103 variables or features in the prediction process.

104 Also, the complex combination of binders used in stabilisation of the soils considered in this
105 study can easily be regarded as a ML classification problem. Hence, the multiclass elements
106 of the proposed ML models are further deployed to provide analysis by considering multiple
107 cross-validation methods, an aspect which has not been considered in most previous studies.

108 The structure and framework of this study are as follows: a statement of the method involving
109 database construction, development and the experimental procedures adopted for soil
110 stabilisation is presented section in 2. This shall consist of a series of steps and processes
111 required for the preparation of collated datasets. Subsequent development and
112 implementation of the proposed models are given in section 3. Detailed analyses and
113 discussions of the performance of ML models are carried out in section 4. This section
114 includes an evaluation of ML regression, classification, and sensitivity analysis of the
115 prediction problems. In section 5, the significance of this study and recommendations for ML
116 model deployment are laid out. Following this is the concluding section where the main
117 points and highlights of the study are given.

118

119

2. Research Method

120

2.1. Database development and stabilisation procedure

121 High quality and original experimental data of unconfined compressive strength (UCS) tests
122 on soils stabilised using OPC and a combination of cementitious by-product materials were
123 compiled from literature and used to train ML models [14]. Soil-binder-water reactions under
124 a prescribed or natural curing environment can play a very significant role in the hardening
125 rate of the stabilised mixed. In this regard, 5 different soil types of varying initial properties
126 were improved by OPC, a blend of OPC and PFA and OPC-PFA-GGBS. OPC contents of
127 5%, 10%, 15% and 20% (by weight of dry soil) were applied to stabilise the weak soils. The
128 OPC-PFA-stabilised soils were first composed of a 50% reduction in the OPC content used.
129 The OPC was then further reduced to 33.33% and substituted by equal amounts of PFA and
130 GGBS to produce a stabilised soil of OPC-PFA-GGBS mixes. In all the soil mixtures, the
131 total proportion of the binder in the stabilised soil remained constant at 5%, 10%, 15% and
132 20%. The ratio of water-OPC in the stabilised soil mixture was unity. Three series of UCS
133 tests (ASTM D2166) were carried out following curing at 7, 14, 28 and 56 days to assess
134 strength developments of the soil with varying initial moisture and plasticity properties. The
135 first series was performed to study the influence of OPC addition alone on UCS while the
136 second and third series of tests were carried out to investigate the effect of OPC-PFA and
137 OPC-PFA-GGBS addition respectively. Compared to the natural soils, there was an
138 improvement in the UCS of the soils when stabilised by the binders details of which are given
139 in [Abbey et al. \[14\]](#). However, in keeping with the main goal of this research, the datasets
140 will be used to train ML models for UCS prediction.

141

142

2.2. Featurization and hyper-parameter optimisation

143 Stabilised soil datasets used in the supervised training of ML models in this study did not
144 contain any missing feature and attributes. Nevertheless, there was need for other forms of
145 feature engineering to be carried out to enable a reduction in unnecessary redundancy and
146 improvement of the integrity of data used for training of the ML models.

147

148 2.2.1. *Data normalization*

149 Without any distortion of the differences in the range of values of the UCS dataset, these
150 values were transformed into a common scale by using the Z-score standardisation method to
151 ensure outliers were avoided. The Z-score transformation is described mathematically as:

152

$$153 \qquad Y_n = \frac{Y_0 - m_y}{\sigma_y} \qquad (1)$$

154 Where Y_n and Y_0 represent the normalised and raw UCS influential factors respectively, while
155 m_y and σ_y denote corresponding values of the mean and standard deviation, respectively.

156

157 2.2.2. *Cross-validation (CV)*

158 For average sized datasets such as those used in this study, it was necessary to apply cross-
159 validation techniques to improve the reliability of the training sets and reduce the chances of
160 certain coincidental features receiving more importance [40]. Moreover, an overfitted model
161 is highly undesirable since it lessens the predictive performance on some “unseen” tested
162 data [41]. The following cross-validation techniques were used to optimise hyper-parameters:
163 k -fold cross-validation (k -FCV), Monte Carlo cross-validation (MCCV) and train-validation
164 split (T-VS) method.

165 The k -fold cross-validation (k -FCV) technique tends to divide the dataset (N) points into
166 some k - subsets of equal sizes. The process then treats one of the k -subsets as a training
167 subset and the remaining as validation subset. This process then repeats k number of times by
168 excluding one of the k - subsets in each cycle. In this study, 10-fold CV was adopted.

169 Monte Carlo cross-validation (MCCV) method splits the dataset (N) points into two subsets
170 by sampling and without replacement of one of the data points. The training is then
171 performed on the subset that was not replaced and validation on the replaced subset. Even
172 though there exist a rather unique training set, MCCV tends to avoid the need to run any form
173 of iterations unlike the k -FCV.

174 In summary, if we consider both MCCV and k -FCV and then assume N to represent the size
175 of the dataset, k denoting the number of k -fold subsets, n_v the size of validation set, and n_t
176 the size of training set then:

177

$$178 \qquad k\text{-FCV, } N = k \times n_v \qquad (2)$$

179

$$\qquad \qquad \qquad \text{MCCV, } N = n_t + n_v. \qquad (3)$$

180 Train-validation split (T-VS) is a simple method of randomised dataset splitting whereby
181 each of the subsets are used for training and testing purposes, respectively. In this research,
182 80% of the parent dataset were used to meticulously train the ML models (i.e., selection and
183 optimisation of hyper-parameters and functions) while the remaining 20% of the dataset were
184 used to test and assess the prediction performance of the ML models. Previous studies have
185 recommended that the testing data subset may not be less than 10% nor more than 30% of the
186 entire data records [32,40].

187
188
189
190

191 3. Machine learning models.

192 3.1. Multiple linear regression (REG)

193 This is one of the most common mathematical methods employed for supervised ML. The
194 least square function of REG technique establishes several correlations between one or more
195 independent or explanatory variables and predictor or dependent variables. Changes in the
196 predictor variable say, Y are often triggered by the nature of the independent variable X as the
197 following general equation shows:

$$198 Y_n = \beta_0 + \beta_1 X_{1n} + \beta_2 X_{2n} + \dots + \beta_m X_{mn} + \epsilon_n \quad (4)$$

200 Where Y_n = predictor variable; $X_{1n}, X_{2n}, \dots, X_{mn}$ = independent variables; β_0 = constant; and $\beta_1,$
201 β_2, \dots, β_m = coefficients of regression; and ϵ = error term.

203 3.2. Logistic regressor (LR)

204 This is a nonlinear model where the deviation or variance of the predictor variable is a
205 function of its mean. In other words, unlike REG, the value of the predictor variable depends
206 on the probability that it belongs to a certain class. LR tends to add an exponential function
207 on top of REG in order to restrain the predictor of response $Y_n \in [0,1]$, instead of $Y_n \in \mathcal{R}$ as in
208 REG [16]. The LR model for say ‘y’ distinct predictors could be represented as:

$$209 y(X) = P_r(Y = 1|X) = \frac{e^{\beta_0 + \beta_1 X_{1n} + \beta_2 X_{2n} + \dots + \beta_m X_{mn} + \epsilon_n}}{1 + e^{\beta_0 + \beta_1 X_{1n} + \beta_2 X_{2n} + \dots + \beta_m X_{mn} + \epsilon_n}} \quad (5)$$

211 Where $\Pr(Y=1|X)$ = probability that the response $Y = 1$ or 0 given $X_{1n}, X_{2n}, \dots, X_{mn} =$
212 independent variables, $X_{1n}, X_{2n}, \dots, X_{mn} =$ independent variables; β_0 = constant; and $\beta_1, \beta_2, \dots,$
213 β_m = coefficients of regression; and ϵ = error term. Because of its extreme validity in
214 classification as well as regression problems, LR unlike its REG counterpart is mostly
215 preferred as the default first option [40].

216 217 3.3. Bayesian linear Regressor (BLR)

218 This is a special case of linear regression whereby the model analysis is undertaken within
219 the context of a statistical inference of the “Bayes” theorem. This is then used to update the
220 probability of a given hypothesis as more information or evidence become available. Bayes
221 theorem describes the probability of an event taking place as a result of having prior
222 knowledge of certain conditions that might be related to such event. If it is considered that a
223 target value say, t_m , is sampled from an experiment then the relationship between this value
224 and the predictor variable $y(x_m; w)$ can be given as [42]:

$$225 t_m = y(x_m; w) + \epsilon_m \quad (6)$$

227 Where ϵ_m = an additive noise (modelled as Gaussian distribution $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ with a random
228 zero-mean variable); w = weight vector; β = precision parameter

229

230 Hence, following Gaussian distribution: we have the target value t_m as:

$$231 \quad p(t_m|\mathbf{w}, \beta) = \mathcal{N}(t_m|y(\mathbf{x}_m; \mathbf{w}), \beta^{-1}) \quad (7)$$

232
233 With the input parameter given as \mathbf{x} , the likelihood or probability for the target vector \mathbf{t} then:

$$234 \quad p(\mathbf{t}|\mathbf{w}, \beta) = \prod_{m=1}^K \mathcal{N}(t_m|y(\mathbf{x}_m; \mathbf{w}), \beta^{-1}) = \left(\frac{\beta}{2\pi}\right)^{N/2} \exp\left\{-\frac{\beta}{2} \|\mathbf{t} - \Phi\mathbf{w}\|^2\right\} \quad (8)$$

235 Where $\mathbf{x} = (x_1, x_2, x_3, \dots, x_N)^T$; $\mathbf{t} = (t_1, t_2, t_3, \dots, t_N)^T$.

236 To prevent any model complexity and over-fitting of the maximum likelihood directed
237 against \mathbf{w} , a prior distribution is then defined as follows:

$$238 \quad p(\mathbf{w}|\alpha) = \prod_{m=1}^L \mathcal{N}(w_m|0, \alpha^{-1}) \quad (9)$$

239 Where α = hyper-parameter controlling the distribution of w_m , $p(w_m|\alpha) = \mathcal{N}(w_m|0, \alpha^{-1})$.
240 the posterior distribution over \mathbf{w} can also be obtained as follows:

$$241 \quad p(\mathbf{w}|\mathbf{t}, \alpha, \beta) = \frac{p(\mathbf{t}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha)}{p(\mathbf{t}|\alpha, \beta)} = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (10)$$

242
243 Where $p(\mathbf{t}|\alpha, \beta)$ = normalising factor; $\boldsymbol{\Sigma}$ = posterior variance; $\boldsymbol{\mu}$ = mean of weights

244 *3.4. Tree-ensembles*

245 Tree-based models or ensemble of decision trees (Fig. 1) are a ML paradigm whereby formal
246 rules are obtained from detected patterns in the datasets hence, the tree-based models must be
247 trained in a rigorous manner on the data in order to be able to predict the properties presented
248 by a query [43]. Depending on the application, differences exist of how the tree-based ML
249 models are built and for the purposes of his study the random decision forest (RDF), boosted
250 decision trees (BDT) and decision jungles (DJ) shall be considered.

251

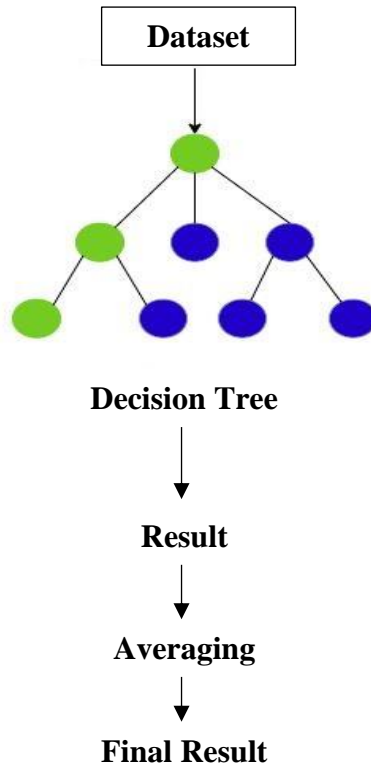


Fig. 1. Structure of a single regression (decision) tree.

3.4.1. *Random decision forest (RDF)*

This is a decision ensemble that is created to reduce the instability or fluctuations of single regression trees. It utilises the “bootstrap aggregation” (bagging) concept to generate various “bootstrap aggregation” (bagging) concept to generate various similar data records that are hitherto sampled from the same parent source. Bagging is simply a technique of aggregating a multiple tree models in a ‘bag’ for data training [37]. One of the disadvantages RDF is that it tends to fall victim of overfitting due to its small biases and wide variance.

3.4.2. *Boosted decision trees (BDT)*

Like RDF, these are an ensemble means of solving the problems of instability and poor performance of a single regression tree. In general, the idea of “boosting” is a strategy that is used to improve the performance of weaker learning regression tree algorithms. The step of boosted model building is often repeated through a set of iterations. Unlike RDF where all the trees are of equal importance, the BDT are rather hierarchical, and each tree layer is created recursively [41]. BDT tend to possess high performance ability especially on nonlinear datasets. However, one of the disadvantages of BDT is that its interpretability capacity is low and that makes it difficult to gain sufficient intuition of the patterns learned by the model.

276 3.4.3. Decision jungles (DJ)

277 These are a somewhat recent extension to RDF. A DJ is composed of an ensemble of rooted
 278 decision directed acyclic graphs (DAGs) as a method to obtain compact and accurate ML
 279 classifiers [44]. By permitting the merging of trees, a decision DAG traditionally has low
 280 memory footprint and are great at generalisation performance. DJs have the advantage that
 281 they are non-parametric ML models that can represent nonlinear decision boundaries. They
 282 are capable of selecting integrated features and performing classifications while also being
 283 very resilient to noisy features.

