

1 **Straightforward prediction for air-entry value of compacted soils using** 2 **machine learning algorithms**

3 Han-Lin Wang¹, Zhen-Yu Yin^{1,*}, Pin Zhang¹, and Yin-Fu Jin¹

4

5 ¹ Department of Civil and Environmental Engineering, The Hong Kong Polytechnic
6 University, Hung Hom, Kowloon, Hong Kong, China.

7 * **Corresponding author:**

8 Dr. Zhen-Yu Yin, Tel: (+852) 3400 8470, Fax: (+852) 2334 6389, Email:

9 zhenyu.yin@polyu.edu.hk; zhenyu.yin@gmail.com

10

11 **Abstract:** The straightforward prediction for the air-entry value of compacted soils is
12 practically useful, but the investigation on this issue is scarce. This study presents three
13 alternative straightforward prediction models for the air-entry value of compacted soils using
14 the representative machine learning algorithms of multi expression programming (MEP),
15 evolutionary polynomial regression (EPR) and random forest (RF). Five known soil
16 properties (i.e. sand content, fines content, plasticity index, initial water content and initial
17 void ratio) are used as input variables. All models are developed based on a large database,
18 covering a wide range of soil classifications. The results show that all the three proposed
19 models are appropriate to predict the air-entry values of different compacted soils, with high
20 prediction accuracies for both the training and the testing data. The monotonicity, the
21 sensitivity and the robustness of the three prediction models are evaluated, showing
22 consistency among different models with a slight difference and providing a strong support
23 for the model feasibility. On the whole, the MEP and the EPR models are recommended for
24 more convenient applications with explicit expression, while higher prediction accuracy may
25 require the RF model although no explicit expression can be derived.

26 **Keywords:** machine learning; air-entry value; multi expression programming; evolutionary
27 polynomial regression; random forest; compacted soils

28 **Introduction**

29 Compacted soils are widely used in geotechnical or geoenvironmental engineering, such as
30 the urban/landfill barrier system (Tinjum et al. 1997; Miller et al. 2002; Osinubi and Nwaiwu
31 2006; Birle et al. 2008; Krisdani et al. 2008), buffer for nuclear waste disposal (Delage et al.
32 1998; Cui et al. 2002; Sun et al. 2009; Ye et al. 2009; He et al. 2019), transportation
33 substructure (Zhao et al. 2016; Zhou et al. 2016; Wang et al. 2017, 2018, 2019a, 2019b; Wang
34 and Chen 2019; Chen et al. 2020b, 2020c, 2020d; de Freitas et al. 2020; Liu et al. 2020), etc.
35 During their servicing period, the soils get saturated during the wetting conditions and the
36 unsaturated state emerges when drying occurs. Through the wetting-drying cycles, the
37 hydro-mechanical behaviours of these compacted soils are highly related to the soil-water
38 characteristic curves (SWCC, Fredlund et al. 2012). Fig. 1 shows a conceptual SWCC for the
39 drying process, plotted by the degree of saturation S_r (the percentage between the volume
40 of liquid water and the volume of voids) against the matric suction ψ (the difference
41 between the pore air pressure and the pore water pressure). In terms of the SWCC, two
42 separation points can be easily identified: air-entry value (*AEV*) and residual degree of
43 saturation. The air-entry value serves as the matric suction which is required to cause
44 desaturation of the largest pores (i.e. beyond which suction value the air starts to enter the
45 pores of the saturated soil). As shown in Fig. 1, the air-entry value can be obtained by
46 extending the constant slope section of the SWCC to intersect the suction axis at $S_r = 100\%$.
47 Before the air-entry value, water fills the pores of the saturated soil as the matric suction
48 increases (I: boundary effect zone). When the suction increases beyond the air-entry value,
49 water in the pores starts to lose with the continuous increase of the matric suction (II:

50 transition zone). After the matric suction reaches the point at the residual state, liquid water in
51 the pores becomes discontinuous. In the residual zone (phase III), water is difficult to lose as
52 the matric suction continues to increase.

53 Due to the importance of the air-entry value in determining the water holding capacity
54 and further the hydro-mechanical behaviours of compacted soils, this parameter is considered
55 in all typical SWCC models (e.g. Brooks and Corey 1964; van Genuchten 1980; Fredlund
56 and Xing 1994). In the engineering practice, the straightforward prediction about the air-entry
57 value of compacted soils is also very important for the design and the maintenance of
58 geotechnical structures. For example, when a subgrade is planned to be constructed in a
59 specific area with known humidity data throughout the year and the subgrade contains a
60 barrier layer to hold the water under heavy precipitation, the prediction of the air-entry value
61 of this layer is needed. During the raining conditions, the barrier layer stores rainwater. Under
62 unfavourable conditions, the rainwater cannot go deeper or be lost if the humidity-induced
63 suction is lower than the air-entry value. Thus, the threshold suction value can be used as a
64 reference for the selection of the barrier materials and initial placement conditions for these
65 materials.

66 In general, the air-entry value is determined through the SWCC, which can be obtained
67 by laboratory tests or analytical models (Fredlund et al. 2012). The laboratory tests to
68 measure the SWCC include pressure plate test, vacuum desiccator, chilled-mirror dew-point
69 method, filter paper, unsaturated oedometer or unsaturated triaxial apparatus, etc. All these
70 tests are costly and time-consuming (lasting for several months for a specific sample). Thus,
71 several models were proposed to predict the SWCC of different soils from the soil properties

72 (Fredlund et al. 2002; Johari et al. 2006; Li et al. 2014; Zhou et al. 2014; Zhai et al. 2020).
73 However, using these methods, the air-entry value still needs to be determined after fitting the
74 SWCC. Thus, the accuracy of determining the air-entry value is highly dependent on the
75 prediction of the SWCC. Furthermore, in these studies, the database used for predicting the
76 SWCC was not large enough to cover a wide range of soils. To the authors' knowledge, the
77 straightforward prediction for the air-entry value of the compacted soils remains scarce.

78 To solve nonlinear and complex problems for the prediction with a large database,
79 several machine learning algorithms have been proven as effective approaches. In the field of
80 geotechnical engineering, the machine learning algorithms have been successfully used to
81 predict cyclic soil response (Shahnazari et al. 2010), creep index (Jin et al. 2019), bearing
82 capacity of composite column (Sarir et al. 2019), spatiotemporal response of rooted soil
83 (Cheng et al. 2020a), suction distribution close to tree (Cheng et al. 2020b), soil liquefaction
84 (Njock et al. 2020), jet grouted diameter in soft soils (Shen et al. 2020), tunneling induced
85 settlement (Zhang et al. 2020b), etc. In these studies, the algorithms cover the representative
86 multi expression programming (MEP), evolutionary polynomial regression (EPR) and
87 random forest (RF), etc. To date, the air-entry value of compacted soils has scarcely been
88 predicted by the machine learning algorithms. To obtain the air-entry value of compacted
89 soils in a fast and accurate manner, a comprehensive understanding of different machine
90 learning algorithms on the prediction of the air-entry value is imperative and worth
91 investigating.

92 In this study, alternative straightforward prediction models for the air-entry value of
93 compacted soils are developed using three commonly used machine learning algorithms:

94 MEP, EPR and RF. A large database of soils with multi classifications is collected from the
95 previous publications. Two-thirds of the data are chosen as the training data, while the
96 remaining are used for testing. Each prediction model is developed with each algorithm using
97 the respective optimum parameters. The prediction accuracy of the three models are verified
98 by the training and the testing data. The feasibility is further examined by the monotonicity,
99 sensitivity and robustness analysis for all three models, along with discussing their
100 advantages and limitations.

101

102 **Machine learning algorithms**

103 *Multi expression programming*

104 Multi expression programming (MEP) is a representative approach to linear-based genetic
105 programming (GP). In a chromosome of the MEP algorithm, multiple solutions (programs)
106 can be encoded, starting with the creation of a random population of computer programs. The
107 first gene of a chromosome must be a terminal randomly selected from the terminal set. In the
108 following genes, a gene with a function has a pointer towards the function arguments. For a
109 specific gene, the expression indices have lower values than the position of this gene in the
110 chromosome. Through the calculation of the fitness of all expressions, the best encoding
111 solution is determined to represent the chromosome by repeating the following steps, until the
112 termination condition is reached (Oltean and Grosan 2003): (i) selecting two parents by a
113 procedure of binary tournament and recombining them with a fixed crossover probability; (ii)
114 obtaining two offspring by recombining two parents; (iii) mutating the offspring and
115 replacing the worst individual in the current population with the best of them when the

116 offspring is better than the worst one. After the identification of the best solution, the explicit
 117 expressions can be generated by reading the chromosome from top to bottom.

118 ***Evolutionary polynomial regression***

119 Evolutionary polynomial regression (EPR) is another type of genetic programming (GP), to
 120 develop symbolic models following two steps: (i) structure identification, and (ii) parameter
 121 estimation (Giustolisi and Savic 2006). In the first step, the genetic algorithm (GA) is adopted
 122 to search for symbolic structures of EPR. During the second step, the values of parameters
 123 are estimated by solving the least squares (LS) linear problem. The advantage of EPR
 124 highlights that a simple explicit expression can be presented in the EPR algorithm, to
 125 describe the correlation between input and output variables. A general EPR expression is
 126 formulated as:

127
$$t = \sum_{j=1}^m a_j \cdot z_j + a_0 \quad (1)$$

128 where t is the predicted output; a_j is an adjustable parameter for the j th term; a_0 is an
 129 optional bias; m is the number of transformed terms; z_j is the j th transformed variable,
 130 which can be obtained by:

131
$$z_j = x_1^{E_{j,1}} \cdot \dots \cdot x_k^{E_{j,k}} \quad (2)$$

132 where x_i is the i th input variable; k is a total number of input variables; $\mathbf{ES}_{m \times k}$ is the
 133 exponent matrix, determined by GA. The key objective of EPR is to identify the best form of
 134 the function: the number of transformed variables and a combination of vectors of
 135 independent input variables. Then, the adjustable parameters and an optional bias can be
 136 determined by the least squares regression. Finally, the optimum explicit expression can thus
 137 be deduced.

138 *Random forest*

139 Random forest (RF) is an ensemble learning algorithm, integrated with the methods of
140 bootstrap aggregating (Breiman 1996) and random subspace (Ho 1998). Due to the
141 integration of numerous decision trees, the prediction of RF shows a strong performance
142 (Zhang et al. 2019, 2020a, 2020b, 2020c). In bagging, n_t bootstrap sets are built by
143 sampling with the replacement of N training examples from the training database. The
144 number of samples in the bootstrap training set is arbitrary, less than the original one. Then,
145 each bootstrap set is used to develop a decision tree. Each node in a decision tree represents a
146 classification criterion, with the leaves of the tree representing the output labels. Hence, a
147 decision tree classifies a bootstrap training sample by testing random features at each node.
148 As a result, a regression space can be determined. The ultimate predicted output t can be
149 obtained by aggregating the outputs of all trees as (Liaw and Wiener 2002):

$$150 \quad t = \frac{1}{n_t} \sum_{i=1}^{n_t} t_i(\mathbf{x}) \quad (3)$$

151 in which $t_i(\mathbf{x})$ is the predicted output for a tree with an input vector \mathbf{x} ; n_t is the number
152 of trees.

153

154 **Model development**

155 *Database collection*

156 To directly predict the air-entry value of compacted soils, 189 relevant samples are collected
157 from the experimental data of previous publications (Han et al. 1995; Tinjum et al. 1997;
158 Huang et al. 1998; Vanapalli et al. 1999; Ng and Pang 2000; Agus et al. 2001; Montanez
159 2002; Khalili et al. 2004; Yang et al. 2004; Indrawan et al. 2006; Puppala et al. 2006; Sun et

160 al. 2006; Thu et al. 2006, 2007; Birle et al. 2008; Krisdani et al. 2008; Rahardjo et al. 2008;
161 Li 2009; Zhang and Chen 2009; Gallage and Uchimura 2010; Zhou and Kong 2011; Mirzaii
162 and Yasrobi 2012; Oh et al. 2012; Rahardjo et al. 2012; Lin and Cerato 2013; Salager et al.
163 2013; Sun et al. 2014; Sun and Gao 2015; Amadi and Osinubi 2016; Cuceoglu 2016; Han and
164 Vanapalli 2016; Hashem and Houston 2016; Priono et al. 2016; Fattah et al. 2017; Jiang et al.
165 2017; Satyanaga et al. 2017; Chen et al. 2019, 2020a; de Freitas et al. 2020). According to
166 these studies, the main influencing factors on the air-entry value of compacted soils include
167 grain size distribution, plasticity and initial placement conditions. Thus, the soil properties of
168 gravel content C_G , sand content C_S , fines content C_F , plasticity index PI , initial water
169 content w_0 , initial void ratio e_0 and air-entry value AEV are collected. The details of the
170 soil properties and the relevant testing methods can be downloaded and referred to the
171 supplementary database. Note that the gravel, sand and fines are separated by the grain size
172 range of 75 mm to 4.75 mm, 4.75 mm to 0.075 mm and smaller than 0.075 mm, respectively
173 (ASTM 2017). In this database, various soil classifications (ASTM 2017) are collected,
174 including lean clay (CL), silty clay (CL-ML), fat clay (CH), silt (ML), elastic silt (MH),
175 well-graded sand (SW), well-graded sand with clay (SW-SC), poorly graded sand (SP),
176 poorly graded sand with clay (SP-SC), clayey sand (SC), silty clayey sand (SC-SM), silty
177 sand (SM), well-graded gravel with silt (GW-GM), clayey gravel (GC) and poorly graded
178 gravel (GP). In the supplementary database, these soils are ordered and numbered by the soil
179 classification. For the same classification, the soils are ordered alphabetically by the name of
180 the authors.

181 Table 1 lists the descriptive statistics of each variable in the database, with the values of

182 minimum, maximum, mean and standard deviation. Fig. 2 shows the detailed frequency
183 histogram of each variable, including the soil classification. As the gravel soil presents a
184 relatively lower water holding capacity, the SWCC of this kind of soil was not widely
185 investigated in the literature. Thus, the majority of the collected soils have the gravel content
186 of less than 20% (Fig. 2a), leading to a mean gravel content of 3.8% (Table 1). Compared to
187 the gravel content, the frequency distribution of the sand and the fines content is more
188 uniform, with the highest value at 0% to 20% (Fig. 2b) and 80% to 100% (Fig. 2c),
189 respectively. Regarding the plasticity index, around 120 soils are in the range of 0 to 20,
190 whereas 19 samples show the plasticity index higher than 40 (Fig. 2d). The initial water
191 content ranges from 0.5% to 48.6% (Table 1), showing the highest frequency at 10% to 20%
192 (Fig. 2e). The majority of the initial void ratio concentrates in the range from 0.4 to 0.8 (Fig.
193 2f), which is also common for the compacted soils. In this database, a wide range of soil
194 classifications from clay to gravel is introduced (Fig. 2g), showing the highest frequency for
195 the samples of lean clay (48) and clayey sand (35). In terms of the air-entry value, most of the
196 values are located in the range from 0 kPa to 20 kPa (Fig. 2h), while the highest air-entry
197 value is 100 kPa (Table 1).

198 To have an overall understanding of the distribution of the air-entry value of compacted
199 soils with different classifications, Fig. 3 is plotted with the air-entry value versus the specific
200 and the general classification, respectively. Table 2 lists the minimum, the maximum, the
201 mean and the median air-entry values for the general soil classifications. It is widely known
202 that at the same initial placement conditions, the fine-grained soil should have a higher water
203 holding capacity than the coarse-grained soil, leading to a higher air-entry value of the

204 fine-grained soil. From Fig. 3 (a), it can be observed that the maximum air-entry value shows
 205 for the fat clay (CH, 100 kPa). However, the highest mean air-entry value (box symbol in the
 206 figure) locates for the well-graded sand with clay (SW-SC). This is because for this soil in the
 207 database, only 4 samples are collected and 3 samples have the fines with a very high
 208 plasticity index (88), resulting in high air-entry values. Regarding the general soil
 209 classification in Fig. 3 (b) and Table 2, the clay has the highest air-entry values of maximum,
 210 median (transverse line symbol) and minimum, and the silt has the highest mean air-entry
 211 value. In other words, the air-entry value of the fine-grained soil shows relatively higher
 212 values than the coarse-grained soil, supporting the feasibility of the supplementary database.

213 Before developing the prediction model, the basic linear fitting between the air-entry
 214 value and each input soil property is depicted to have an overall view of the monotonic
 215 variation trend, as shown in Fig. 4. The coefficient of determination R^2 is used to evaluate
 216 the fitting accuracy as:

$$217 \quad R^2 = 1 - \frac{\sum_{i=1}^n (h_i - t_i)^2}{\sum_{i=1}^n (h_i - \bar{h}_i)^2} \quad (4)$$

218 where h_i and t_i are the actual and predicted output values for the i th output; n is the
 219 number of outputs; \bar{h}_i is the average value of the actual outputs. The R^2 ranges from 0 to
 220 1. The higher the R^2 value, the higher the fitting accuracy is obtained. It can be seen from
 221 Fig. 4 that the linear fitting between the air-entry value and the input soil property all shows a
 222 relatively low fitting accuracy ($R^2 \leq 0.173$). For the grain size distribution, the air-entry
 223 value decreases as the gravel or the sand content increases (Figs. 4a and 4b), while the
 224 air-entry value increases with the increase of the fines content (Fig. 4c). When the plasticity

225 index increases, the air-entry value increases accordingly (Fig. 4d). Regarding the initial
226 placement conditions, the air-entry value increases as the initial water content increases (Fig.
227 4e) or as the initial void ratio decreases (Fig. 4f). On the whole, the linear fittings in Fig. 4 are
228 not accurate enough to state the relationship between the air-entry value and a single soil
229 property. Hence, a more comprehensive prediction model is needed to connect the air-entry
230 value and the known input variables.

231 ***Model development***

232 To develop the prediction model accurately and to check the validity of the model, about
233 two-thirds of the samples (126 samples) in the supplementary database are chosen as the
234 training data, and the rest (63 samples) are used for testing. As the samples are ordered by
235 soil classification in the database, the samples with the line number as the multiple of three
236 are picked out for testing. In this way, both the training and the testing data cover all kinds of
237 soil classifications in the database. Note that the solid soil particles are constituted by gravel,
238 sand and fines, with the summation of their contents as 100%. Hence, the gravel content is
239 not considered in the training process. With the contents of sand and fines, the gravel content
240 can be calculated accordingly. It is also worth mentioning that the slight discrepancy of the
241 air-entry value induced by different testing methods (Lin and Cerato 2013; Sun et al. 2016) is
242 not considered during the model development. More studies are needed to clarify this issue.

243 To evaluate the prediction precision, three indicators of mean absolute error MAE , root
244 mean squared error $RMSE$ and coefficient of determination R^2 [see Eq. (4)] are
245 introduced as:

246
$$MAE = \frac{\sum_{i=1}^n |h_i - t_i|}{n} \quad (5)$$

247
$$RMSE = \sqrt{\frac{\sum_{i=1}^n (h_i - t_i)^2}{n}} \quad (6)$$

248 The lower values of *MAE* or *RMSE*, or the higher value of R^2 indicate the higher
 249 precision of the model.

250 For the MEP, the source code of Oltean (2004) is used for training the data, with various
 251 pre-set parameters (Oltean and Dumitrescu 2002). The population size defines the number of
 252 programs in the population. The number of generations is the number of calculations before
 253 the run of a program terminates. The crossover and the mutation probability indicate the
 254 probability of an offspring that is imposed on the crossover and the mutation operator,
 255 respectively. When the crossover type is set as uniform, the offspring genes are randomly
 256 taken from one parent to another. The code length represents the number of genes in each
 257 chromosome. The replication number is the number of runs or the number of developed
 258 chromosomes. In this study, the determination of the optimum code parameters follows the
 259 trial and error method used by Wang and Yin (2020), keeping the prediction *MAE* at the
 260 minimum level. Table 3 lists the initial parameter setting, as suggested by several previous
 261 studies (Oltean and Dumitrescu 2002; Shahnazari et al. 2010; Wang and Yin 2020). Using the
 262 optimum parameter setting, explicit expressions are developed accordingly.

263 To obtain the best solution using EPR, the number of transformed terms in the explicit
 264 expression, the values and the interval of elements in the exponent matrix need to be
 265 optimised. Table 3 lists the initial setting of the EPR parameters. Following Zhang et al.