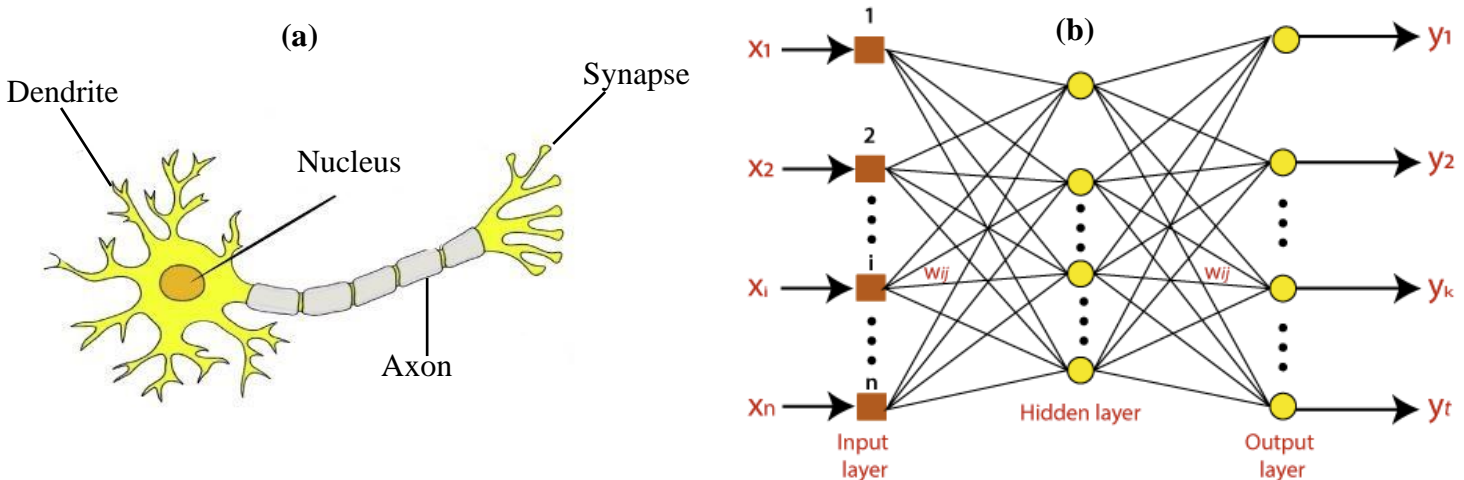
284
 285 3.5. Artificial Neural Network (ANN)

286 Just like the tree-ensemble ML models, ANN will have to be trained on a dataset to be able to
 287 predict the properties of a presented query. ANN are a family of data-processing ML models
 288 that are basically inspired by the human brain or neural networks whereby the neurons are
 289 interconnected through synapses (Fig. 2a) [45]. This network of neurons receive inputs,
 290 processes them and then make decisions or predictions. The neuron which is the processing
 291 element has the ability of filtering functions to make sure that inputted data to a specific node
 292 does not affect the network. The neuron also has an adaptive learning capability to adjust the
 293 weights that are connected between the nodes. ANN has a basic input layer, a hidden layer,
 294 and an output layer (Fig. 2b). When a processing neuron or element provides an input to
 295 another unit, the output is received as an input by the successive processing unit. This
 296 interconnectedness can be mathematically expressed as:

297
 298
$$\alpha_m = \sigma(\sum_j \omega_{ij}y_j), \quad \sigma(x) = \frac{1}{1+e^{-x}} \quad (11)$$

299 Where α_m = ANN activities; w_{ij} = weight between two neurons; y_n = output signal; x =
 300 activation of n th neuron; $\sigma(x)$ = activation function facilitating input transformation to output
 301 by multiplication of the inputs from the processing neuron by corresponding weights.

302
 303



304

305 **Fig. 2.** Neural network architecture (a) human neuron (b) artificial neural network (ANN).

306 3.6. Meta-ensembles

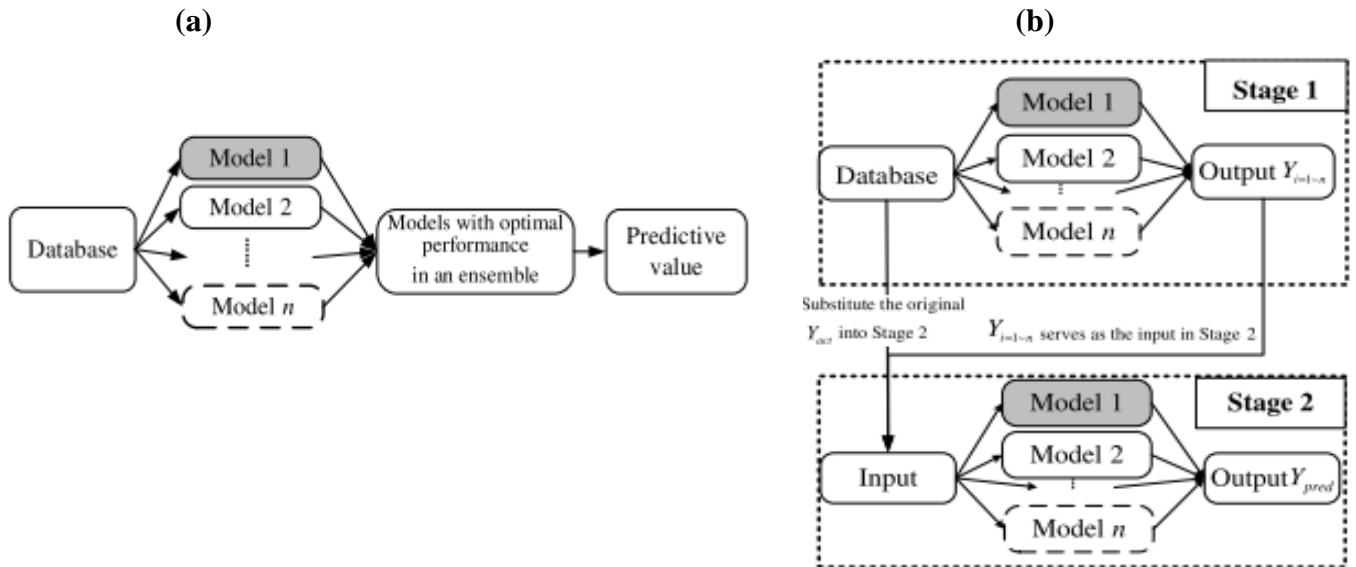
307 Compared to stand-alone and tree-based models stated in the forgoing, meta-ensembles (or
 308 model of models) are used to further improve the accuracy of prediction by combining some
 309 of the above-mentioned models. Both voting and stacking meta-ensembles are proposed for
 310 this study.

311
 312 3.6.1. Voting (VE)

313 When considering a regression problem, VE calculates the average predictions from the
 314 combined models. Meanwhile, for classification problems, predictions for each class label are
 315 added and the label that has a majority vote is predicted [24]. Fig. 3a depicts the VE ML
 316 model whereby the mean inputs are obtained from a combination of several models as the
 317 value of prediction.

318
 319 3.6.2. Stacking (SE)

320 This is an extension of the VE whereby the ML model learns how much and when to rely on
 321 each model to make generalised multistage predictions. The result of predictions obtained
 322 from say previous combined models ($X_{m=1\sim j}$) serves as input Y of the next stage as further
 323 predictions are being made (X_{pred}) as shown in Fig. 3b [24,46].
 324



325
 326 **Fig. 3.** Structure of meta-ensemble models (a) voting ensemble (b) stacking ensemble.

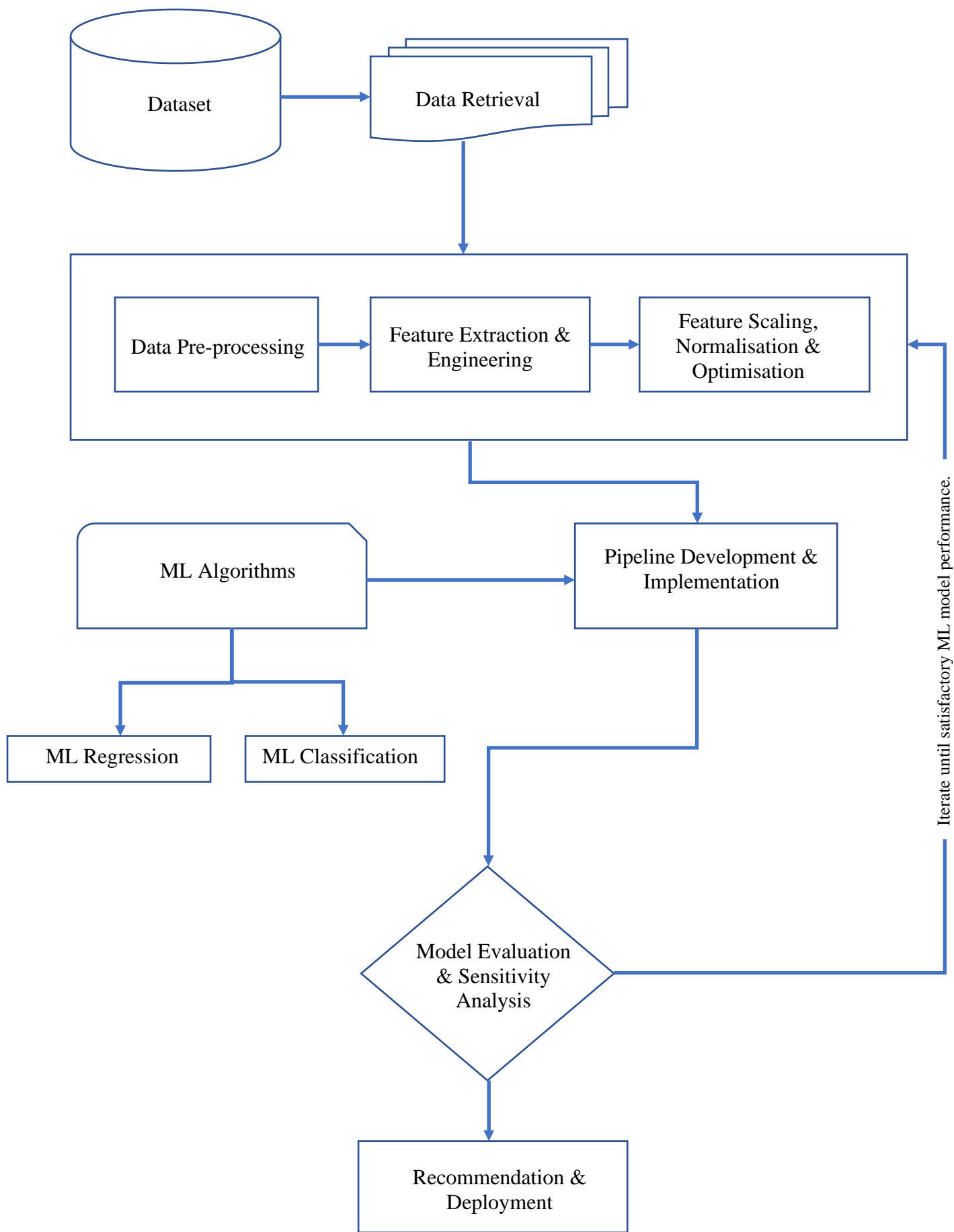
327 3.7. ML model pipeline development and implementation

328 Implementation of the algorithms was carried out on a designer platform which supports
 329 Python programming language (that includes the numpy, scipy and scikit-learn libraries) and
 330 ML pipeline developments. Considering the nature of the dataset used in this research, both

331 ML regression and classification were conducted to investigate and evaluate the performance
332 of the algorithms. Important features and parameters used in the ML models and datasets
333 respectively are given in [Table 1](#) and [Table 2](#). [Fig. 4](#) depicts the flowchart of the methodology
334 followed in developing the desired ML pipeline and subsequent evaluation of the models in
335 this study.

336

337



338

339

Fig. 4. Machine learning methodology flowchart.

340 **Table 1**
 341 Parameter settings of ML models

Model	Parameter	Option/value
REG	Method of solution	Ordinary least squares
	L2 regularization weight	0.001
LR	Trainer mode	Single parameter
	L1 regularization weight	1
	L2 regularization weight	1
BLR	Regularization weight	1
RDF	Resampling method	Bagging
	Trainer mode	Single parameter
	Number of decision trees	8
	Max. decision tree depth	32
	Number of random splits per node	128
	Min. sample no. per leaf node	1
BDT	Trainer mode	Single parameter
	Max. no of leaves per tree	20
	Training instances to form a tree	10
	Rate of learning	0.2
	Number of trees constructed	100
DJ	Resampling method	Bagging
	Trainer mode	Single parameter
	No. of DAGs	8
	DAGs max. depth	32
	DAGs max. width	128
	No. of optimization steps per DAG layer	2048
ANN	Hidden layer spec.	Fully connected
	Number of hidden nodes	100
	Rate of learning	0.005
	Number of learning iterations	100
	Initial learning weight diameter	0.1
	Normaliser	min-max

342

343

344

345

346 **Table 2**
 347 Data features and attributes

Feature	Attributes	Data type
Binder combinations	C	String
	C-PFA	
	C-GGBS-PFA	
Duration	7	Integer
	14	
	28	
	56	
Soil type	Soil 1	String
	Soil 2	
	Soil 3	
	Soil 4	
	Soil 5	
Binder quantity	5	Integer
	10	
	15	
Plasticity Index	20	Integer
	PI	

348

349 *3.8. Performance evaluation of ML models*

350 For an examination and evaluation of the precision of prediction and subsequent performance
 351 of the ML models studied in this research, three indicators are considered namely, Mean
 352 Absolute Error (MAE), Root Mean Squared Error (RMSE) and coefficient of determination
 353 (R^2). Detailed mathematical formulae of these performance criteria are given in literature
 354 [31,47]. Further assessment of the robustness and integrity of the regression models are
 355 considered by considering their prediction intervals. Prediction interval (Eq. 12) (unlike
 356 confidence interval) defines a range that may have a likelihood of containing the value of a
 357 dependent variable for a single future observation given some specific values of the
 358 independent variables.

359

$$360 \quad X = Y \pm t_{\alpha} SE \sqrt{1 + \frac{1}{n} + \frac{(X - X_m)^2}{SS_{xx}}} \quad (12)$$

361

362 Where Y = predicted values; t = hypothesis test value based on the percentage confidence; SE
 363 = standard error; n = size of dataset or number of observations; X_m = mean of actual
 364 observations ; SS_{xx} = explained variation.

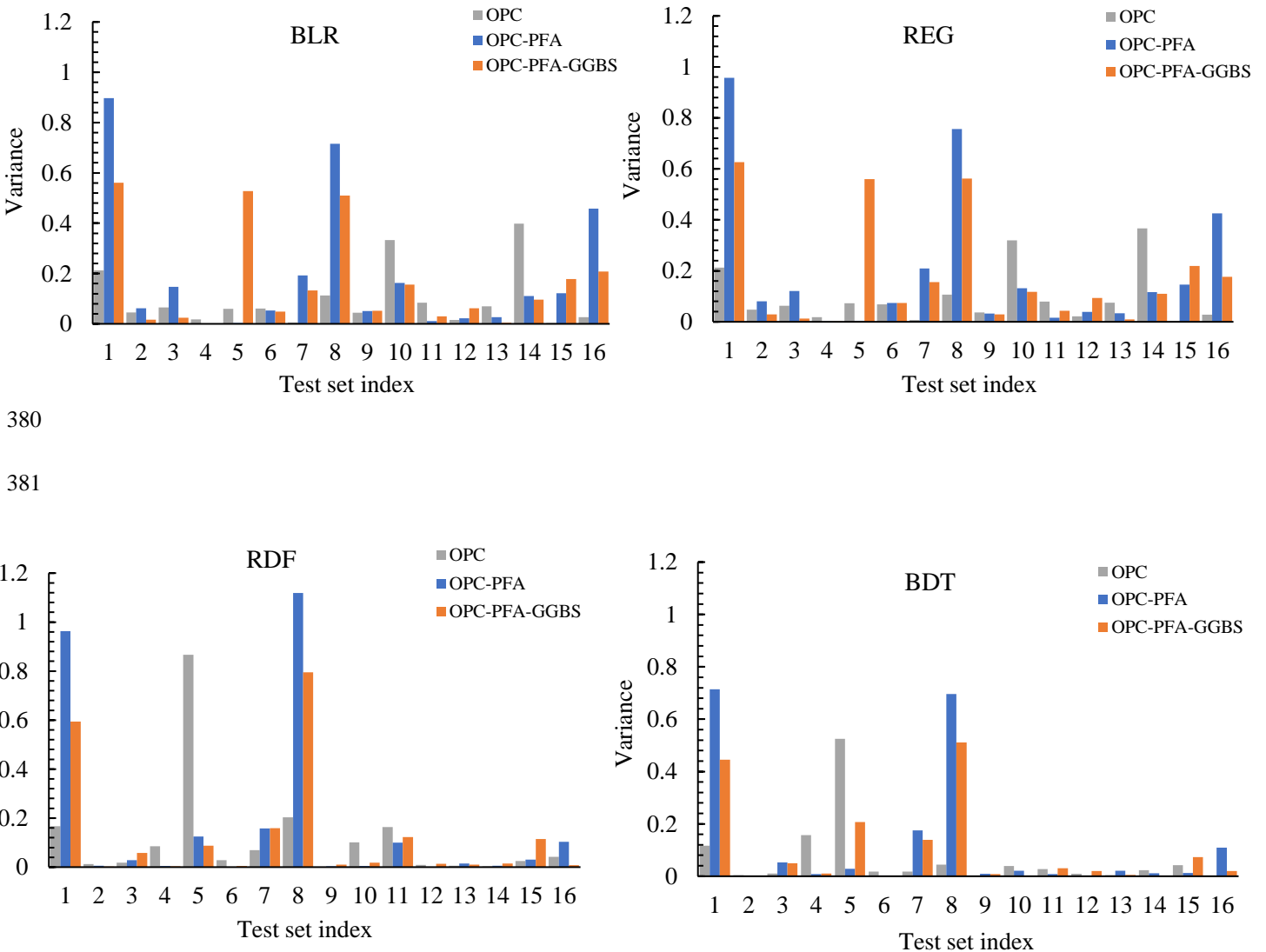
365

4. Results and discussion

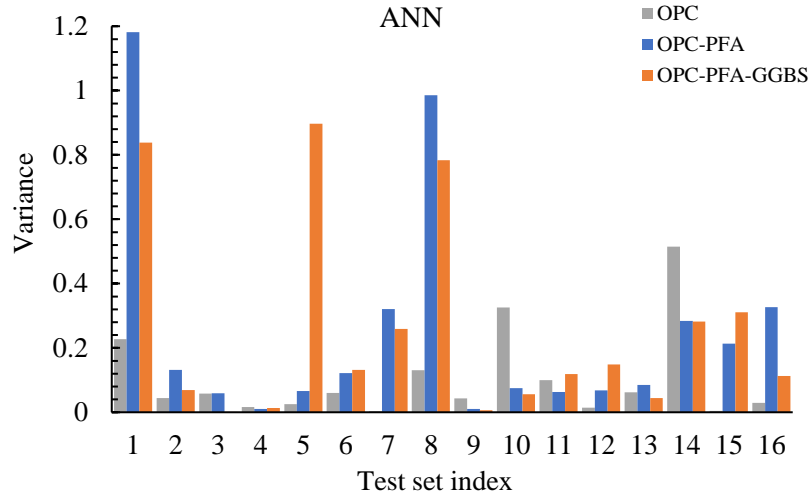
367 The ML models (stand-alone algorithms – REG, BLR, ANN; tree-based algorithms – BDT
 368 and RDF and the meta-ensemble models, VE and SE) proposed for regression analysis will
 369 be considered first in this section. Subsequently, analyses of classification problems using the
 370 multiclass ML models (LR, ANN, RDF, DJ and the meta-ensembles) will be given.

371 **Fig. 5** compares the statistical variance (from the mean of the target variable) of each
 372 predictor components (OPC, OPC-PFA and OPC-PFA-GGBS) for ML test data sets.
 373 Generally, across the algorithms tested on, higher deviations are experienced by both the
 374 OPC-PFA and OPC-PFA-GGBS predictor variables compared to the soil stabilised by OPC
 375 alone. However, compared to the other models, RDF seems to register the highest possible
 376 variance (about 0.87) followed by BDT (about 0.53) for the soil stabilised using only OPC.
 377 Further examination of the performance of the ML models and an analysis of the independent
 378 variables are given in the following sections.

379



382



383

384 **Fig. 5.** Statistical variance of predictor components of improved soils (a) BLR (b) REG (c)
 385 RDF (d) BDT (e) ANN

386 *4.1. ML Regression*

387 *4.1.1. Quality Assessment of ML regression models*

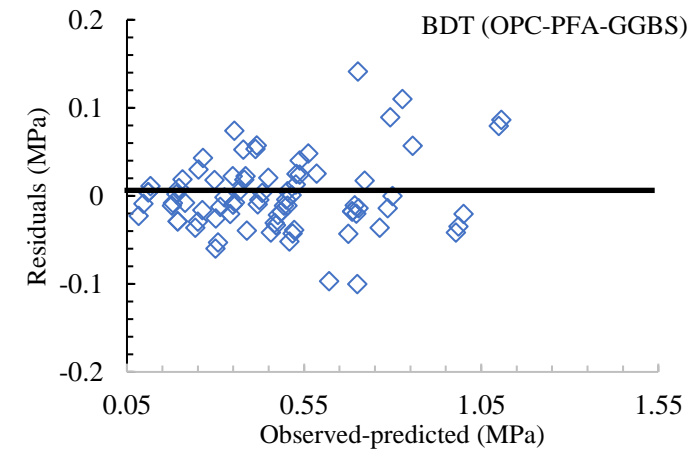
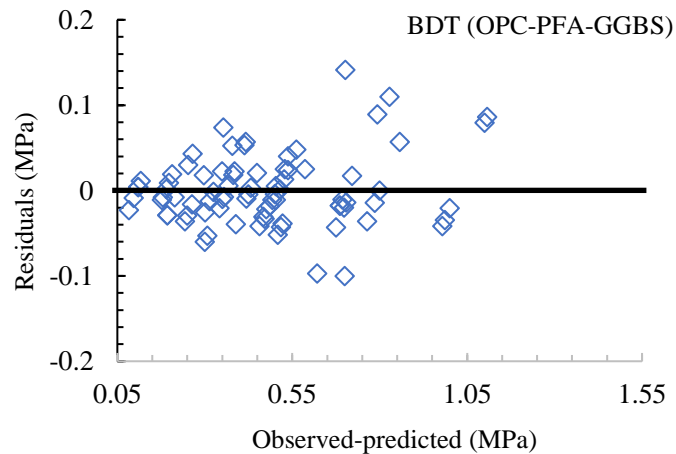
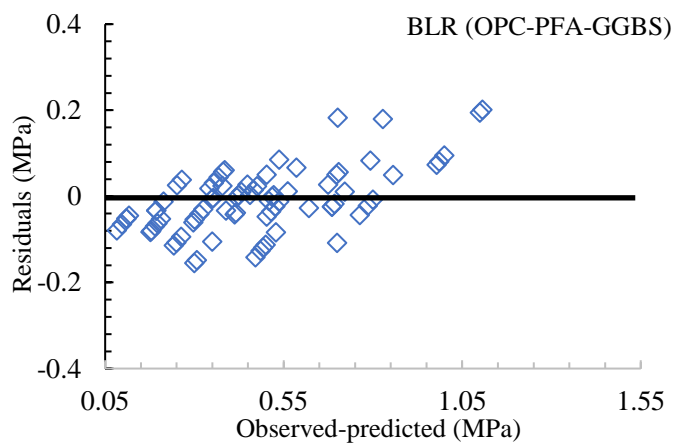
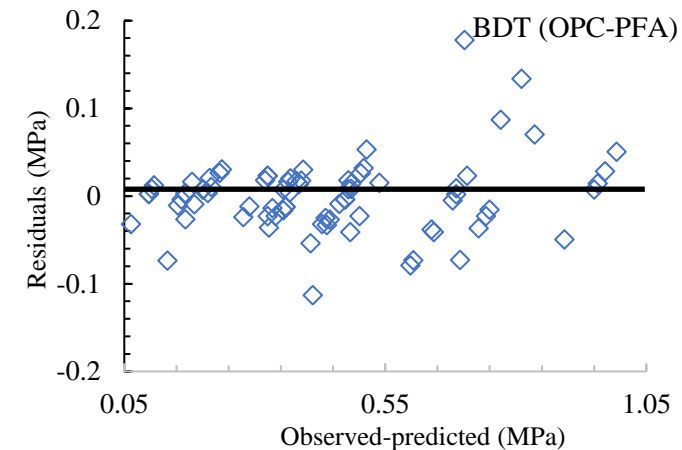
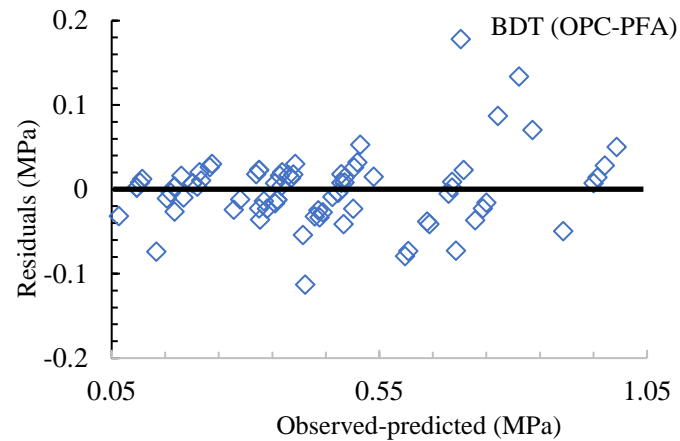
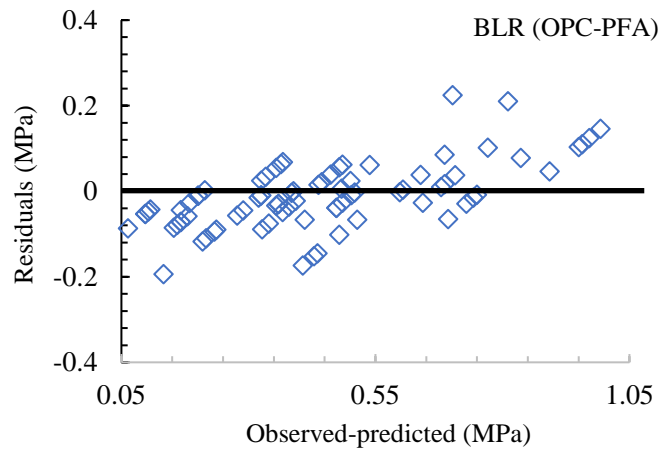
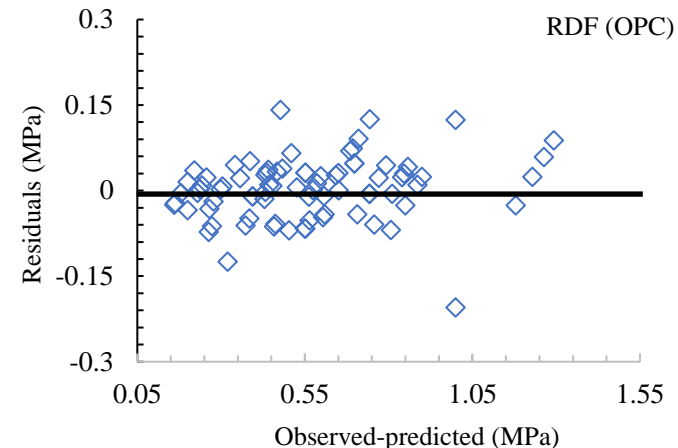
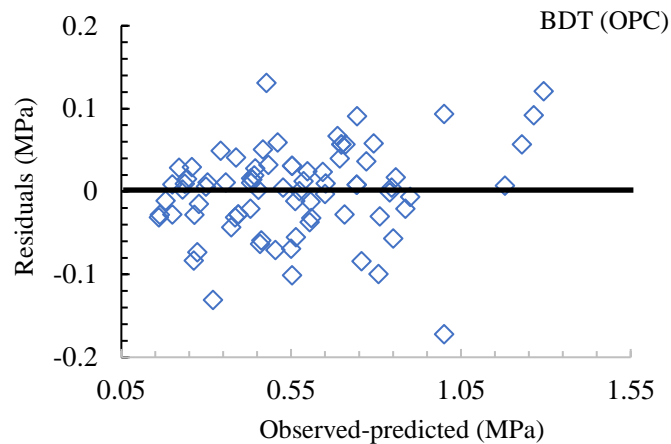
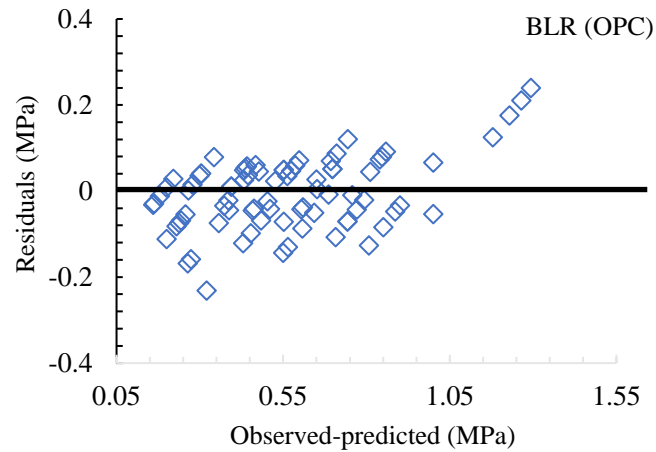
388 Residual lag plot provides the basis for an evaluation of the quality of the algorithms used to
 389 perform the regression analysis in this research. Furthermore, it does allow for an
 390 examination of any underlying statistical assumptions especially when considering the
 391 independence of features or variables and normality of distribution [48]. In order for any
 392 assumptions to hold true for a given regression model, the residuals will have to be
 393 distributed randomly around zero [49]. For a good model, the residual’s scatter plot will show
 394 a disorderly pattern of the data hence, without indicating any form of trends. In other words,
 395 if there are any forms of trends observed in the data, this will indicate that the residuals are
 396 not entirely independent.