266 (2020d), the interval of elements is set as 0.5. In the EPR calculation, the number of terms
267 needs to be pre-set at a fixed value. Then, the values of elements for each input variable is
268 determined by GA, with the parameters shown in Table 3. These parameters are verified to be
269 optimum for the EPR method (Zhang et al. 2020d). During each calculation, the values of
270 elements are randomly assigned to machine learning algorithms by GA. The performance of
271 the machine learning algorithms with these parameters is assessed by the fitness value until
272 the terminal condition is satisfied. As the explicit EPR expression contains the input
273 components with index and some of them serve as the denominators (Jin et al. 2019; Jin and
274 Yin 2020), the input variable with the value of 0 is not applicable and may influence the
275 prediction accuracy if directly used. To improve the accuracy of EPR and to guarantee the
276 feasibility of the EPR expression, the input variables expressed as percentages are written as
277 decimals, and then the input variables with the value of 0 are converted to a small value of
278 0.001. Note that the conversions take place for a minor proportion of the input variables of
279 sand content, fines content and plasticity index, showing negligible conversion value
280 compared to the variation range of these variables (see Table 1). After the calculation, the
281 optimum combination of the values of elements is determined. By comparing the prediction
282 *MAE* for the cases with a different number of terms, the optimum EPR explicit expression is
283 deduced.

284 Regarding the RF, the numbers of trees and features at each node need to be optimised.
285 In this study, the initial setting of these two parameters is listed in Table 3. GA is also used to
286 search for the optimum setting of these two parameters, using the method by Zhang et al.
287 (2020d). The parameter setting for GA is listed in Table 3, verified to be optimum for the RF

288 method (Zhang et al. 2020d). By randomly assigning the parameters to machine learning
 289 algorithms and evaluating the fitness value, the optimum number of trees and features is
 290 determined. Accordingly, the air-entry values can be predicted using this optimum parameter
 291 setting.

292 *Straightforward prediction of AEV*

293 Following the respective procedure of each machine learning algorithm, the optimum
 294 parameters are determined. For the MEP, the optimum combination of code parameters shows
 295 population size of 3000, code length of 50, crossover probability of 0.9, crossover type of
 296 uniform, mutation probability of 0.01, number of generation of 3000, function set of +, -, ×, /,
 297 pow and replication number of 10. The optimum number of transformed terms for the EPR is
 298 6. The optimum number of trees and features at each node for the RF is 116 and 3,
 299 respectively.

300 Using the optimum parameters, the MEP explicit expression for the air-entry value of
 301 compacted soils is derived as:

$$302 \quad AEV = 3(A_4 + A_3 A_4^{A_3}) + A_6 + A_4^{A_3} + e_0 \quad (7)$$

303 in which A_1 , A_2 , A_3 , A_4 , A_5 , A_6 and A_7 are parameters as:

$$304 \quad A_1 = (2e_0)^{A_7} \quad (8)$$

$$305 \quad A_2 = PI + w_0 + e_0^{e_0} \quad (9)$$

$$306 \quad A_3 = \frac{2PI}{C_S + C_F - A_7} \quad (10)$$

$$307 \quad A_4 = \frac{A_2 - 2e_0}{C_S + C_F} \quad (11)$$

$$308 \quad A_5 = \frac{2e_0 A_7 (C_S + C_F - A_7)}{A_1 + A_2} \quad (12)$$

309
$$A_6 = \frac{A_5 + e_0 A_2}{e_0^{e_0} + 2C_S e_0} \quad (13)$$

310
$$A_7 = 2e_0 + C_S \quad (14)$$

311 The EPR explicit expression for the air-entry value of compacted soils is deduced as:

312
$$AEV = 1.8116 - 0.0069 \left(\frac{C_{S1} PI_1}{w_{01}} \right)^3 \left(\frac{C_{F1}}{e_0} \right)^{2.5} + 29.2526 \left(\frac{C_{F1}}{e_0} \right)^{2.5} \frac{C_{S1}^3}{w_{01}^{0.5}} + 27.4022 \frac{C_{S1}^{0.5} C_{F1}^3}{e_0} \quad (15)$$

313
$$+ 470.6056 \frac{(C_{S1} PI_1)^3 w_{01}^{1.5}}{C_{F1}^{0.5} e_0} + 1.4266 \cdot 10^{-7} \frac{w_{01}^2}{(C_{S1} e_0)^3 PI_1^{0.5}} + 1.2272 \frac{C_{F1}^{2.5} PI_1 e_0}{C_{S1}^{0.5} w_{01}}$$

313 where C_{S1} , C_{F1} , PI_1 , w_{01} and e_{01} are sand content, fines content, plasticity index,
 314 initial water content and initial void ratio all in decimals, respectively. Note also that the
 315 explicit expression cannot be derived from the RF model.

316 Fig. 5 presents the comparison between the predicted and the reference air-entry values
 317 for the training data by the three prediction models. The prediction indicators are also shown
 318 in each figure for evaluation. It can be observed from this figure that the predicted air-entry
 319 values have a good agreement with the reference data for all the three models, with the R^2
 320 higher than 0.85. For the training data, the RF presents a higher prediction accuracy, while
 321 the prediction accuracies of MEP and EPR stay close to each other. Nevertheless, the high
 322 prediction accuracies of the training data indicate the appropriate development of all the
 323 prediction models.

324

325 **Results and discussions**

326 *Model validation*

327 The soil samples with the line number as the multiple of three are considered as the testing
 328 data in the supplementary database. Using the soil properties and the prediction models, the

329 air-entry values of these testing samples are calculated. Fig. 6 plots the comparison between
330 the predicted and the reference air-entry values of the testing data by each prediction model.
331 Note that to stay consistent with the training data, when using the EPR model, the input
332 variables with the percentage expressions are written as decimals and the input variables with
333 the value of 0 are also converted to 0.001. From Fig. 6, it is observed that although the RF
334 model shows a higher prediction accuracy for the training data than the other two models,
335 similar testing accuracy can be identified for the three models, with the R^2 varying from
336 0.84 to 0.88. On the whole, satisfactory prediction accuracies are verified by all the three
337 models for the air-entry values of compacted soils.

338 ***Monotonicity analysis***

339 To check the availability of the prediction models, the monotonicity analysis is performed,
340 using the method from Jin et al. (2019), Jin and Yin (2020), Wang and Yin (2020). To conduct
341 the monotonicity analysis, the investigated soil property changes, while the other properties
342 stay constant. With the input soil properties and the prediction model, the air-entry value of
343 compacted soils can be calculated. Table 4 lists the basic setting of the soil properties for the
344 monotonicity analysis. These soil properties are chosen as their respective mean values from
345 the supplementary database (see Table 1). During the monotonicity analysis, the variation of
346 the input soil property cannot exceed the threshold defined by their respective maximum and
347 minimum values (see Table 1). Note that as the soils are constituted by gravel, sand and fines,
348 the content of each material influences each other and the summation of their contents equals
349 to 100%. Thus, for the monotonicity analysis of the sand and the fines content, the gravel
350 content is fixed at 3.8%. For the monotonicity analysis of the other soil properties, the sand

351 and the fines content are set as the values shown in Table 4.

352 Fig. 7 depicts the results of the monotonicity analysis using the prediction models,
353 showing the variation of the predicted air-entry value with each input soil property. In general,
354 the monotonic variation trend predicted by the three models agrees with each other, also
355 showing consistency with that by the basic linear fitting in Fig. 4. Despite the overall
356 consistent trend, the smooth correlation between the RF predicted output variable and the
357 input variable is difficult to obtain, because the RF model is developed strictly following the
358 actual data from the database with no smooth relationships (Fig. 4). The predicted air-entry
359 value increases with the decrease of the sand content or the increase of the fines content (Figs.
360 7a and 7b). As the plasticity index increases, the air-entry value increases accordingly (Fig. 7c;
361 available for the RF prediction only when $PI > 15$). In terms of the initial placement
362 conditions, the air-entry value increases when the initial water content increases or when the
363 initial void ratio decreases (Figs. 7d and 7e). In spite of the general consistency of the
364 monotonicity results for each model, some differences still exist, especially for the results
365 regarding the initial placement conditions. With the increase of the initial water content, the
366 increasing amplitude of the predicted air-entry value by EPR and RF is relatively small (Fig.
367 7d). Besides, the increasing trend only shows for the EPR model when the initial water
368 content is higher than about 12%. With the initial void ratio lower than around 0.5, the
369 monotonicity predicted by the RF shows the opposite trend (Fig. 7e). However, the general
370 consistent monotonic variation trend of the predicted air-entry value with each input soil
371 property between the prediction results and the original database directly verifies the validity
372 of the developed models.

373 ***Sensitivity analysis***

374 To have a better understanding of the contribution of the input soil property on the predicted
375 air-entry value, the sensitivity analysis is conducted on the whole database. For a specific
376 input variable x_i , the sensitivity R_{sen} is determined as (Wang and Yin 2020):

377
$$R_{sen} = \frac{\sum_{i=1}^N (x_i t_i)}{\sqrt{\sum_{i=1}^N x_i^2 \sum_{i=1}^N t_i^2}} \quad (16)$$

378 where t_i is the predicted output air-entry value using the proposed prediction models; N
379 is the number of soil samples in the supplementary database ($N = 189$). The sensitivity R_{sen}
380 ranges from 0 to 1, indicating the relevance between the predicted air-entry value and each
381 input soil property. With the R_{sen} value closer to 1, the specific input soil property has a
382 more remarkable influence on the predicted air-entry value. Note that to stay consistent with
383 the prediction setup using the EPR algorithm, the input variables with the percentage
384 expressions are written as decimals, also with the input variables of 0 value converted to
385 0.001 in this analysis.

386 Fig. 8 presents the distribution of the sensitivity value for each input soil property on the
387 predicted air-entry value using the three prediction models. For a specific input soil property
388 except for the plasticity index, the sensitivity shows a slightly higher value by the RF model,
389 while the sensitivity value by the MEP model is relatively lower. However, this difference is
390 not significant for a given input soil property. On the whole, the sand content has the least
391 influence on the prediction of the air-entry value by all the three models. By contrast, the
392 sensitivity value of the fines content and the initial water content rank the first- and the
393 second-highest, respectively. Except for the sand content, the sensitivity value of the other

394 four soil properties stays close to each other, showing the values between 0.6 and 0.8.