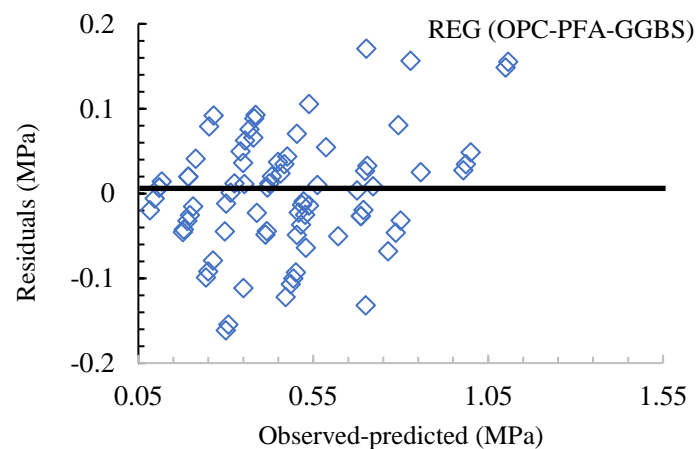
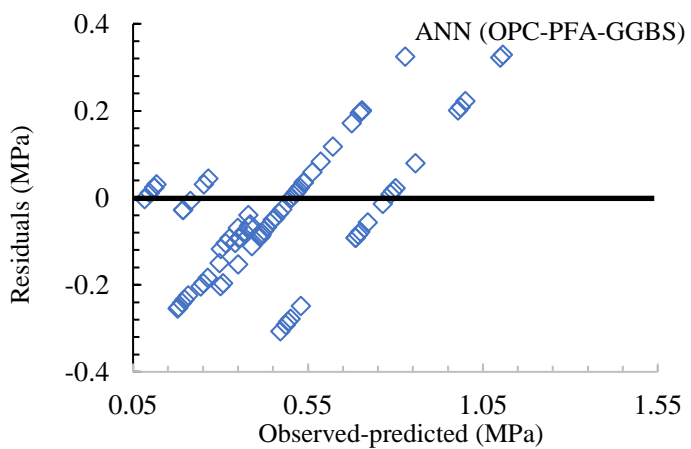
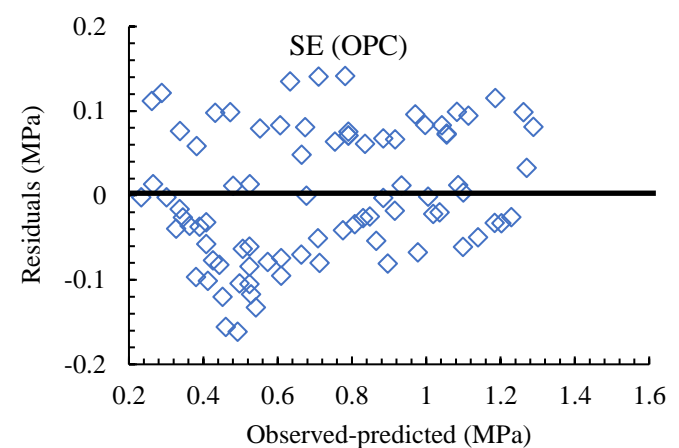
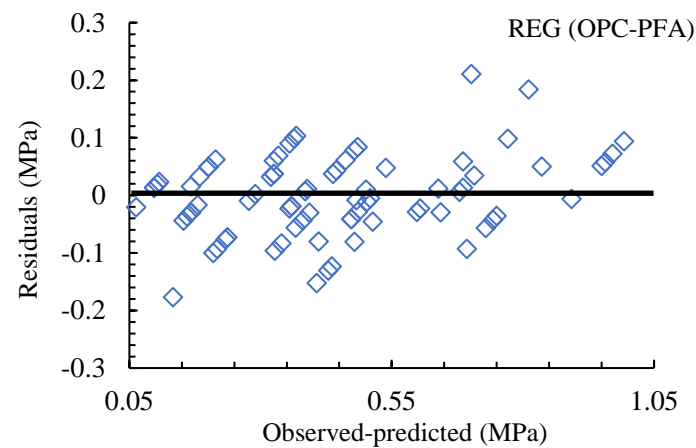
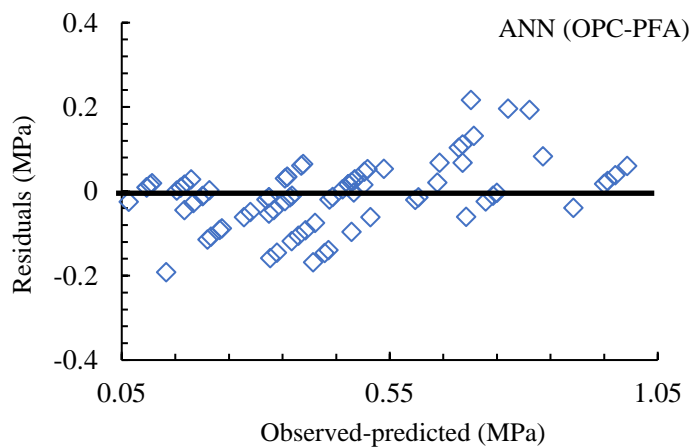
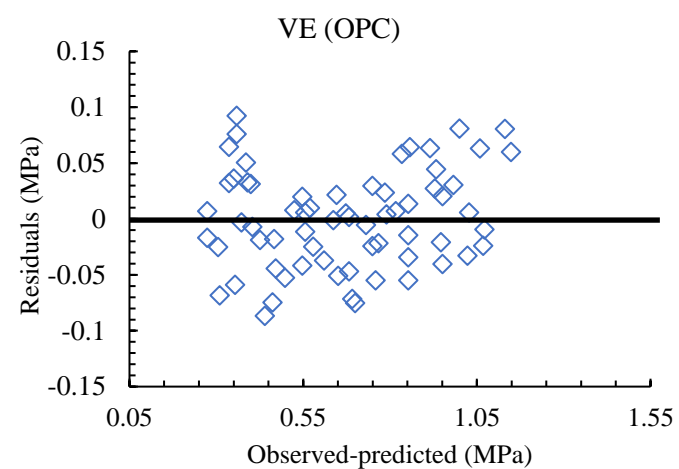
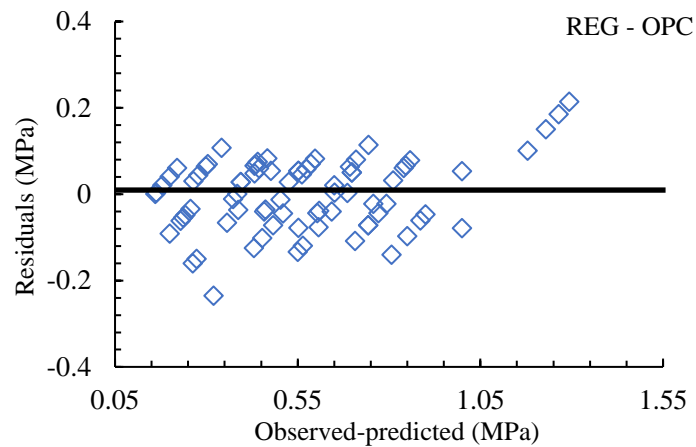
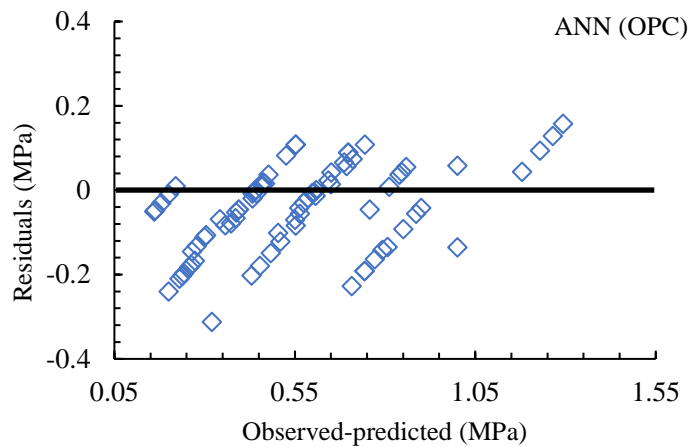
397 While performing feature engineering as discussed previously, it is important to note that the
 398 data used in the regression modelling were first normalised before being validated (train-
 399 validation split method) in order to prevent overfitting and imbalance. From Fig. 6, it is
 400 observed that the models used to perform the regression analysis all seem to indicate some
 401 measure of distribution about zero. However, a closer examination shows that the tree-
 402 ensemble and meta-ensembles all appear to exhibit more scatter and rather better randomness
 403 in the positioning of the data compared to the stand-alone models. Also, much better
 404 independence of error terms is exhibited by the tree-based and meta-ensembles. However, in
 405 terms of the features or the dependent variables used for the improvement of the soils, it is
 406 observed across the models that there is little or no difference in the degree of randomness of
 407 the 3 different combinations of the binders used (Fig. 6).

408

409

410





1 **Fig. 6.** Residual lag plots on trained dataset of ML models

2 *4.1.2. Performance forecast*

3 Indicators of the ML models' predictive performance are presented in **Table 3** however,
4 during the discussions herein, more attention will be given to the RMSE and R^2 metrics all of
5 which are highlighted (bold-face fonts) in **Table 3**. It is observed that all ML algorithms
6 provided predictions with some measured degrees of accuracy. The coefficient of
7 determination ranges from approximately 0.78 to 0.96, the RMSE varies from as low as 0.02
8 kPa to a high of 0.50 kPa. Interestingly, it could be observed that the REG model gives
9 predictions of the mixed soil's UCS with higher accuracy compared to the ANN, but with
10 broadly lower accuracy compared to the BLR, RDF and BDT. The later outcome may not be
11 entirely surprising with the REG model given that the data may not have fitted this
12 algorithm's underlying assumptions as they would with the BLR and the tree-ensemble
13 models (RDF and BDT). This behaviour was previously highlighted from the non-normality
14 of the residual lag plots showing that this model may be incapable of approximating some
15 unobserved phenomena of the mixed soil materials. It is pertinent to state here that the R^2
16 values obtained using the REG model are quite comparable to those of a previous study
17 which relied on similar dataset for its prediction [14]. Albeit the R^2 values of Table 3 are only
18 slightly lower because unlike the methodology of prediction adopted in the said previous
19 study where all the datasets were used in the training of the model, 80% and 20% of the
20 overall dataset were set aside and used in the training and testing of the REG model
21 respectively in this research following the train-validation split method. Of particular interest,
22 is the relatively worst prediction performance exhibited by the ANN model even though
23 several previous studies have indicated that this algorithm can predict the strength of
24 stabilised soils with reasonably high accuracy [24,27]. However, the ANN's inferior
25 performance could be explained by also leveraging some of the theories advanced in prior
26 studies [45,50], some of which are relevant to this research.

27 One major drawback of ANN is that the process of training is performed by relying on a
28 search and optimisation algorithm (Levenberg-Marquardt or gradient descent) to constantly
29 update its weights and biases over an error space that includes or that converges to local
30 minima instead of a more global one, an approach regarded as backpropagation (BP) [32].
31 Low performances in strength prediction that is sometimes demonstrated by ANN could be
32 mitigated by using Extreme ML (EML) algorithms whose training process also involves
33 single hidden-layer feed-forward mechanism [50]. In other words, ELMs provide even more
34 simplicity given that stopping criteria and learning rates (as given in **Table 1** for ANN) may
35 not have to be taken into consideration. But ELMs must also be used with caution because
36 such models could require many more hidden neurons due to the random need to determine
37 input weights and biases [51,52]. Moreover, excessive number of hidden layers of neurons
38 (or black-boxes) could lead to overfitting meaning that the complex nature of a stabilised
39 soil's mechanical property can be overestimated by feed-forward mechanisms like the ANN
40 [53].

41 In order to rectify or overcome some of the deficiencies mentioned above, the tree-based
42 ensemble ML algorithms could be used. As depicted in **Table 3**, these models seem to
43 produce relatively much better predictions. The RDF has R^2 of 0.89 and RMSE values of
44 0.34 and 0.35 for soil stabilised by OPC alone and OPC substituted by PFA and GGBS
45 respectively. Hence a slightly low prediction accuracy is given for the stabilised soil using a
46 combination of OPC and PFA. A further explanation and appreciation of the difference in
47 class predictions are given in subsequent sections in ML classification. On the other hand,
48 BDT seems to produce the most superior prediction performance of both the stand-alone and

49 tree-ensemble models with the highest R^2 of 0.94 and a corresponding RMSE of 0.19 for the
50 soil stabilised by using OPC alone. Unlike the RDF, the “boosting” strategy aided in an
51 improved performance of weaker regression tree learner algorithms. Also, it may be
52 suggested that the RDF performed slightly below the capacity of the BDT because for RDF,
53 all trees are of equal importance, hence it has the potential of being subjected to overfitting
54 given its minimum biases and wider variance.

55 It may not be completely unexpected that the tree-ensemble methods (BDT and RDF) have
56 provided higher degrees of accuracy and performance compared to the stand-alone models in
57 terms of the statistical metrics used in their assessment. The superior accuracy obtained from
58 the tree-based models are mainly attributed to their structure and architecture. The tree-based
59 ensemble methods are simply an aggregation or composition of single regression trees. This
60 combination or boosting of trees are needed because low predictions and overfitting could be
61 down to instability of a single regression tree when used alone [41]. Like the ANN models,
62 the formal rules needed for training and subsequent predictions by tree-ensembles are learnt
63 from patterns in the data. However, for the tree-ensembles, a series of tests or training are
64 required to be performed on the data in order to logically partition them and by so doing,
65 inconsistent variable features are learnt. To put it differently, the predictor variables are so
66 repeatedly partitioned such that each successive final partition generates different sets of
67 output value. Moreover, without having to smoothen or prune so-called “deep” trees,
68 generalisation errors are thus also reduced, and overfitting mitigated [32].