395 ***Robustness analysis***

396 To validate the prediction model, the robustness analysis is another key aspect (Jin et al. 2019;
397 Jin and Yin 2020), to guarantee that the output values are reasonable with the appropriate
398 input variables. To assess the robustness of the prediction model, a robustness ratio r is
399 defined as:

$$400 \quad r = \frac{\text{Samples in the reasonable range}}{\text{Total testing samples}} \quad (17)$$

401 From the present database, the reasonable range for the air-entry value of compacted soils
402 locate between 0.06 kPa and 100 kPa (Table 1). To generate the testing samples, the five
403 input variables (C_S , C_F , PI , w_0 , e_0) are first assumed to be independent to each other
404 and to obey the lognormal distribution (Jin et al. 2019; Jin and Yin 2020). From the statistics
405 in Table 1, 80,000 testing samples are randomly generated using the mean and the standard
406 deviation of each input variable with the lognormal distribution. Then, the values exceeding
407 the minimum and the maximum thresholds (Table 1) are deleted. Besides, the generated
408 samples with the summation of sand and fines contents exceeding 100% are deleted. Finally,
409 10,000 samples are chosen as the testing samples, still showing close mean value and
410 standard deviation for each input variable. Using the 10,000 samples, the robustness analysis
411 is conducted for each prediction model.

412 Fig. 9 depicts the distribution of the predicted air-entry values for each prediction model,
413 showing the robustness ratio in the legend. It can be seen that the majority of the predicted
414 air-entry values locate in the range from 0 kPa to 20 kPa. For the models of MEP and EPR,
415 some negative values exist. For the three models, some air-entry values exceed 100 kPa.

416 However, the robustness ratios are higher than 98% for all the three methods, suggesting the
417 feasibility of the prediction models. By comparing the robustness ratio of each model, the RF
418 and the MEP models are slightly more robust.

419 *Limitations of the present algorithms*

420 In this study, three representative machine learning algorithms are used to predict the
421 air-entry values of compacted soils. All the prediction models by the present algorithms show
422 relatively high accuracies, consistent monotonicity and strong robustness. However, there are
423 still some limitations for each algorithm. The prediction model by MEP has a good
424 performance for the complex and nonlinear problems. But the length of the expressions and
425 the feasibility of the model needs to be balanced for practical applications. The EPR model is
426 more convenient to use, with only one explicit expression. While due to the expression form
427 of EPR, the limitation exists when the input variables include some values equalling to 0. The
428 prediction accuracy of the training data and the robustness ratio of RF are the highest in the
429 present study. Nevertheless, no explicit expressions can be generated using this algorithm,
430 causing some inconvenience for the applications. In addition, the smooth monotonicity
431 correlation between the predicted output and each input variable is difficult to obtain by RF.
432 Hence, according to the specific scenario, a suitable machine learning algorithm needs to be
433 selected to develop the prediction model with high accuracy, feasibility and convenience.

434

435 **Conclusions**

436 In this study, three alternative straightforward prediction models for the air-entry value of
437 compacted soils have been developed using three representative machine learning algorithms:

438 multi expression programming (MEP), evolutionary polynomial regression (EPR) and
439 random forest (RF). A large database with a wide range of soil classifications has been
440 collected, covering clay, silt, sand and gravel.

441 The optimum parameter setting for each algorithm was determined firstly. Using their
442 optimum parameter settings, three prediction models for the air-entry value of compacted
443 soils were developed based on the training data, showing reasonable prediction accuracies.
444 By comparison between the predicted air-entry values and the reference ones of the testing
445 data, the prediction precisions were validated for all the three models.

446 The monotonicity, the sensitivity and the robustness analysis were conducted using the
447 proposed models, showing consistent results among different models. From the monotonicity
448 analysis, the predicted air-entry value increases monotonically as the fines content, the
449 plasticity index or the initial water content increases, while it decreases with the increase of
450 either the sand content or the initial void ratio. The variation trend of the monotonicity
451 analysis shows a good agreement with the original database. The sensitivity analysis indicates
452 that the sand content has a relatively lower relevance on the prediction of the air-entry value,
453 whereas the influence of the other four soil properties stays close to each other at a high value.
454 From the robustness analysis, the high robustness ratios of the predicted air-entry values
455 strongly support the feasibility of the three prediction models.

456 Although three models have slight differences in terms of performance, the MEP model
457 is first recommended due to its advantages of both explicit formulation and good
458 monotonicity. For convenient engineering practice, the EPR model is also recommended
459 because of the simple explicit formulation, despite of its deficiency to treat the input variable

460 with 0 value. Without the explicit expression, the RF model is not as convenient as the other
461 two models, but it can be used for the cases demanding a higher prediction accuracy.

462

463 **Notations**

464 a_j Adjustable parameter for j th item

465 a_0 Optional bias

466 $A_1, A_2, A_3, A_4, A_5, A_6, A_7$ Parameters for the prediction model

467 AEV Air-entry value

468 C_F Fines content

469 C_{F1} Fines content in decimal

470 C_G Gravel content

471 C_S Sand content

472 C_{S1} Sand content in decimal

473 e_0 Initial void ratio

474 e_{01} Initial void ratio in decimal

475 $\mathbf{ES}_{m \times k}$ Exponent matrix

476 h_i Actual output variable

477 \bar{h}_i Average value of actual outputs

478 k Total number of input variables

479 m Number of transformed terms

480 MAE Mean absolute error

481 n Number of outputs

482 N Number of soil samples

483 n_t Number of trees

484 PI Plasticity index

485 PI_1 Plasticity index in decimal

486 r Robustness ratio

487 R^2 Coefficient of determination

488 R_{sen} Sensitivity

489 $RMSE$ Root mean squared error

490 S_r Degree of saturation

491 t_i Predicted output variable

492 $t_i(\mathbf{x})$ Predicted output for a tree with an input vector \mathbf{x}

493 \mathbf{t} , t Predicted output

494 w_0 Initial water content

495 w_{01} Initial water content in decimal

496 x_i Input variable

497 \mathbf{x}_i i th input variable

498 \mathbf{x} Input vector

499 z_j j th transformed variable

500 ψ Matric suction

501

502 **Acknowledgement**

503 The financial supports provided by the RIF project (Grant No. R5037-18F) from Research
 504 Grants Council (RGC) of Hong Kong are gratefully acknowledged.

505 **References**

- 506 Agus, S.S., Leong, E.C., Rahardjo, H., 2001. Soil-water characteristic curves of Singapore
507 residual soils. *Geotech. Geol. Eng.* 19, 285-309.
- 508 Amadi, A.A., Osinubi, K.J., 2016. Soil-water characteristic curves for compacted lateritic
509 soil-bentonite mixtures developed for landfill liner applications. In: *Geo-Chicago 2016:*
510 *Sustainability and Resiliency in Geotechnical Engineering*, Chicago, Illinois, US, pp.
511 488-497.
- 512 ASTM., 2017. Standard practice for classification of soils for engineering purposes (unified
513 soil classification system). ASTM D2487-17, West Conshohocken, PA.
- 514 Birle, E., Heyer, D., Vogt, N., 2008. Influence of the initial water content and dry density on
515 the soil–water retention curve and the shrinkage behavior of a compacted clay. *Acta*
516 *Geotech.* 3(3), 191.
- 517 Breiman, L., 1996. Bagging Predictors. *Mach. Learn.* 24 (2), 123-140.
- 518 Brooks, R., Corey, T., 1964. Hydraulic properties of porous media. *Hydrology Papers*,
519 Colorado State University, 24, 37.
- 520 Chen, R., Tan, R., Chen, Z., Ping, Y., Mei, Z., 2020a. Influence of degree of compaction on
521 unsaturated hydraulic properties of a compacted completely decomposed granite.
522 *Geofluids* 2020, 7615361.
- 523 Chen, R.P., Qi, S., Wang, H.L., Cui, Y. J., 2019. Microstructure and hydraulic properties of
524 coarse-grained subgrade soil used in high-speed railway at various compaction degrees. *J.*
525 *Mater. Civil Eng.* 31 (12), 04019301.
- 526 Chen, W.B., Feng, W.Q., Yin, J.H., 2020b. Effects of water content on resilient modulus of a
527 granular material with high fines content. *Constr. Build. Mat.* 236, 117542.
- 528 Chen, W.B., Feng, W.Q., Yin, J.H., Chen, J.M., Borana, L., Chen, R.P., 2020c. New model for
529 predicting permanent strain of granular materials subjected to cyclic loadings. *J. Geotech.*
530 *Geoenviron. Eng.* 146 (9), 04020084.
- 531 Chen, W.B., Liu, K., Feng, W.Q., Borana, L., Yin, J.H., 2020d. Influence of matric suction on
532 nonlinear time-dependent compression behavior of a granular fill material. *Acta Geotech.*
533 15 (3), 615-633.
- 534 Cheng, Z.L., Zhou, W.H., Ding, Z., Guo, Y.X., 2020a. Estimation of spatiotemporal response
535 of rooted soil using a machine learning approach. *J. Zhejiang Uni.-SCI. A (Appl. Phys.*
536 *Eng.)*, 21 (6), 462-477.
- 537 Cheng, Z.L., Zhou, W.H., Garg, A., 2020b. Genetic programming model for estimating soil
538 suction in shallow soil layers in the vicinity of a tree. *Eng. Geol.* 268, 105506.
- 539 Cuceoglu, F., 2016. An experimental study on soil water characteristics and hydraulic
540 conductivity of compacted soils. MSc dissertation, Virginia Tech.
- 541 Cui, Y.J., Yahia-Aissa, M., Delage, P., 2002. A model for the volume change behavior of
542 heavily compacted swelling clays. *Eng. Geol.* 64 (2-3), 233-250.