69 It may be possible that lower performance is achieved in tree-ensemble models. In this case,
70 the ensemble may have learnt from some interference due to the noise that ensues from the
71 residuals rather than signals that emanate from within the data [41]. This phenomenon is
72 attributed to a somewhat “greedy” construction process whereby at each step, an aggregation
73 of single best performing variable and optimal point of split is selected which invariably
74 means that there may also be a multi-step lookahead which takes into account variable
75 combinations with even better results. Another drawback which Dreiseitl and Ohno-Machado
76 [54] appear to agree with may have been the loss of information in the process resulting from
77 continuous discretisation of variables by the splitting process,

78 As demonstrated in this research, such weakness could be slightly reduced by creating a
79 further ensemble of models – meta-ensemble models (or model of models) through VE and
80 SE techniques. This method was used to aggregate other models (stand-alone and tree-
81 ensembles) with classifiers to enable an optimisation of the overall machine learning
82 predictive performance. The statistical performance metrics from the meta-ensembles were
83 obtained by calculating the weighted average of predictions from the combined models.

84 From Table 3, it is observed that the RMSE of the meta-ensemble models (VE and SE)
85 fluctuates between approximately 0.04 and 0.09. Hence, compared to the stand-alone and
86 tree-based models, RMSE values for the meta-ensemble models are about 4-5 folds lower.
87 The meta-ensemble models (most especially the VE) also seem to have very high accuracy of
88 prediction as observed from the R^2 values with the lowest being 0.80 and the highest, 0.96.
89 Regarding the coefficient of determination, the meta-ensemble models (though slightly
90 higher), seem slightly comparable with the algorithms derived from the tree-ensembles.
91 Comparing both meta-ensemble models, Table 3 indicates that the VE has a slightly better
92 performance than the SE given the later has higher RMSE and lower R^2 values in general.
93 With better decrease in the component variance of the prediction errors, an aggregation of the
94 ML algorithms is thus able to improve performance through the voting mechanism. Table 3
95 also reveals that predictions given for the stabilised soils with OPC substituted by equal
96 amounts of PFA and GGBS combinations are the worst.

97 The improvement observed in the performance of the meta-ensemble ML models stems from
 98 these models' capacity to incorporate predictions from the stand-alone and tree-ensemble
 99 algorithms meta-heuristically so that the outcome is even more accurate.

100

101 **Table 3.**

102 Statistical metrics indicating the performance of ML models utilised in strength prediction.

103

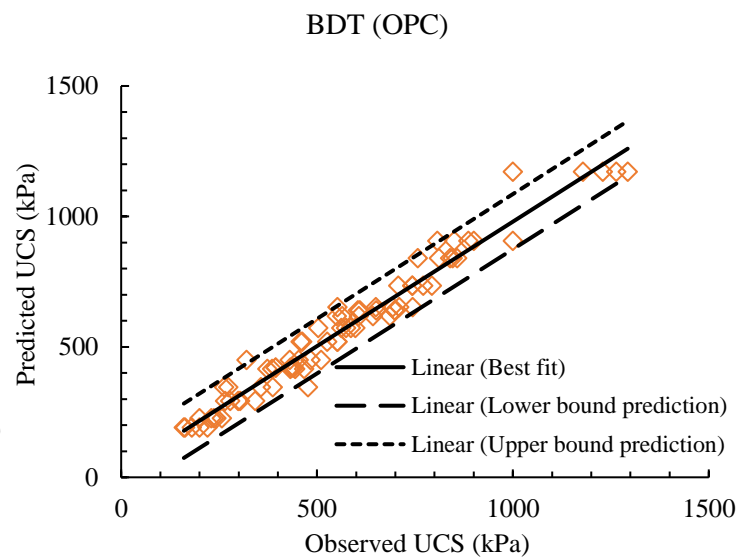
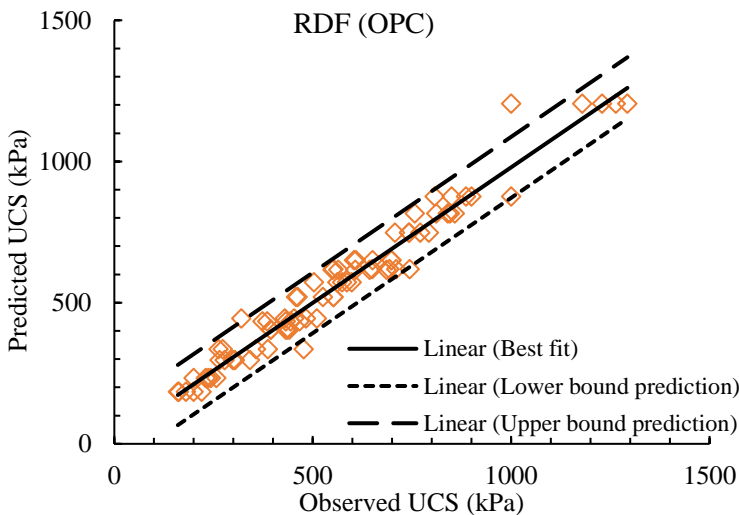
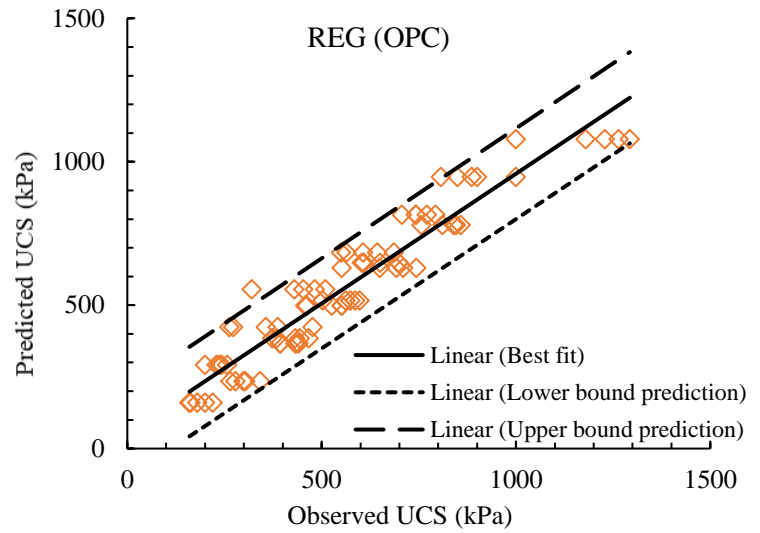
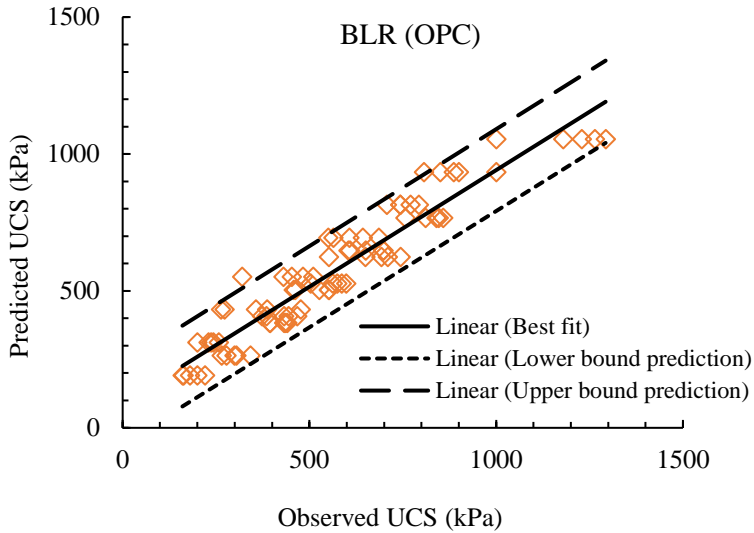
Model	Binder mix	Statistical metrics		
		R ²	RMSE	MAE
		-	kPa	kPa
BDT	OPC	0.94	0.26	0.19
	OPC-PFA	0.90	0.34	0.23
	OPC-PFA-GGBS	0.92	0.31	0.22
RDF	OPC	0.89	0.34	0.26
	OPC-PFA	0.85	0.41	0.27
	OPC-PFA-GGBS	0.89	0.35	0.26
BLR	OPC	0.91	0.31	0.27
	OPC-PFA	0.83	0.44	0.35
	OPC-PFA-GGBS	0.86	0.40	0.34
REG	OPC	0.91	0.31	0.27
	OPC-PFA	0.83	0.44	0.36
	OPC-PFA-GGBS	0.85	0.42	0.35
ANN	OPC	0.90	0.32	0.27
	OPC-PFA	0.78	0.50	0.42
	OPC-PFA-GGBS	0.78	0.50	0.42
SE	OPC	0.89	0.07	0.06
	OPC-PFA	0.94	0.05	0.04
	OPC-PFA-GGBS	0.79	0.09	0.08
VE	OPC	0.93	0.05	0.04
	OPC-PFA	0.93	0.05	0.04
	OPC-PFA-GGBS	0.96	0.04	0.03

104

105 *4.1.3. Uncertainty checks*

106 An assessment of each of the best performing regression ML models is established herein by
 107 considering their prediction intervals with 95% confidence. Following a transformation of the
 108 datasets using Eq. 1 to enable uniformity and linearity of prediction, Eq. 12 was further
 109 applied to derive both upper bound and lower bound interval of prediction based on the 2-

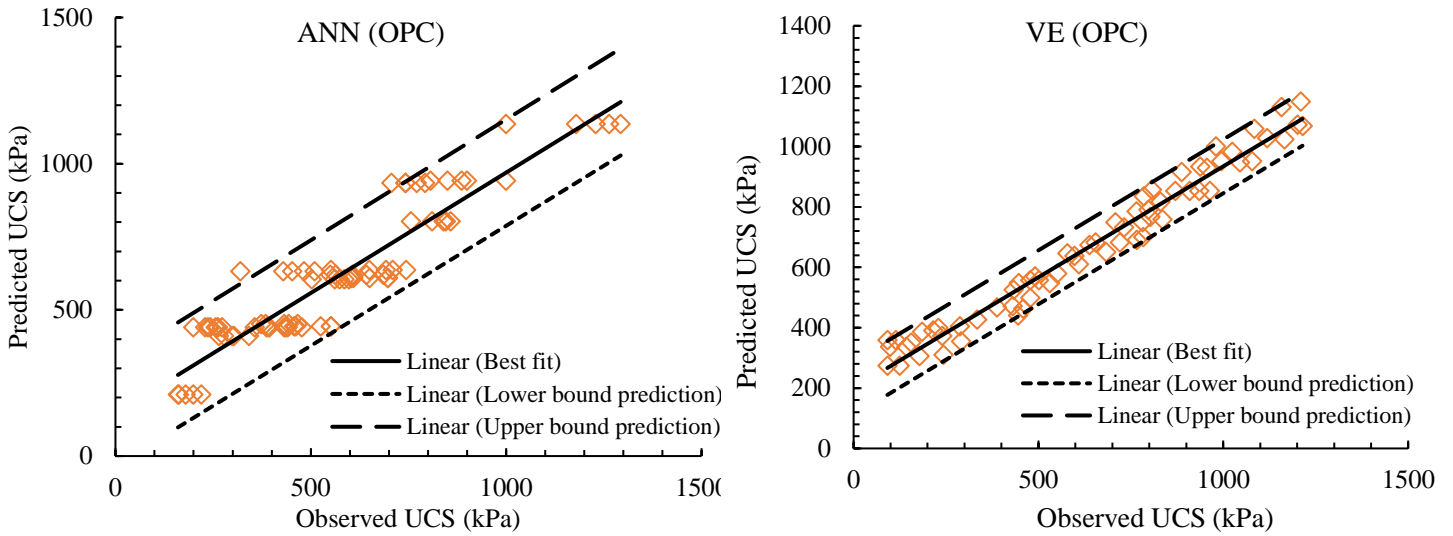
110 tailed t-test (at 2 degrees of freedom). Regardless of the differences in the range of data and
111 their distribution around the true predictor (the trendline), as indicated in Fig. 7, virtually all
112 the models (both stand-alone and ensembles) possess the interval that contains the dependent
113 variable with a confidence level of 95%. This observation mirrors those of the coefficient of
114 determination whereby the models exhibited relatively high values with the least being about
115 0.78. Moreover, the narrower ranges of prediction exhibited by the ensemble models (BDT,
116 RDF, VE and SE) proves further, their superior prediction capabilities.
117



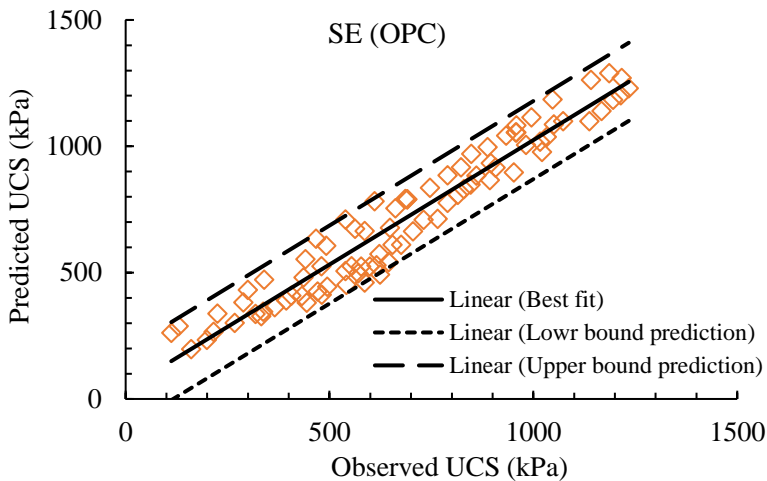
120

121

122



123

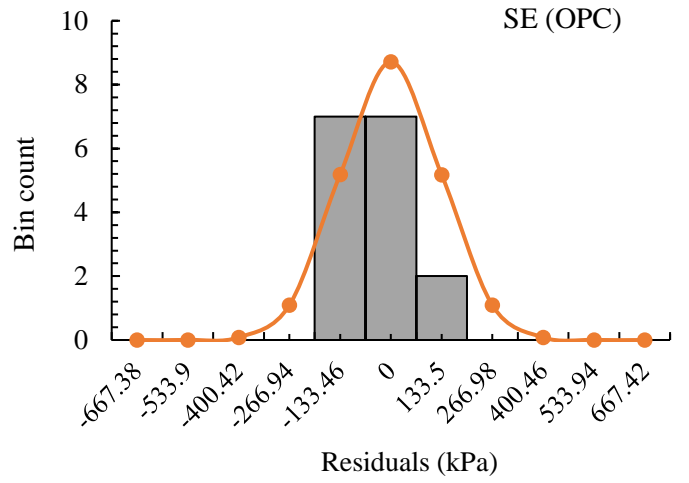
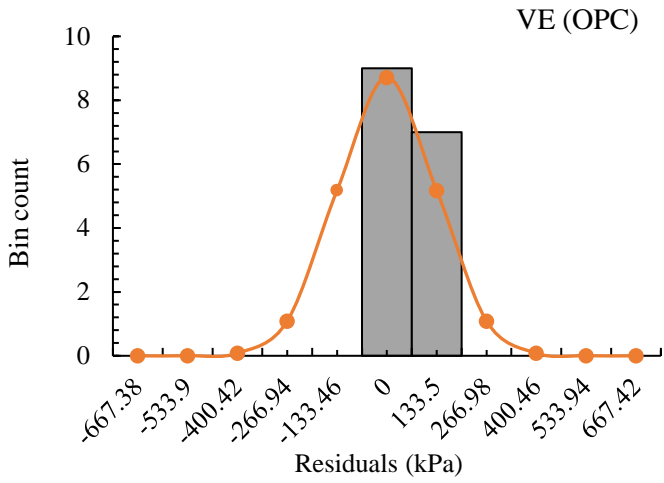


124

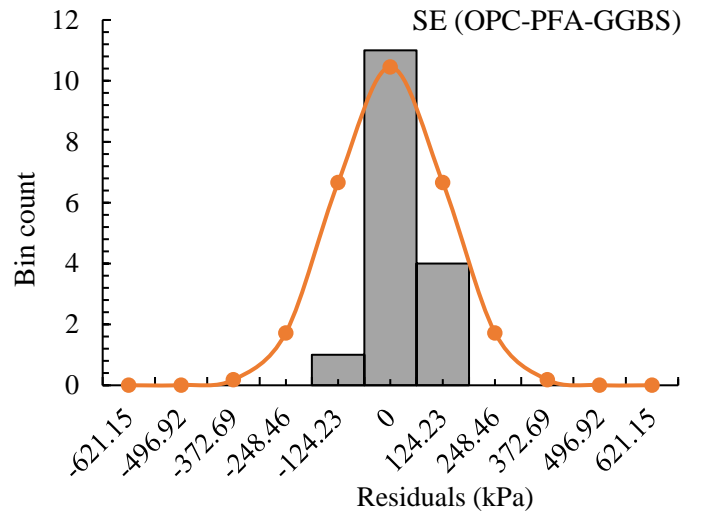
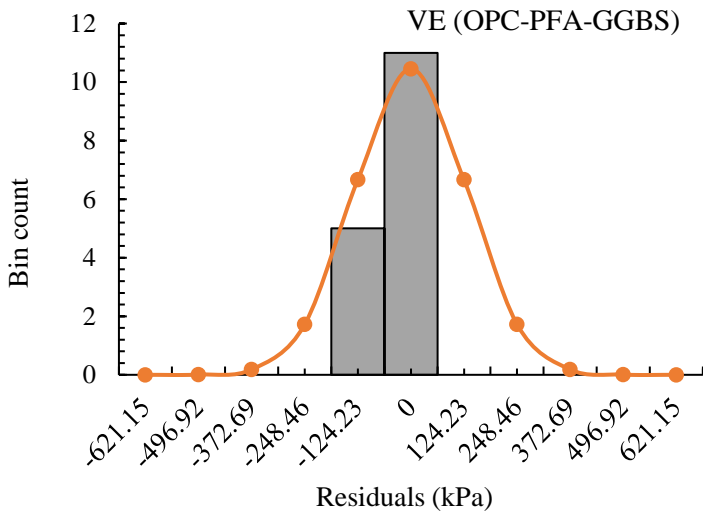
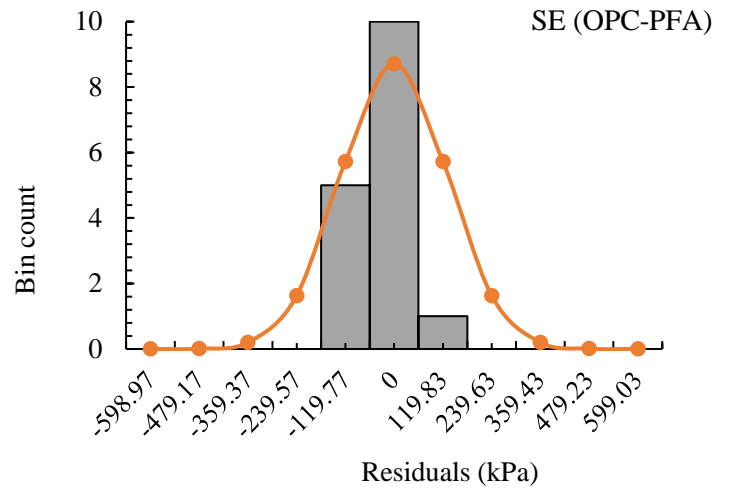
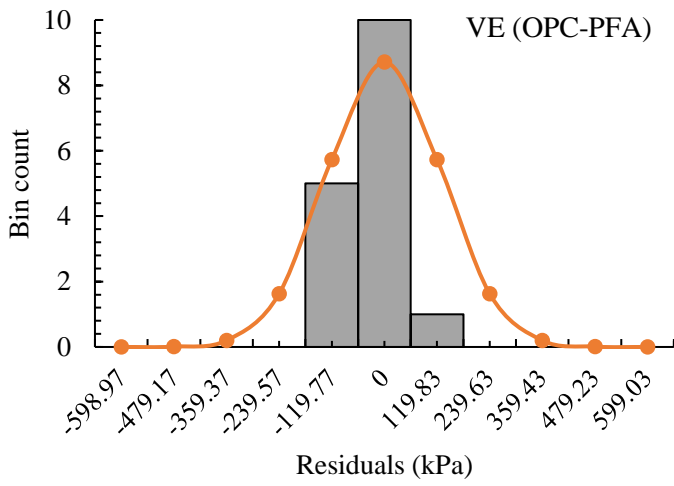
125 **Fig. 7.** Comparison of UCS strength predictions of the ML algorithms.

126 *4.1.4. Normality of meta-ensembles*

127 A diagnosis of normality of variance for the best performing meta-ensemble models were
 128 considered on the stabilised soils. The normality of assumptions for random error should hold
 129 true if on a histogram of the residuals, a symmetric bell-shaped curve or distribution is
 130 obtained [48]. It is observed from Fig. 8 that the residuals for both VE and SE prediction
 131 considering all the binder combinations appear to peak at zero and with less adjacent
 132 residuals. This phenomenon should indicate better performance; however, it is interesting to
 133 note that the distribution experienced from predictions on the soils stabilised by using only
 134 OPC seem slightly skewed and with lower bin counts compared to the soils stabilised with
 135 the OPC substituted by PFA and GGBS. Reasons for this behaviour cannot be advanced here
 136 except that the predictive features may not have been sufficient, thus making any further
 137 explanations rather inconclusive.



138



139

140

Fig. 8. Residuals and normal distribution plots of performance of meta-ensemble models

141

142 4.2. ML classification

143 The regression analyses presented in the foregoing were carried out on the data with three
 144 class features of binder combinations serving as dependent variables. Given the nature of
 145 these class features (each predicting the same output, the unconfined compressive strength),
 146 the result could be an arbitrary value. It is then very necessary that a classifier boundary
 147 between classes be determined by a threshold value.

148 First and foremost, looking at the classification metrics (Table 4) of the multiclass models
 149 employed for prediction of the unconfined compressive strength of the stabilised soils, it is
 150 observed that on average, the meta-ensemble models seem to perform better than some of the
 151 remaining 4 multiclass models. The multiclass logistics regression (*mLR*) model has the
 152 worst performance with an average accuracy of about 0.61, average precision of 0.42 and
 153 recall average of 0.42 probably due to an assumption of linearity even though there are
 154 instances of multi-collinearity between the dependent and predictor labels. Also, inaccurate
 155 predictions may have been increasing given the inability of *mLR* to sufficiently learn the
 156 categorical features. Interestingly, among the tree-based and stand-alone models, the
 157 multiclass *mANN* seems to produce the highest accuracy (0.78), precision (0.67) and recall
 158 (0.67). It should be recalled that when used previously in the regression analysis, the neural
 159 network algorithm gave the least performance which in the case of classification; that is,
 160 within the context of this study, the previously stated setbacks may have been rectified thus
 161 making *mANN* more suited to the complexities of non-linearity compared to *mLR*. The tree-
 162 ensembles (*mRDF* and *mDJ*) seem to have performed well on the multiclass problem due to
 163 their relatively high degree of accuracy in the prediction. Overall, the VE model appears to
 164 outperform all the other models given it possesses the highest classification metrics presented
 165 in Table 4. Further discussions of the meta-ensembles are provided below.

166
 167 **Table 4.**
 168 Classification metrics of multiclass ML models

Multiclass ML models	Average accuracy	Average precision	Average recall
<i>mRDF</i>	0.72	0.58	0.58
<i>mDJ</i>	0.67	0.50	0.5
<i>mLR</i>	0.61	0.42	0.42
<i>mANN</i>	0.78	0.67	0.67
VE	0.88	0.79	0.83
SE	0.72	0.60	0.58

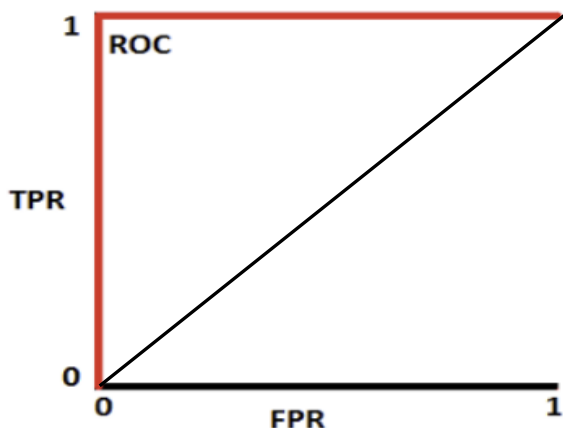
169

170 4.2.1. Sensitivity and multiclass prediction capability of the best algorithms

171 In order to assess the performance of the best meta-ensemble models, both the
 172 Receiver Operating Characteristic (ROC) and Lift curves were applied for some diagnostic
 173 tests of sensitivity and prediction capacity respectively. The best models are compared across
 174 some validation techniques applied to the normalised data before training by adopting 10-
 175 fold, Monte Carlo and Train-validation split cross-validation methods. This was done to
 176 provide an unbiased evaluation and estimate of the algorithm's calibration and
 177 discrimination during the classification process.

178 4.2.2. Receiver operating characteristic (ROC)

179 Within the context of analysis in this section, sensitivity or recall may be regarded as a
180 measure of how well the best prediction models can identify the true positives belonging to
181 either of the three independent class variables. Receiver operating characteristics (ROC)
182 curve is one of the most important probability evaluation metrics for assessing the best
183 performing classification model by indicating the relationship between true positive rate
184 (TPR) and false positive rate (FPR) during the course of any change in the decision threshold.
185 Along with AUC (Area under the curve) which is a measure of separability, the ROC
186 indicates how much an algorithm is able to distinguish between classes. The higher the AUC,
187 the better the performance of a given model. For the ROC curve, an excellent or perfect
188 classification is indicated by a point on the upper left corner with coordinates of 0 and 1 on
189 the TPR vs FPR graph (Fig. 9). That is also to say that this represents 100% sensitivity or
190 recall (no false negatives). A rather random act of guessing would produce a point along the
191 blue diagonal line (i.e., line of no-discrimination) that runs from the origin (0,0) to the top
192 right corners irrespective of the negative and positive base rates). This will indicate the worst
193 possible situation. That will mean the AUC is approximately 0.5 and that the model possesses
194 no discrimination capacity to distinguish between negative and positive classes.



195

196 **Fig. 9** Receiver operating curve (ROC)

197 **Table 5** shows the AUC for the meta-ensemble models across three hyper-parameter tuning
198 techniques. It is observed that the meta-ensemble models have an AUC value above 0.5. This
199 means that the meta-ensemble models at the least, do have a discrimination ability to
200 distinguish between negative and positive classes. The corresponding ROC curves (Fig. 10)
201 confirms this characteristic even though not exactly perfect or an ideal situation. Stated in
202 another way, it could be concluded that the nature of overlap has minimised any type 1 or
203 type 2 error. Across the VE ML model, adopting the train-validation split technique appears
204 to produce the best performance compared to SE. However, using the k -fold cross-validation
205 technique seems to give the worst performance for the VE ML model. On average, it could be
206 concluded from **Table 5** that the Monte Carlo cross-validation method provides a middle
207 ground for hyper-parameter tuning between the other two techniques.