- 543 de Freitas, J.B., de Rezende, L.R., de FN Gitirana Jr, G., 2020. Prediction of the resilient
544 modulus of two tropical subgrade soils considering unsaturated conditions. *Eng. Geol.*
545 105580.
- 546 Delage, P., Howat, M.D., Cui, Y.J., 1998. The relationship between suction and swelling
547 properties in a heavily compacted unsaturated clay. *Eng. Geol.* 50 (1-2), 31-48.
- 548 Fattah, M.Y., Salim, N.M., Irshayid, E.J., 2017. Determination of the soil–water
549 characteristic curve of unsaturated bentonite–sand mixtures. *Environ. Earth Sci.* 76 (5),
550 201.
- 551 Fredlund, D.G., Rahardjo, H., Fredlund, M.D., 2012. *Unsaturated soil mechanics in*
552 *engineering practice*. John Wiley & Sons, Inc.
- 553 Fredlund, D.G., Xing, A., 1994. Equations for the soil-water characteristic curve, *Can.*
554 *Geotech. J.* 31 (3), 521–532.
- 555 Fredlund, M.D., Wilson, G.W., Fredlund, D.G., 2002. Use of the grain-size distribution for
556 estimation of the soil-water characteristic curve. *Can. Geotech. J.* 39 (5), 1103-1117.
- 557 Gallage, C.P.K., Uchimura, T., 2010. Effects of dry density and grain size distribution on
558 soil-water characteristic curves of sandy soils. *Soils. Found.* 50 (1), 161–172.
- 559 Giustolisi, O., Savic, D.A., 2006. A symbolic data-driven technique based on evolutionary
560 polynomial regression. *J. Hydroinform.* 8 (4), 235-237.
- 561 Han, K.K., Rahardjo, H., and Broms, B.B. 1995. Effect of hysteresis on the shear strength of
562 a residual soil. In: *Proceedings of the First International Conference on Unsaturated soil*
563 *(UNSAT 95), Paris, France*, pp. 6–8.
- 564 Han, Z., Vanapalli, S.K., 2016. Relationship between resilient modulus and suction for
565 compacted subgrade soils. *Eng. Geol.* 211, 85-97.
- 566 Hashem, E.B., Houston, S.L., 2016. Volume change consideration in determining unsaturated
567 soil properties for geotechnical applications. *Int. J. Geomech.* 16 (6), D4015003.
- 568 He, Y., Ye, W.M., Chen, Y.G., Cui, Y.J., 2019. Effects of K⁺ solutions on swelling behavior of
569 compacted GMZ bentonite. *Eng. Geol.* 249, 241-248.
- 570 Ho, T.K., 1998. The random subspace method for constructing decision forests. *IEEE T.*
571 *Pattern Anal.*, 20 (8), 832-844.
- 572 Huang, S., Barbour, S.L., Fredlund, D.G., 1998. Development and verification of a
573 coefficient of permeability function for a deformable unsaturated soil. *Can. Geotech. J.*
574 35 (3), 411-425.
- 575 Indrawan, I.G.B., Rahardjo, H., Leong, E.C., 2006. Effects of coarse-grained materials on
576 properties of residual soil. *Eng. Geol.* 82 (3), 154-164.
- 577 Jiang, Y., Chen, W., Wang, G., Sun, G., Zhang, F., 2017. Influence of initial dry density and
578 water content on the soil–water characteristic curve and suction stress of a reconstituted
579 loess soil. *B. Eng. Geol. Environ.* 76 (3), 1085-1095.
- 580 Jin, Y.F., Yin, Z.Y., 2020. An intelligent multi-objective EPR technique with multi-step model

581 selection for correlations of soil properties. *Acta Geotech.* 15, 2053–2073.

582 Jin, Y.F., Yin, Z.Y., Zhou, W.H., Yin, J.H., Shao, J.F., 2019. A single-objective EPR based
583 model for creep index of soft clays considering L2 regularization. *Eng. Geol.* 248,
584 242-255.

585 Johari, A., Habibagahi, G., Ghahramani, A., 2006. Prediction of soil–water characteristic
586 curve using genetic programming. *J. Geotech. Geoenviron. Eng.* 132 (5), 661-665.

587 Khalili, N., Geiser, F., Blight, G.E., 2004. Effective stress in unsaturated soils: Review with
588 new evidence. *Int. J. Geomech.* 4 (2), 115-126.

589 Krisdani, H., Rahardjo, H., Leong, E.C., 2008. Effects of different drying rates on shrinkage
590 characteristics of a residual soil and soil mixtures. *Eng. Geol.* 102 (1-2), 31-37.

591 Li, X., 2009. Dual-porosity structure and bimodal hydraulic property functions of coarse
592 granular soils. Ph. D. Thesis, Hong Kong University of Science and Technology, Hong
593 Kong.

594 Li, X., Li, J. H., Zhang, L.M., 2014. Predicting bimodal soil–water characteristic curves and
595 permeability functions using physically based parameters. *Comput. Geotech.* 57, 85-96.

596 Liaw, A., Wiener, M., 2002. Classification and regression by random forest. *R News*, 23 (23),
597 18-21.

598 Lin, B., Cerato, A.B., 2013. Hysteretic soil water characteristics and cyclic swell–shrink paths
599 of compacted expansive soils. *B. Eng. Geol. Environ.* 72 (1), 61-70.

600 Liu, K., Yin, J., Chen, W., Feng, W.Q., Zhou, C., 2020. The stress–strain behaviour and
601 critical state parameters of an unsaturated granular fill material under different suctions.
602 *Acta Geotech.* in press. doi: 10.1007/s11440-020-00973-1.

603 Miller, C.J., Yesiller, N., Yaldo, K., Merayyan, S., 2002. Impact of soil type and compaction
604 conditions on soil water characteristic. *J. Geotech. Geoenviron. Eng.* 128 (9), 733-742.

605 Mirzaii, A., Yasrobi, S.S., 2012. Influence of initial dry density on soil-water characteristics
606 of two compacted soils. *Geotechnique Let.* 2 (4), 193-198.

607 Montanez, J.E.C., 2002. Suction and volume changes of compacted sand-bentonite mixtures.
608 PhD dissertation, University of London.

609 Njock, P.G.A., Shen, S.L., Zhou, A., Lyu, H.M., 2020. Evaluation of soil liquefaction using
610 AI technology incorporating a coupled ENN/t-SNE model. *Soil Dyn. Earthq. Eng.* 130,
611 105988.

612 Ng, C.W.W., Pang, Y.W., 2000. Experimental investigations of the soil-water characteristics
613 of a volcanic soil. *Can. Geotech. J.* 37 (6), 1252-1264.

614 Oh, S., Lu, N., Kim, Y.K., Lee, S.J., Lee, S.R., 2012. Relationship between the soil-water
615 characteristic curve and the suction stress characteristic curve: Experimental evidence
616 from residual soils. *J. Geotech. Geoenviron. Eng.* 138 (1), 47-57.

617 Oltean, M., 2004. Multi expression programming source code. Available at:
618 https://www.mepx.org/source_code.html

- 619 Oltean, M., Dumitrescu, D., 2002. Multi expression programming. Technical Report,
620 UBB-01-2002, Babes-Bolyai University, Cluj-Napoca.
- 621 Oltean, M., Grosan, C., 2003. A comparison of several linear genetic programming
622 techniques. *Complex Syst.* 14 (4), 285-314.
- 623 Osinubi, K.J., Nwaiwu, C.M., 2006. Design of compacted lateritic soil liners and covers. *J.*
624 *Geotech. Geoenviron. Eng.* 132 (2), 203-213.
- 625 Priono, Rahardjo, H., Chatterjea, K., Leong, E.C., Wang, J.Y., 2016. Effect of hydraulic
626 anisotropy on soil–water characteristic curve. *Soils Found.* 56 (2), 228-239.
- 627 Puppala, A.J., Punthutaecha, K., Vanapalli, S.K., 2006. Soil-water characteristic curves of
628 stabilized expansive soils. *J. Geotech. Geoenviron. Eng.* 132 (6), 736-751.
- 629 Rahardjo, H., Indrawan, I.G.B., Leong, E.C., Yong, W.K., 2008. Effects of coarse-grained
630 material on hydraulic properties and shear strength of top soil. *Eng. Geol.* 101 (3-4),
631 165-173.
- 632 Rahardjo, H., Satyanaga, A., D'Amore, G.A., Leong, E.C., 2012. Soil–water characteristic
633 curves of gap-graded soils. *Eng. Geol.* 125, 102-107.
- 634 Salager, S., Nuth, M., Ferrari, A., Laloui, L., 2013. Investigation into water retention
635 behaviour of deformable soils. *Can. Geotech. J.* 50 (2), 200-208.
- 636 Sarir, P., Shen, S.L., Wang, Z.F., Chen, J., Horpibulsuk, S., Pham, B.T., 2019. Optimum
637 model for bearing capacity of concrete-steel columns with AI technology via
638 incorporating the algorithms of IWO and ABC. *Eng. Comput.* doi:
639 10.1007/s00366-019-00855-5
- 640 Satyanaga, A., Rahardjo, H., Zhai, Q., 2017. Estimation of unimodal water characteristic
641 curve for gap-graded soil. *Soils Found.* 57 (5), 789-801.
- 642 Shahnazari, H., Dehnavi, Y., Alavi, A.H., 2010. Numerical modeling of stress–strain behavior
643 of sand under cyclic loading. *Eng. Geol.* 116 (1-2), 53-72.
- 644 Shen, S., Njock, P.G.A., Zhou, A., Lyu, H.M., 2020. Dynamic prediction of jet grouted
645 column diameter in soft soil using Bi-LSTM deep learning. *Acta Geotech.* in press. doi:
646 10.1007/s11440-020-01005-8
- 647 Sun, D.A., Cui, H., Sun, W., 2009. Swelling of compacted sand–bentonite mixtures. *Appl.*
648 *Clay Sci.* 43 (3-4), 485-492.
- 649 Sun, D.A., Gao, Y., 2015. Water retention behaviour of soils with different preparations.
650 *Chinese J. Geotech. Eng.* 37 (1), 91-97 (in Chinese).
- 651 Sun, D.A., Gao, Y., Zhou, A.N., Sheng, D.C., 2016. Soil-water retention curves and
652 microstructures of undisturbed and compacted Guilin lateritic clay. *Bull. Eng. Geol.*
653 *Environ.* 75 (2), 781–791.
- 654 Sun, D.A., Liu, W.J., Lü, H.B., 2014. Soil-water characteristic curve of Guilin lateritic clay.
655 *Rock Soil Mech.* 35 (12), 3345-3351 (in Chinese).
- 656 Sun, D.A., Sheng, D.C., Cui, H.B., Li, J., 2006. Effect of density on the soil-water-retention

657 behaviour of compacted soil. In Proceedings of the Fourth International Conference on
658 Unsaturated Soils, Carefree, Arizona, US, pp. 1338-1347.

659 Thu, T.M., Rahardjo, H., and Leong, E.C. 2006. Shear strength and pore-water pressure
660 characteristics during constant water content triaxial tests. *J. Geotech. Geoenviron. Eng.*
661 132(3): 411–419.