208

209

210

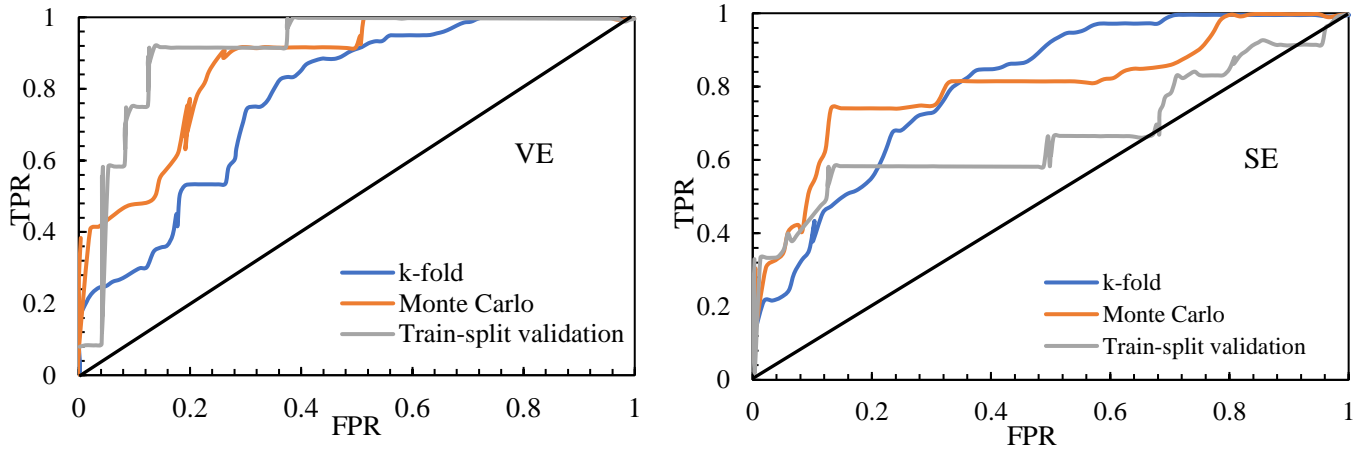
211

212

214 **Table 5**
 215 Comparison of Area under curve (AUC) values for different cross-validation methods

Meta-ensemble models	Cross-validation methods		
	<i>k</i> -fold	Monte Carlo	Train-validation split
VE	0.79	0.86	0.91
SE	0.81	0.83	0.67

216



217

218 **Fig. 10.** Receiver operating curve of meta-ensemble models

219 *4.2.3. Lift curve*

220 Since this study deals with more than a binary (two-class) classification problem, further
 221 proves of how a better model could perform when compared to a random model is provided
 222 by the lift curve. This relative performance stems from the theory that a random model is
 223 likely to make an incorrect prediction of a multiclass classification problem compared to a
 224 better model with higher fractions of the sampled data. Hence, given a random model, the lift
 225 curve is a visual representation of the ratio of cumulative gains to the cumulative gains for
 226 that random model. The corresponding baseline lift curve is the horizontal or percentile axis.
 227 The greater the area between the baseline and the lift curve, the better the model. **Fig. 11**
 228 represents the lift curve of the stabilised soils using the meta-ensemble models for predictive
 229 classifications following the application of *k*-fold, Monte Carlo and train-validation split
 230 cross-validations. Again, it is observed here that the overall area of the curves rising from the
 231 baseline indicates VE is the best performing model. Also, when comparing the performance
 232 of the hyper-parameter tuning techniques, it could be observed that train-validation split
 233 method provides the best validation. It is interesting to note how the *k*-fold method compares
 234 to its Monte Carlo counterpart given that both do possess almost the same trend and area
 235 under the lift curve rising from the percentile baseline.

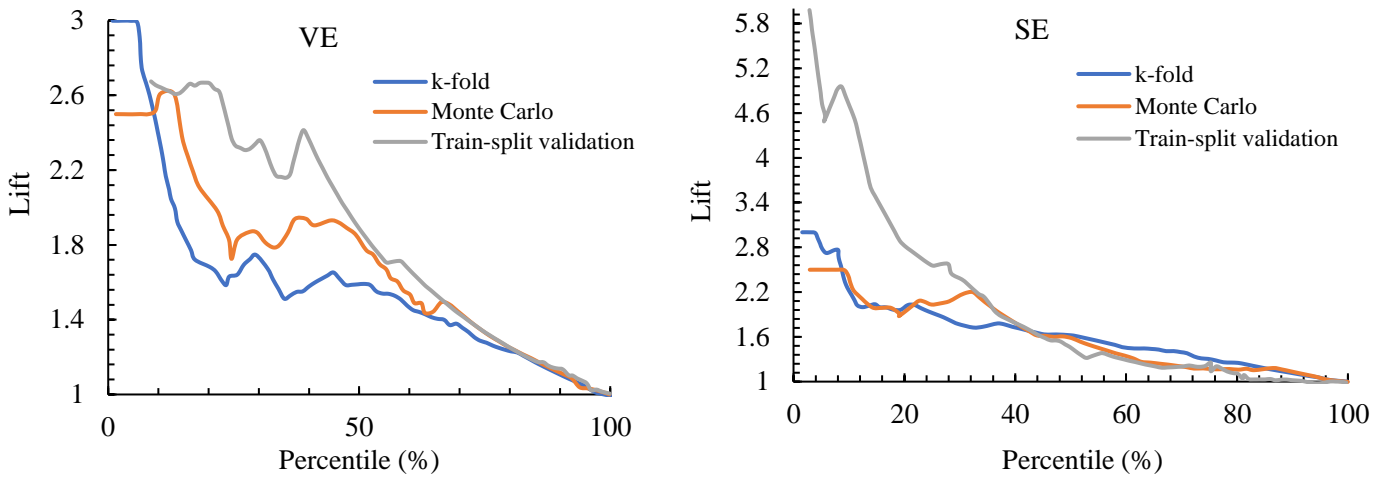
236

237

238

239

240



241

242 **Fig. 11.** A comparison of baseline lift performance of meta-ensemble models.

243

244 **5. Significance of study, recommendations, and deployment of ML models**

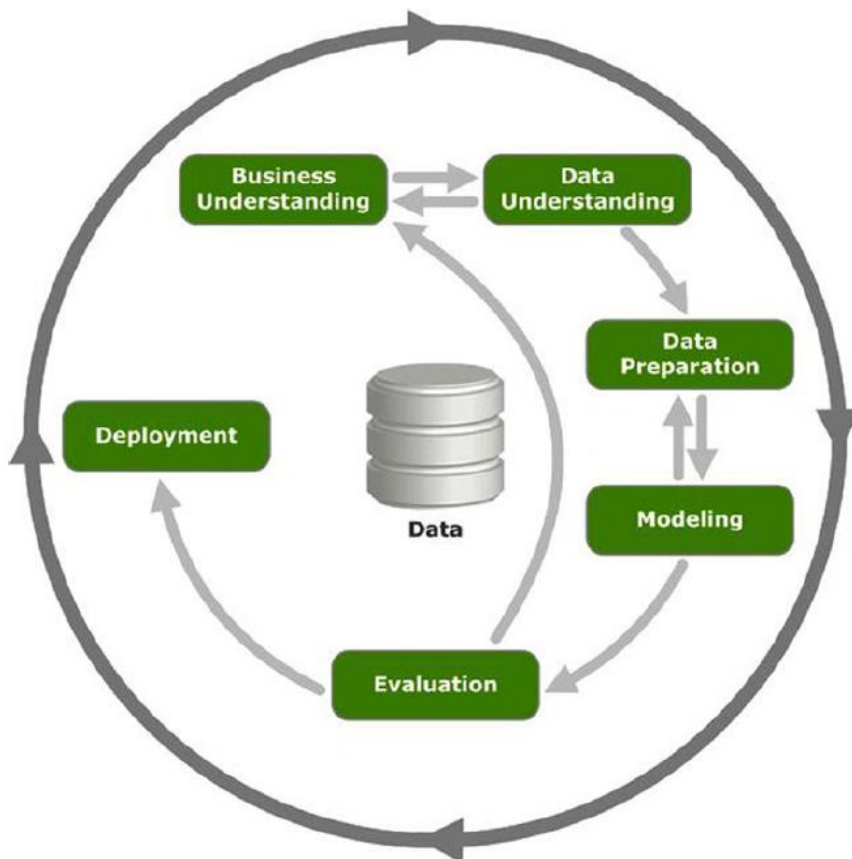
245 Computers with better processing speeds, higher computation power, and larger storage are
 246 some of the factors that characterise what is now termed the “age of information” or the “age
 247 of data”. Accordingly, researchers, data scientists, developers, and engineers have been
 248 working assiduously to study and develop tools, algorithms, techniques, frameworks, and
 249 methodologies to build intelligent systems and models that can predict events, perform
 250 complex analyses, automate tasks, detect anomalies, ensure autonomous or self-healing
 251 failures, and even understand as well as respond to human inputs. Hence, data-driven
 252 decision making by leveraging machine learning paradigms is quite beneficial in modern
 253 times for the following reasons [55] :

- 254 • Insufficient human knowledge and expertise in a domain (e.g., simulating navigations
 255 in unknown or uncharted territories or even spatial planets).
- 256 • The rapid flux in system behaviour over time (e.g., availability of organisational
 257 infrastructure, network connectivity, etc).
- 258 • The inability for humans to formally explain or translate a well-known domain
 259 expertise into computational tasks (e.g., speech recognition, transformation, cognitive
 260 tasks, scene recognition, etc).
- 261 • Addressing domain specific challenges at scale with large volumes of data
 262 characterised by lots of complex conditions and constraints.

263

264 The present study has successfully built upon while also enhancing the evolving concepts and
 265 ideas of artificial intelligence and ML most especially those reported in recent research
 266 within the realms of soil stabilisation and geotechnical engineering in general. Only but a few
 267 of the properties or features known to influence the unconfined compressive strength of
 268 stabilised soils have been considered herein. Hence, it is recommended that for even better
 269 and effective data-driven decision, an evaluation of various other factors (e.g., compaction
 270 condition, polymers, wastes, mineralogy, soil-water chemistry, soil structure, fabric, etc) and
 271 environmental constraints (temperature, groundwater movement, drainage, and other climatic
 272 conditions) which could potentially affect soil strength should be taken into account in the
 273 future based on the techniques and framework already proposed in this research.

274 Furthermore, it is pertinent to state that the methods and techniques of evaluation adopted in
275 this study represent a significant aspect of the end-to-end data mining lifecycle as suggested
276 by a typical CRISP-DM model which is depicted in Fig. 12. The CRISP-DM model is an
277 abbreviation for Cross Industry Standard Process for Data Mining. CRISP-DM indicates the
278 necessary processes, steps, and workflows for implementing any project right from
279 formalising business requirements up to and including testing and deployment of a solution to
280 transform data into valuable insights. This model does serve as a pointer to the tremendous
281 amount of interest and investments in the Data Science discipline across industries,
282 enterprises, companies, and domains. It also reinforces the earlier stated proposition that
283 intelligent ML systems and data-driven organisations are becoming a reality with the
284 advancements in tools and techniques only aiding in their further expansion. Hence, within
285 the context of this study, it is suggested that for a practical application of the concepts
286 developed, the ‘Deployment’ phase will ensure that the insights proposed are seamlessly
287 transferred to production in a real-life setting. Accordingly, the models and their predictions
288 as well as the background coding derived from this research can be deployed as saved files on
289 an organisation’s server, hardware or software resource and the proposed best meta-ensemble
290 models reloaded while predictions are offered for new data samples on both the studied
291 regression and classification problems. This can be applied either during preliminary stages
292 of a geotechnical site investigation or the design and construction phases to predict and assess
293 the strength performance of a stabilised soil.



294

295 **Fig. 12.** Data mining lifecycle of the CRISP-DM model [55].

296

297

6. Conclusions

In this study, an analysis of ML algorithms applied to regression and multiclass classification problems of soil improvement was carried out. The summary of strength prediction of soil stabilised by OPC and part-substitution of OPC with equal amounts and combinations of PFA and GGBS using stand-alone, tree-based and meta-ensemble ML algorithms are as follows:

- Using the stand-alone (REG, BLR, ANN) and tree-ensemble models (RDF and BDT), higher statistical variance are experienced by both the OPC-PFA and OPC-PFA-GGBS predictor variables compared to the soil stabilised by OPC alone. However, RDF appears to register the highest possible variance (about 0.82) followed by BDT (about 0.53) for the soil improved using only OPC.
- Quality assessment of the ML algorithms indicated that the tree-based and meta-ensembles (VE and SE) produced much better independence of error terms. However, in terms of the features or the dependent variables used, it was observed across the models that there was relatively little difference in the degree of randomness about the zero axis of the residuals plot of 3 different combinations of the binders.
- With regards an analysis of regression, on average, REG model produced predictions of the mixed soil's UCS with higher accuracy (RMSE of 0.39 and R^2 of 0.86) compared to the ANN (RMSE of 0.44 and R^2 of 0.82), but with comparatively lower accuracy compared to the tree-based models (average RMSE of 0.33 and R^2 of 0.90) and meta-ensemble models (average RMSE of 0.06 and R^2 of 0.91).
- For ML multiclass classification, multiclass neural network algorithm (*mANN*) gave the highest accuracy (0.78), precision (0.67) and recall (0.67) compared to tree-based and the remaining stand-alone models while only falling short to the meta-ensemble models (average accuracy of 0.80, precision of 0.70 and recall of 0.71).
- Sensitivity analysis from the receiver operating curve (ROC) and lift curves carried out across different validation techniques showed further prove of better performance of the meta-ensemble (VE) ML model compared to its SE ML counterpart when adopting the train-validation split technique as against the *k*-fold and Monte Carlo cross-validation methods.