662 Thu, T.M., Rahardjo, H., Leong, E.C., 2007. Soil-water characteristic curve and consolidation
663 behavior for a compacted silt. *Can. Geotech. J.* 44 (3), 266-275.

664 Tinjum, J.M., Benson, C.H., Blotz, L.R., 1997. Soil-water characteristic curves for
665 compacted clays. *J. Geotech. Geoenviron. Eng.* 123 (11), 1060-1069.

666 van Genuchten, M.T., 1980. A closed-form equation for predicting the hydraulic conductivity
667 of unsaturated soils, *J. Soil Sci. Soc. America* 44, 892–898.

668 Vanapalli, S.K., Fredlund, D.G., Pufahl, D.E., 1999. The influence of soil structure and stress
669 history on the soil–water characteristics of a compacted till. *Geotechnique* 49 (2),
670 143-159.

671 Wang, H.L., Cui, Y.J., Lamas-Lopez, F., Dupla, J.C., Canou, J., Calon, N., Saussine, G.,
672 Aïmediou, P., Chen, R.P., 2017. Effects of inclusion contents on resilient modulus and
673 damping ratio of unsaturated track-bed materials. *Can. Geotech. J.* 54 (12), 1672-1681.

674 Wang, H.L., Cui, Y.J., Lamas-Lopez, F., Dupla, J.C., Canou, J., Calon, N., Saussine, G.,
675 Aïmediou, P., Chen, R. P., 2018. Permanent deformation of track-bed materials at various
676 inclusion contents under large number of loading cycles. *J. Geotech. Geoenviron. Eng.*
677 144 (8), 04018044.

678 Wang, H.L., Chen, R.P., 2019. Estimating static and dynamic stresses in
679 geosynthetic-reinforced pile-supported track-bed under train moving loads. *J. Geotech.*
680 *Geoenviron. Eng.* 145 (7), 04019029.

681 Wang, H.L., Chen, R.P., Cheng, W., Qi, S., Cui, Y.J., 2019a. Full-scale model study on
682 variations of soil stress in the geosynthetic-reinforced pile-supported track-bed with
683 water level change and cyclic loading. *Can. Geotech. J.* 56 (1), 60–68.

684 Wang, H.L., Chen, R.P., Liu, Q.W., Kang, X., 2019b. Investigation on geogrid reinforcement
685 and pile efficacy in geosynthetic-reinforced pile-supported track-bed. *Geotext.*
686 *Geomembranes* 47 (6), 755-766.

687 Wang, H.L., Yin, Z.Y., 2020. High performance prediction of soil compaction parameters
688 using multi expression programming. *Eng. Geol.* 276, 105758.

689 Yang, H., Rahardjo, H., Leong, E. C., Fredlund, D.G., 2004. Factors affecting drying and
690 wetting soil-water characteristic curves of sandy soils. *Can. Geotech. J.* 41 (5), 908-920.

691 Ye, W.M., Cui, Y.J., Qian, L.X., Chen, B., 2009. An experimental study of the water transfer
692 through confined compacted GMZ bentonite. *Eng. Geol.* 108 (3-4), 169-176.

693 Zhai, Q., Rahardjo, H., Satyanaga, A., Dai, G., 2020. Estimation of the soil-water
694 characteristic curve from the grain size distribution of coarse-grained soils. *Eng. Geol.*
695 267, 105502.

- 696 Zhang, F., Chen, X., 2009. Experimental study on characteristics of deformation and strength
697 of unsaturated clay. *Chinese J. Rock Mech. Eng.* 28 (2), 3808-3814 (in Chinese).
- 698 Zhang, P., Chen, R.P., Wu, H.N., 2019. Real-time analysis and regulation of EPB shield
699 steering using Random Forest. *Automat. Constr.* 106, 102860.
- 700 Zhang, P., Jin, Y.-F., Yin, Z.-Y., Yang, Y., 2020a. Random forest based artificial intelligent
701 model for predicting failure envelopes of caisson foundations in sand. *Appl. Ocean Res.*
702 101, 102223.
- 703 Zhang, P., Wu, H.-N., Chen, R.-P., Chan, T.H.T., 2020b. Hybrid meta-heuristic and machine
704 learning algorithms for tunneling-induced settlement prediction: A comparative study.
705 *Tunnell. Undergr. Space Technol.* 99, 103383.
- 706 Zhang, P., Yin, Z.-Y., Jin, Y.-F., Chan, T.H.T., 2020c. A novel hybrid surrogate intelligent
707 model for creep index prediction based on particle swarm optimization and random
708 forest. *Eng. Geol.* 265, 105328.
- 709 Zhang, P., Yin, Z.Y., Jin, Y.F., Chan, T., Gao, F.P., 2020d. Intelligent modelling of clay
710 compressibility using hybrid meta-heuristic and machine learning algorithms. *Geosci.*
711 *Front. in press.* doi: 10.1016/j.gsf.2020.1002.1014.
- 712 Zhao, L.S., Zhou, W.H., Fatahi, B., Li, X.B., Yuen, K.V., 2016. A dual beam model for
713 geosynthetic-reinforced granular fill on an elastic foundation. *Appl. Math. Model.* 40
714 (21-22), 9254-9268.
- 715 Zhou, B.C., Kong, L.W., 2011. Effect of volume changes on soil-water characteristics of
716 unsaturated expansive soil. *J. Hydr. Eng.* 42 (10), 1152-1160 (in Chinese).
- 717 Zhou, W.H., Xu, X., Garg, A. 2016. Measurement of unsaturated shear strength parameters of
718 silty sand and its correlation with unconfined compressive strength. *Measurement* 93,
719 351-358.
- 720 Zhou W.H., Yuen K.V., Tan F., 2014. Estimation of soil-water characteristic curve for
721 granular soils with different initial dry densities. *Eng. Geol.* 179, 1–9.

722

Table 1. Descriptive statistics of each variable

Variable	Minimum	Maximum	Mean	Standard deviation
C_G (%)	0	86.5	3.8	13.41
C_S (%)	0	100	40.3	29.93
C_F (%)	0	100	55.9	30.38
PI (%)	0	88	18.9	16.45
w_0 (%)	0.5	48.6	18.0	9.55
e_0	0.24	1.55	0.73	0.28
AEV (kPa)	0.06	100	15.37	17.82

723

724

Table 2. Air-entry value of different soils in the database

Soil	Minimum (kPa)	Maximum (kPa)	Mean (kPa)	Median (kPa)
Clay	2.5	100	18.33	14.63
Silt	0.43	73	22.62	12
Sand	0.2	71.16	10.44	4.07
Gravel	0.06	14.65	3.59	0.77

725

Table 3. Parameter setting for determination of the optimum combination

Algorithm	Parameter	Setting
All algorithms	Terminal set	C_S, C_F, PI, w_0, e_0
MEP	Population size	1000, 2000, 3000
	Number of generation	1000, 2000, 3000
	Crossover probability	0.1, 0.5, 0.9
	Crossover type	Uniform
	Mutation probability	0.01, 0.1, 0.9
	Code length	50, 100
	Function set	+, -, ×, /, pow
	Replication number	10
EPR	Number of terms	2-10
	Values of elements	[-3, 3]
	Interval of elements	0.5
RF	Number of trees	1-500
	Number of features	1-5
GA for EPR and RF	Population size	20
	Number of generation	500
	Crossover probability	0.7
	Mutation probability	0.1

Table 4. Basic setting of input soil properties for the monotonicity analysis

Soil parameter	Value
C_G (%)	3.8
C_S (%)	40.3
C_F (%)	55.9
PI (%)	18.9
w_0 (%)	18.0
e_0	0.73

List of Figures

- Fig. 1. Conceptual soil-water characteristic curve (I: boundary effect zone; II: transition zone; III: residual zone)
- Fig. 2. Frequency histograms of the variables: (a) gravel content; (b) sand content; (c) fines content; (d) plasticity index; (e) initial water content; (f) initial void ratio; (g) soil classification; (h) air-entry value
- Fig. 3. Comparison about air-entry value of soils with different classifications in the database: (a) specific classification; (b) general classification
- Fig. 4. Basic linear fittings between air-entry value and each input soil property: (a) gravel content; (b) sand content; (c) fines content; (d) plasticity index; (e) initial water content; (f) initial void ratio
- Fig. 5. Comparison between predicted and reference air-entry values for the training data: (a) MEP; (b) EPR; (c) RF
- Fig. 6. Comparison between predicted and reference air-entry values for the testing data: (a) MEP; (b) EPR; (c) RF
- Fig. 7. Monotonicity analysis of the predicted air-entry value versus (a) sand content; (b) fines content; (c) plasticity index; (d) initial water content; (e) initial void ratio
- Fig. 8. Sensitivity analysis about the relevance of the input variables on the predicted air-entry value
- Fig. 9. Distribution of the predicted air-entry value in robustness analysis

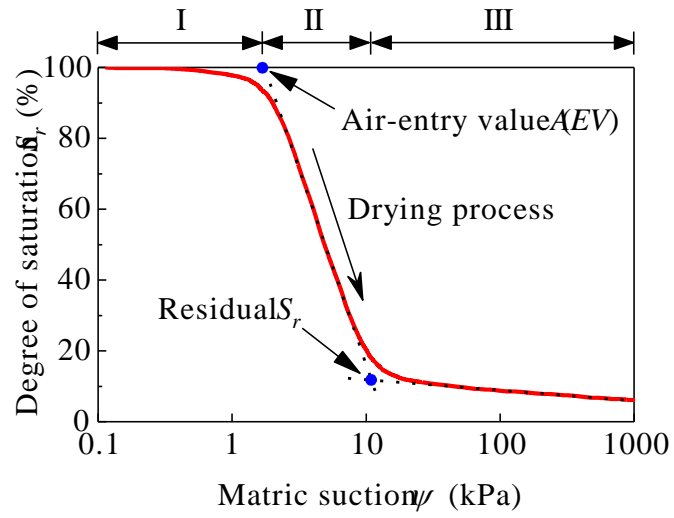


Fig. 1. Conceptual soil-water characteristic curve (I: boundary effect zone; II: transition zone; III: residual zone)

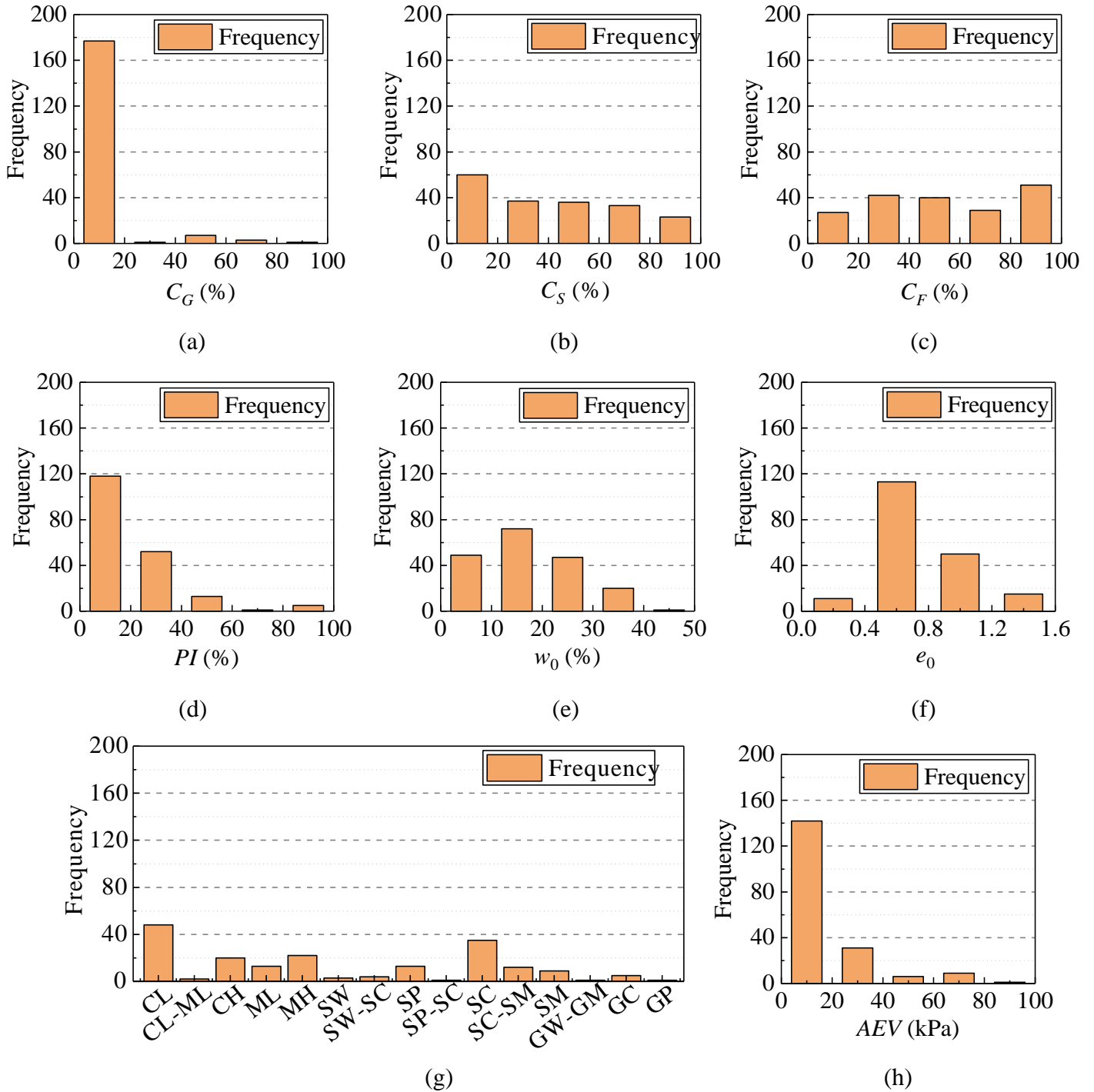


Fig. 2. Frequency histograms of the variables: (a) gravel content; (b) sand content; (c) fines content; (d) plasticity index; (e) initial water content; (f) initial void ratio; (g) soil classification; (h) air-entry value

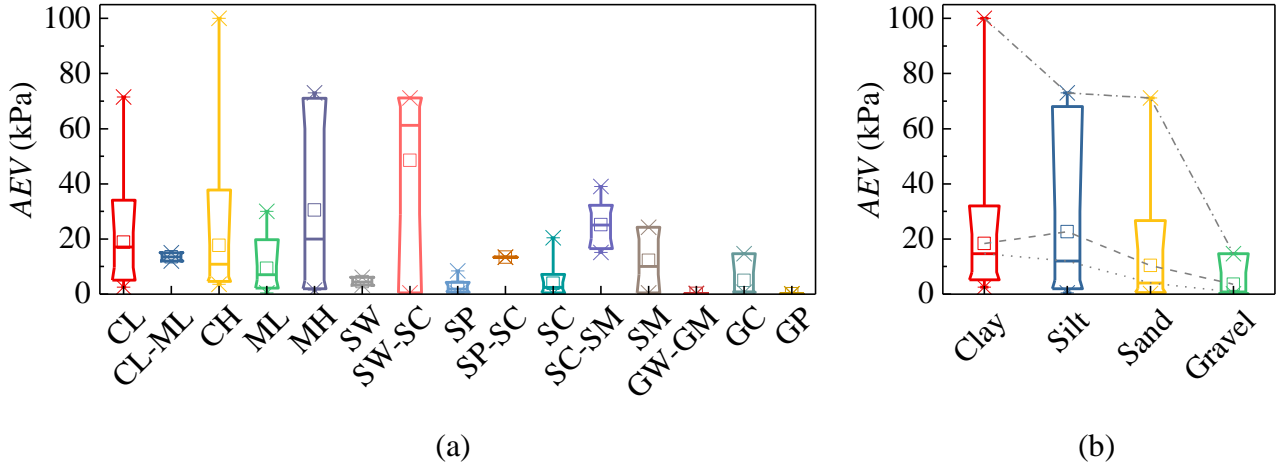


Fig. 3. Comparison about air-entry value of soils with different classifications in the database:
 (a) specific classification; (b) general classification

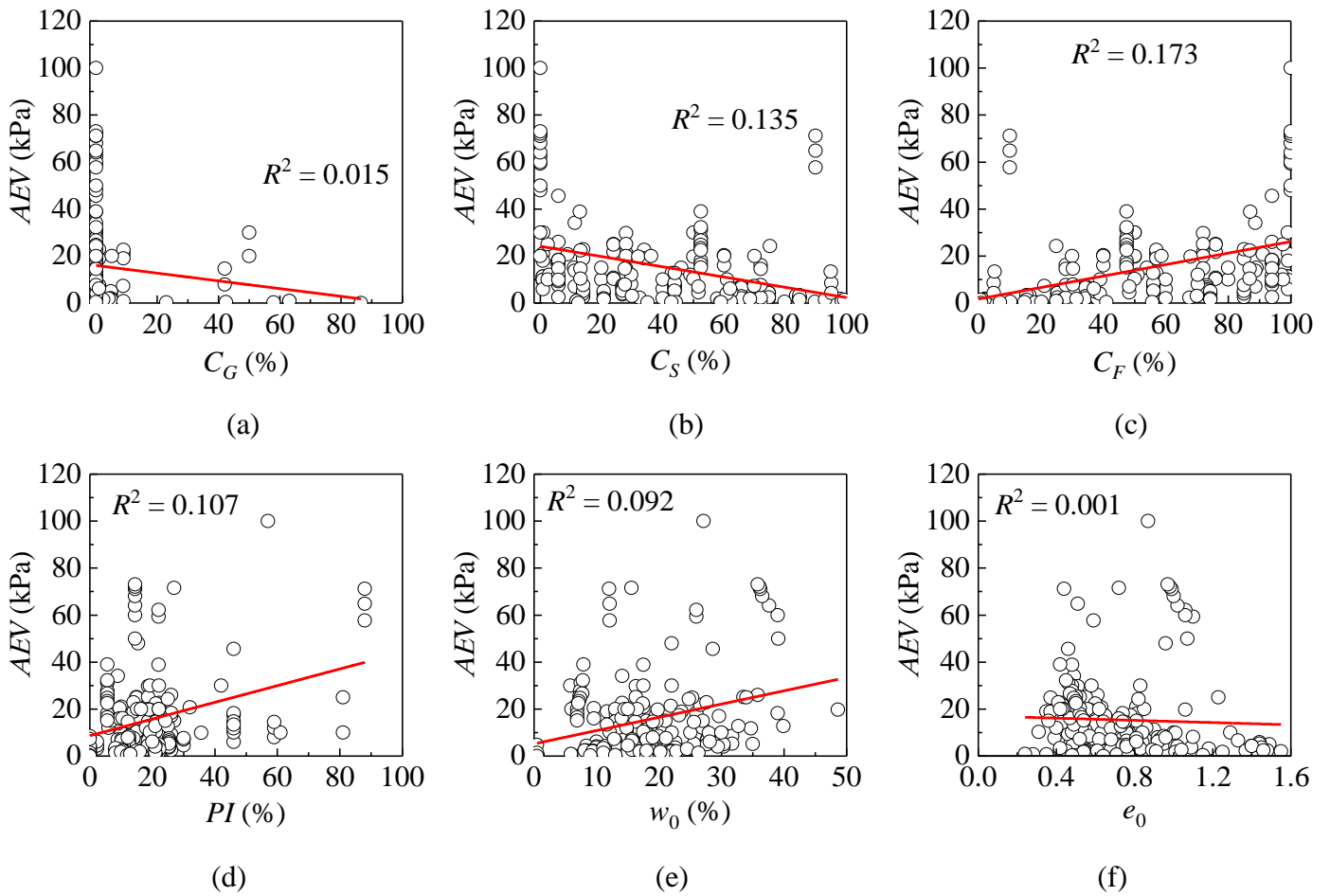


Fig. 4. Basic linear fittings between air-entry value and each input soil property: (a) gravel content; (b) sand content; (c) fines content; (d) plasticity index; (e) initial water content; (f) initial void ratio