References

- [1] European Commission, Reference Document on Best Available Techniques in Cement, Lime and Magnesium Oxide Manufacturing Industries, (2010) 459.
- [2] E.U. Eyo, S. Ng'ambi, S.J. Abbey, Incorporation of a nanotechnology-based additive in cementitious products for clay stabilisation, *Journal of Rock Mechanics and Geotechnical Engineering*. (2020). <https://doi.org/10.1016/j.jrmge.2019.12.018>.
- [3] E.U. Eyo, S.J. Abbey, S. Ngambi, E. Ganjian, E. Coakley, Incorporation of a nanotechnology-based product in cementitious binders for sustainable mitigation of sulphate-induced heaving of stabilised soils, *Engineering Science and Technology, an International Journal*. (2020). <https://doi.org/10.1016/j.jestch.2020.09.002>.
- [4] E.U. Eyo, S. Ng'ambi, S.J. Abbey, Performance of clay stabilized by cementitious materials and inclusion of zeolite/alkaline metals-based additive, *Transportation Geotechnics*. 23 (2020) 100330. <https://doi.org/10.1016/j.trgeo.2020.100330>.
- [5] E.U. Eyo, S. Ngambi, S.J. Abbey, Investigative study of behaviour of treated expansive soil using empirical correlations, in: *International Foundation Congress and Equipment Expo 5-10 March, Orlando, Florida, 2018*: pp. 373–384.
- [6] E.U. Eyo, S. Ngambi, S.J. Abbey, Investigative modelling of behaviour of expansive soils

- 348 improved using soil mixing technique., *International Journal of Applied Engineering Research*.
349 12 (2017) 3828–3836.
- 350 [7] S.J. Abbey, E.U. Eyo, J. Oti, S.Y. Amakye, S. Ngambi, *Mechanical Properties and*
351 *Microstructure of Fibre-Reinforced Clay Blended with By-Product Cementitious Materials*,
352 *Geosciences:MPDI*. 10 (2020).
- 353 [8] S.J. Abbey, A.O. Olubanwo, S. Ngambi, E.U. Eyo, B. Adeleke, *Effect of Organic Matter on*
354 *Swell and Undrained Shear Strength of Treated Soils*, *Journal of Civil, Construction and*
355 *Environmental Engineering*. 4 (2019) 48–58. <https://doi.org/10.11648/j.jccee.20190402.12>.
- 356 [9] S.J. Abbey, S. Ngambi, E. Coakley, *Effect of cement and by-product material inclusion on*
357 *plasticity of deep mixing improved soils*, *International Journal of Civil Engineering and*
358 *Technology*. 7 (2016) 265–274.
- 359 [10] S.J. Abbey, S. Ngambi, A.O. Olubanwo, *Effect of overlap distance and chord angle on*
360 *performance of overlapping soil-cement columns*, *International Journal of Civil Engineering*
361 *and Technology*. 8 (2017) 627–637.
- 362 [11] M. Kitazume, M. Terashi, *The deep mixing method*, 1st editio, CRC Press, 2013.
- 363 [12] EuroSoilStab, *Design Guide Soft Soil Stabilisation*, Ct97-0351. (1997) 95.
- 364 [13] A.E. Emeka, A.J. Chukwumeka, M.B. Okwudili, *Deformation behaviour of erodible soil*
365 *stabilized with cement and quarry dust*, *Emerging Science Journal*. 2 (2018) 383–387.
366 <https://doi.org/10.28991/esj-2018-01157>.
- 367 [14] S.J. Abbey, S. Ng'ambi, E. Ganjian, *Development of strength models for prediction of*
368 *unconfined compressive strength of cement/by-product material improved soils*, *Geotechnical*
369 *Testing Journal*. 40 (2017) 928–935. <https://doi.org/10.1520/GTJ20160138>.
- 370 [15] S. Motamedi, K. Il Song, R. Hashim, *Prediction of unconfined compressive strength of*
371 *pulverized fuel ash–cement–sand mixture*, *Materials and Structures/Materiaux et*
372 *Constructions*. 48 (2015) 1061–1073. <https://doi.org/10.1617/s11527-013-0215-1>.
- 373 [16] A. Gajurel, P.S. Mukherjee, B. Chittoori, *Estimating Optimal Additive Content for Soil*
374 *Stabilization Using Machine Learning Methods*, in: *Geo-Congress 2019*, 2019: pp. 662–672.
375 <https://doi.org/10.1061/9780784482124.067>.
- 376 [17] S. Soleimani, S. Rajaei, P. Jiao, A. Sabz, S. Soheilinia, *New prediction models for unconfined*
377 *compressive strength of geopolymer stabilized soil using multi-gen genetic programming*,
378 *Measurement: Journal of the International Measurement Confederation*. 113 (2018) 99–107.
379 <https://doi.org/10.1016/j.measurement.2017.08.043>.
- 380 [18] J. Tinoco, A. Gomes Correia, P. Cortez, *Support vector machines applied to uniaxial*
381 *compressive strength prediction of jet grouting columns*, *Computers and Geotechnics*. 55
382 (2014) 132–140. <https://doi.org/10.1016/j.compgeo.2013.08.010>.
- 383 [19] K.W. Liao, J.C. Fan, C.L. Huang, *An artificial neural network for groutability prediction of*
384 *permeation grouting with microfine cement grouts*, *Computers and Geotechnics*. 38 (2011)
385 978–986. <https://doi.org/10.1016/j.compgeo.2011.07.008>.
- 386 [20] E. Tekin, S.O. Akbas, *Artificial neural networks approach for estimating the groutability of*
387 *granular soils with cement-based grouts*, *Bulletin of Engineering Geology and the*
388 *Environment*. 70 (2011) 153–161. <https://doi.org/10.1007/s10064-010-0295-x>.
- 389 [21] E. Bachtiar, Mustaan, F. Jumawan, M. Artayani, Tahang, M.J. Rahman, A. Setiawan, M.
390 Ihsan, *Examining polyethylene terephthalate (Pet) as artificial coarse aggregates in concrete*,
391 *Civil Engineering Journal (Iran)*. 6 (2020) 2416–2424. [https://doi.org/10.28991/cej-2020-](https://doi.org/10.28991/cej-2020-03091626)
392 03091626.
- 393 [22] Z.N. Canbolat, G. Silahtaroglu, Ö. Doğuç, N. Yılmaztürk, *A machine learning approach to*
394 *predict creatine kinase test results*, *Emerging Science Journal*. 4 (2020) 283–296.
395 <https://doi.org/10.28991/esj-2020-01231>.
- 396 [23] R.A. Mozumder, A.I. Laskar, M. Hussain, *Empirical approach for strength prediction of*
397 *geopolymer stabilized clayey soil using support vector machines*, *Construction and Building*
398 *Materials*. 132 (2017) 412–424. <https://doi.org/10.1016/j.conbuildmat.2016.12.012>.
- 399 [24] J.-S. Chou, K.-H. Yang, J.-Y. Lin, *Peak Shear Strength of Discrete Fiber-Reinforced Soils*
400 *Computed by Machine Learning and Metaensemble Methods*, *Journal of Computing in Civil*
401 *Engineering*. 30 (2016) 04016036. [https://doi.org/10.1061/\(asce\)cp.1943-5487.0000595](https://doi.org/10.1061/(asce)cp.1943-5487.0000595).
- 402 [25] S. Suman, M. Mahamaya, S.K. Das, *Prediction of Maximum Dry Density and Unconfined*

- 403 Compressive Strength of Cement Stabilised Soil Using Artificial Intelligence Techniques,
404 International Journal of Geosynthetics and Ground Engineering. 2 (2016) 1–11.
405 <https://doi.org/10.1007/s40891-016-0051-9>.
- 406 [26] S. Hanandeh, A. Ardah, M. Abu-Farsakh, Using artificial neural network and genetics
407 algorithm to estimate the resilient modulus for stabilized subgrade and propose new empirical
408 formula, *Transportation Geotechnics*. 24 (2020) 100358.
409 <https://doi.org/10.1016/j.trgeo.2020.100358>.
- 410 [27] R.A. Mozumder, A.I. Laskar, Prediction of unconfined compressive strength of geopolymer
411 stabilized clayey soil using Artificial Neural Network, *Computers and Geotechnics*. 69 (2015)
412 291–300. <https://doi.org/10.1016/j.compgeo.2015.05.021>.
- 413 [28] A.H. Alavi, A.H. Gandomi, A. Mollahassani, A.A. Heshmati, A. Rashed, Modeling of
414 maximum dry density and optimum moisture content of stabilized soil using artificial neural
415 networks, *Journal of Plant Nutrition and Soil Science*. 173 (2010) 368–379.
416 <https://doi.org/10.1002/jpln.200800233>.
- 417 [29] S.K. Das, P. Samui, A.K. Sabat, Application of Artificial Intelligence to Maximum Dry
418 Density and Unconfined Compressive Strength of Cement Stabilized Soil, *Geotechnical and
419 Geological Engineering*. 29 (2011) 329–342. <https://doi.org/10.1007/s10706-010-9379-4>.
- 420 [30] J. Tinoco, C. António Alberto Santos, P. Da Venda, A.G. Correia, L. Lemos, A Data-driven
421 Approach for qu Prediction of Laboratory Soil-cement Mixtures, *Procedia Engineering*. 143
422 (2016) 566–573. <https://doi.org/10.1016/j.proeng.2016.06.073>.
- 423 [31] U.K. Sevim, H.H. Bilgic, O.F. Cansiz, M. Ozturk, C.D. Atis, Compressive strength prediction
424 models for cementitious composites with fly ash using machine learning techniques,
425 *Construction and Building Materials*. (2020) 121584.
426 <https://doi.org/10.1016/j.conbuildmat.2020.121584>.
- 427 [32] T. Han, A. Siddique, K. Khayat, J. Huang, A. Kumar, An ensemble machine learning approach
428 for prediction and optimization of modulus of elasticity of recycled aggregate concrete,
429 *Construction and Building Materials*. 244 (2020) 118271.
430 <https://doi.org/10.1016/j.conbuildmat.2020.118271>.
- 431 [33] M. Jalal, Z. Grasley, C. Gurganus, J.W. Bullard, Experimental investigation and comparative
432 machine-learning prediction of strength behavior of optimized recycled rubber concrete,
433 *Construction and Building Materials*. 256 (2020) 119478.
434 <https://doi.org/10.1016/j.conbuildmat.2020.119478>.
- 435 [34] H.L. Wang, Z.Y. Yin, P. Zhang, Y.F. Jin, Straightforward prediction for air-entry value of
436 compacted soils using machine learning algorithms, *Engineering Geology*. 279 (2020) 105911.
437 <https://doi.org/10.1016/j.enggeo.2020.105911>.
- 438 [35] H. Nguyen, T. Vu, T.P. Vo, H.T. Thai, Efficient machine learning models for prediction of
439 concrete strengths, *Construction and Building Materials*. 266 (2021) 120950.
440 <https://doi.org/10.1016/j.conbuildmat.2020.120950>.
- 441 [36] A. Marani, M.L. Nehdi, Machine learning prediction of compressive strength for phase change
442 materials integrated cementitious composites, *Construction and Building Materials*. 265 (2020)
443 120286. <https://doi.org/10.1016/j.conbuildmat.2020.120286>.
- 444 [37] M.C. Kang, D.Y. Yoo, R. Gupta, Machine learning-based prediction for compressive and
445 flexural strengths of steel fiber-reinforced concrete, *Construction and Building Materials*. 266
446 (2021) 121117. <https://doi.org/10.1016/j.conbuildmat.2020.121117>.
- 447 [38] A. Milad, S.A. Majeed, N.I.M. Yusoff, Comparative study of utilising neural network and
448 response surface methodology for flexible pavement maintenance treatments, *Civil
449 Engineering Journal (Iran)*. 6 (2020) 1895–1905. <https://doi.org/10.28991/cej-2020-03091590>.
- 450 [39] A. Sujatha, L. Govindaraju, N. Shivakumar, V. Devaraj, Fuzzy Knowledge Based System for
451 Suitability of Soils in Airfield Applications, *Civil Engineering Journal*. 7 (2021) 140–152.
- 452 [40] A. Joshi, *Machine Learning and Artificial Intelligence*, 1st ed., Springer International
453 Publishing, 2020. <https://doi.org/10.1007/978-3-030-26622-6>.
- 454 [41] M.A. DeRousseau, E. Laftchiev, J.R. Kasprzyk, B. Rajagopalan, W. V. Srubar, A comparison
455 of machine learning methods for predicting the compressive strength of field-placed concrete,
456 *Construction and Building Materials*. 228 (2019) 116661.
457 <https://doi.org/10.1016/j.conbuildmat.2019.08.042>.

- 458 [42] D. Kong, J. Zhu, C. Duan, L. Lu, D. Chen, Bayesian linear regression for surface roughness
459 prediction, *Mechanical Systems and Signal Processing*. 142 (2020) 106770.
460 <https://doi.org/10.1016/j.ymssp.2020.106770>.
- 461 [43] M.A. DeRousseau, J.R. Kasprzyk, W. V. Srubar, Computational design optimization of
462 concrete mixtures: A review, *Cement and Concrete Research*. 109 (2018) 42–53.
463 <https://doi.org/10.1016/j.cemconres.2018.04.007>.
- 464 [44] J. Shotton, S. Nowozin, T. Sharp, J. Winn, P. Kohli, A. Criminisi, Decision jungles: Compact
465 and rich models for classification, *Advances in Neural Information Processing Systems*.
466 (2013) 1–9.
- 467 [45] S.K. Das, *Artificial Neural Networks in Geotechnical Engineering: Modeling and Application*
468 *Issues*, First Edit, Elsevier Inc., 2013. <https://doi.org/10.1016/B978-0-12-398296-4.00010-6>.
- 469 [46] D.H. Wolpert, Stacked generalization, *Neural Networks*. 5 (1992) 241–259.
470 [https://doi.org/https://doi.org/10.1016/S0893-6080\(05\)80023-1](https://doi.org/https://doi.org/10.1016/S0893-6080(05)80023-1).
- 471 [47] K. Mamudur, M.R. Kattamuri, Application of Boosting-Based Ensemble Learning Method for
472 the Prediction of Compression Index, *Journal of The Institution of Engineers (India): Series A*.
473 101 (2020) 409–419. <https://doi.org/10.1007/s40030-020-00443-7>.
- 474 [48] N. Galwey, *Introduction to mixed modelling : beyond regression and analysis of variance*, 2nd
475 ed., Chichester, [England] : Wiley, 2014.
- 476 [49] S. Chatterjee, A. Hadi, *Regression analysis by example*, 5th ed., Hoboken, New Jersey :
477 Wiley, 2012.
- 478 [50] W. Ben Chaabene, M. Flah, M.L. Nehdi, Machine learning prediction of mechanical properties
479 of concrete: Critical review, *Construction and Building Materials*. 260 (2020) 119889.
480 <https://doi.org/10.1016/j.conbuildmat.2020.119889>.
- 481 [51] Q.Y. Zhu, A.K. Qin, P.N. Suganthan, G. Bin Huang, Evolutionary extreme learning machine,
482 *Pattern Recognition*. 38 (2005) 1759–1763. <https://doi.org/10.1016/j.patcog.2005.03.028>.
- 483 [52] L. Zhang, D. Zhang, Evolutionary Cost-Sensitive Extreme Learning Machine, *IEEE*
484 *Transactions on Neural Networks and Learning Systems*. 28 (2017) 3045–3060.
- 485 [53] A. Behnood, E.M. Golafshani, Predicting the compressive strength of silica fume concrete
486 using hybrid artificial neural network with multi-objective grey wolves, *Journal of Cleaner*
487 *Production*. 202 (2018) 54–64. <https://doi.org/10.1016/j.jclepro.2018.08.065>.
- 488 [54] S. Dreiseitl, L. Ohno-Machado, Logistic regression and artificial neural network classification
489 models: A methodology review, *Journal of Biomedical Informatics*. 35 (2002) 352–359.
490 [https://doi.org/10.1016/S1532-0464\(03\)00034-0](https://doi.org/10.1016/S1532-0464(03)00034-0).
- 491 [55] S. Dipanjan, B. Raghav, S. Tushar, *Practical Machine Learning with Python: A Problem-*
492 *Solver’s Guide to Building Real-World Intelligent Systems*, Apress, Berkeley, CA, 2018.
493 <https://doi.org/10.1007/978-1-4842-6222-1>.
- 494