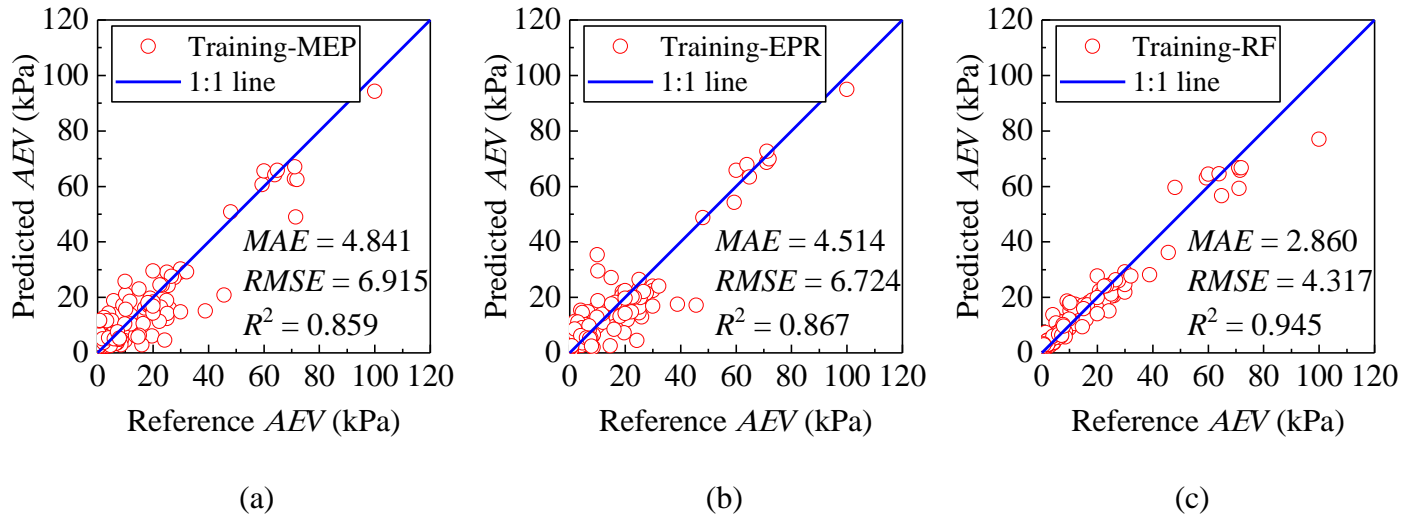


Fig. 5. Comparison between predicted and reference air-entry values for the training data: (a) MEP; (b) EPR; (c) RF

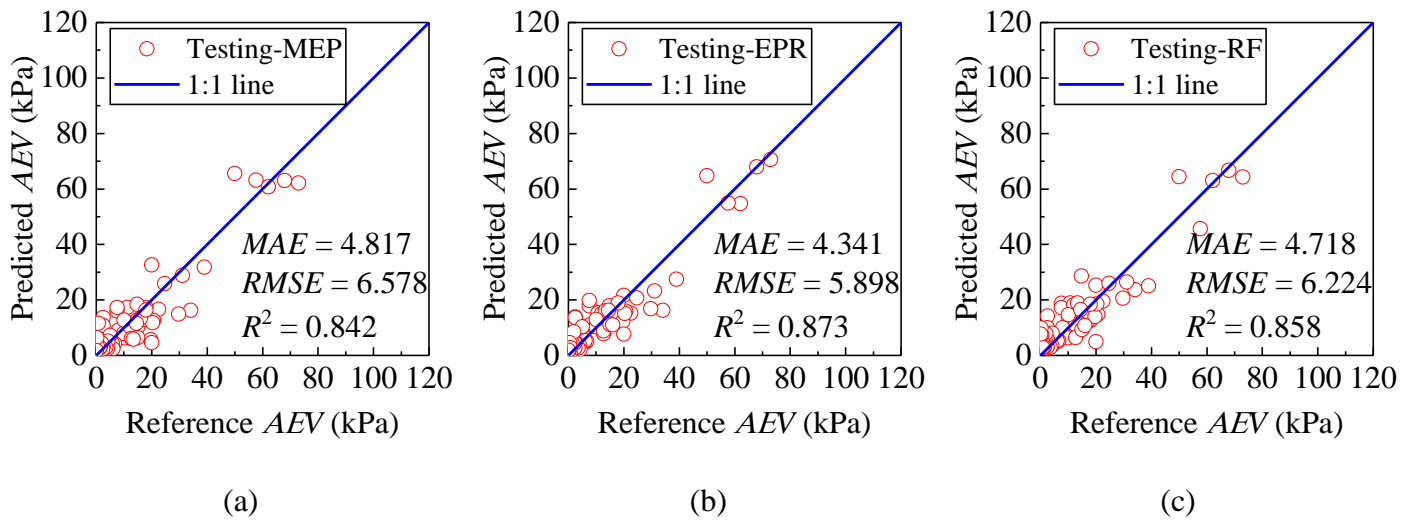
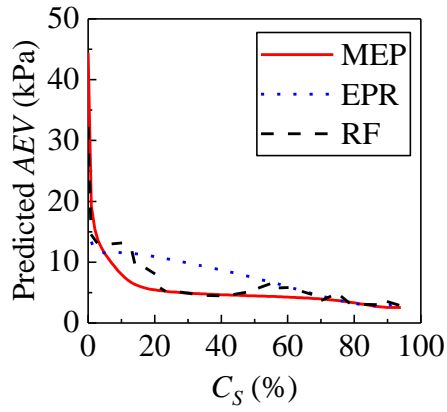
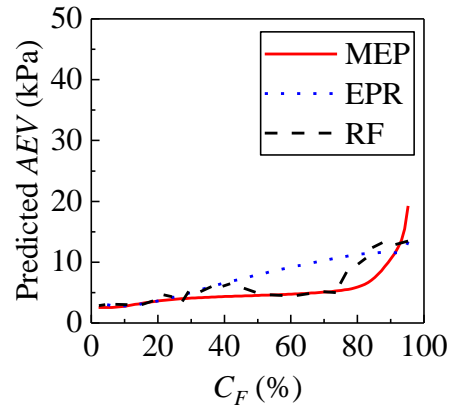


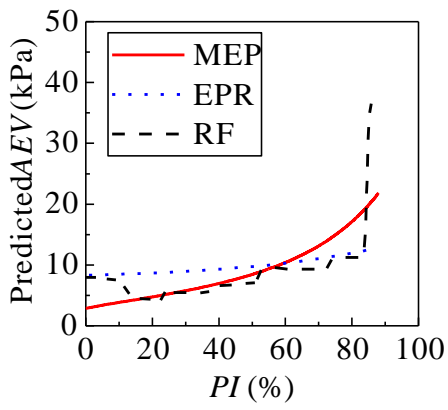
Fig. 6. Comparison between predicted and reference air-entry values for the testing data: (a) MEP; (b) EPR; (c) RF



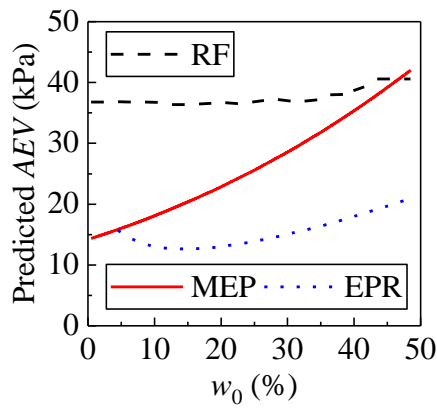
(a)



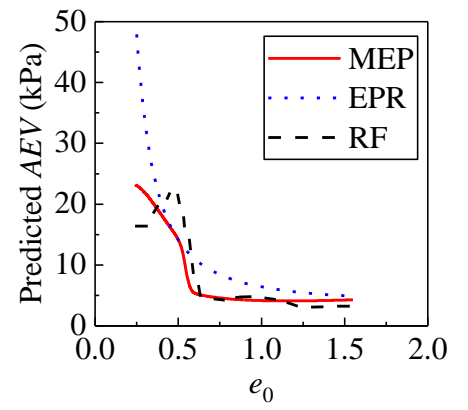
(b)



(c)



(d)



(e)

Fig. 7. Monotonicity analysis of the predicted air-entry value versus (a) sand content; (b) fines content; (c) plasticity index; (d) initial water content; (e) initial void ratio

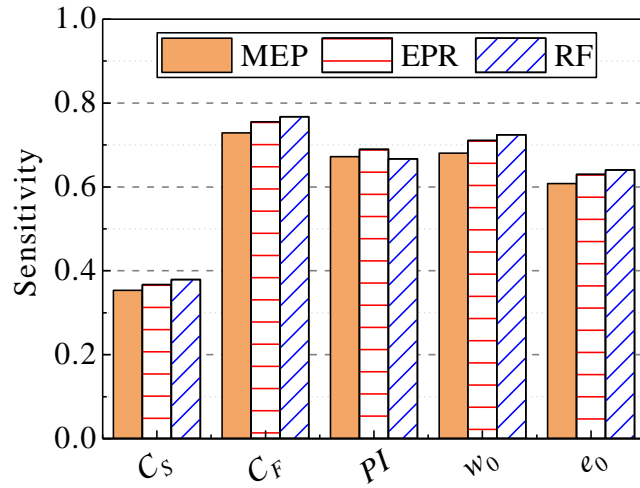


Fig. 8. Sensitivity analysis about the relevance of the input variables on the predicted air-entry value

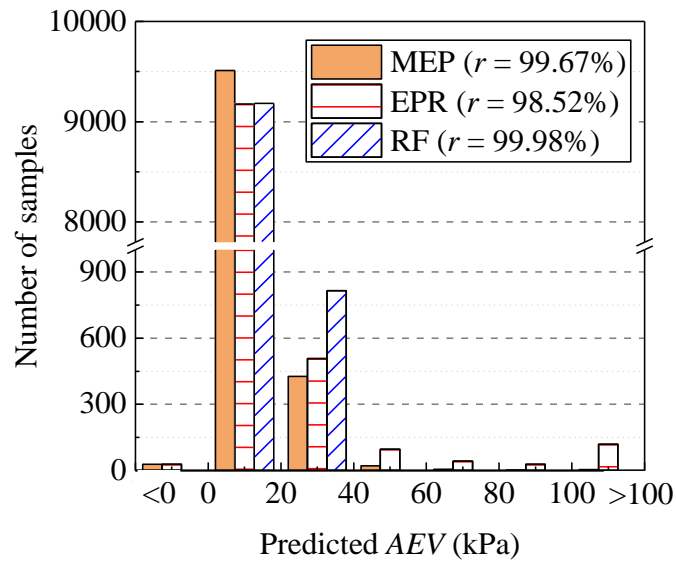


Fig. 9. Distribution of the predicted air-entry value in robustness analysis