# A Novel Hybrid Surrogate Intelligent Model for Creep Index Prediction based on Particle Swarm Optimization and Random Forest

Pin ZHANG[1, 2], Zhen-Yu YIN[2*], Yin-Fu JIN[2], Tommy H.T. Chan[1]

1 School of Civil Engineering & Built Environment, Science and Engineering Faculty, Queensland University of Technology (QUT), Brisbane, Qld 4001, Australia

2 Department of Civil and Environmental Engineering, Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong, China

* Corresponding author: Dr Zhen-Yu YIN, Tel: +852 3400 8470; Fax: +852 2334 6389; E-mail: zhenyu.yin@polyu.edu.hk; zhenyu.yin@gmail.com

**Abstract:** Long-term settlement issues in engineering practice are controlled by creep index $C_\alpha$. But current empirical models of $C_\alpha$ are not enough reliable. Different from previous correlations, this study proposes a hybrid surrogate intelligent model for predicting $C_\alpha$. The new combined model integrates the meta-heuristic particle optimization swarm (PSO) in the random forest (RF) to overcome the problem of user experience-dependent and local optimum. A total number of 151 datasets with four parameters (liquid limit $w_L$, plasticity index $I_p$, void ratio $e$ and clay content $CI$) and one output variable $C_\alpha$ are collected from published works. 11 combinations of these four parameters (one combination with four parameters, four combinations with three parameters and six combinations with two parameters) are set as input variables in the RF algorithm for determining the optimum combination of variables. In this novel model, PSO is employed to determine the optimum hyper-parameters in RF algorithm, and the fitness function in the PSO is defined as the mean prediction error for ten cross-validation sets for enhancing the robustness of RF models. The performance of RF model is particularly compared with the existing empirical formulae. The results indicate that the combinations of $I_P$–$e$, $CI$–$I_P$–$e$, and $CI$–$w_L$–$I_P$–$e$ are optimum RF models in respective group, and these models are recommended to predict $C_\alpha$ in engineering practice. Meanwhile, these three proposed models obviously outperform the empirical methods with lower prediction error. Parametric investigation indicates that the relationships between the $C_\alpha$ and four input variables in the proposed RF models harmonize with the physical explanation. Gini index generated in RF process indicates $C_\alpha$ is much more sensitive to $e$ than the remaining three input variables, followed by $CI$, $I_p$ and $w_L$, but the difference among later three variables can be negligible.

**Keywords:** Creep; Soft clay; Machine learning; Optimization; Physical properties; Correlation

## 1.  Introduction

Natural soft clays exhibit significant creep under both laboratory and in situ conditions after primary consolidation, which significantly influences the long-term stability of slope and safety of infrastructures in various fields, such as tunneling (Meng et al. 2018; Shen et al. 2014), slope (Jin et al. 2003; Yang et al. 2019) and embankment (Karstunen and Yin 2010; Yao et al. 2018; Yin et al. 2015; Zhu et al. 2019), etc. Engineers have to calculate the long-term settlement before construction in order to control the post-construction settlement to a tolerated value. Actually, time-dependent behavior of soft clays has been studied for a long history, various methods in standard and advanced elastic viscoplastic (EVP) models have also been proposed to estimate the creep settlement (Tan et al. 2018a; Yin and Graha 1989; Yin et al. 2010a; Yin et al. 2010b; Zhou et al. 2018), for which the measurement of viscous parameters takes time, and thus it is not convenient for engineers and researchers.

Calculation of long-term settlement using methods recommended in standard or EVP models indicates the corresponding parameters have to be determined in advance. Creep index $C_\alpha = \Delta e / \Delta \log(t)$ generally determined by the one-dimensional oedometer test is a key parameter to calculate long-term settlement in engineering practice, and it is also applied in most standard and EVP models (Yin et al. 2014a; Yin et al. 2011). Although most clays in engineering practice are intact rather than reconstituted, the $C_\alpha$ of intact clays is not an intrinsic property because bonds in these natural soils are progressively destroyed under various loading or unloading formations, which causes the apparent $C_\alpha$ highly nonlinear (Karstunen and Yin 2010; Yin et al. 2017). The varying value of $C_\alpha$ depending on stress level for intact clay is thus not suitable for use in actual engineering problems, which may also result in wrong predictions. However, the $C_\alpha$ of reconstituted clay without the interruption of soil structures is an intrinsic property, which provides the base for understanding the creep behavior of soils and it is thus more adaptable in engineering design (Jin et al. 2019). Because of these factors, this research merely focuses on the $C_\alpha$ of reconstituted clays.

Currently, the value of $C_\alpha$ is primarily determined by the curve-fitting technique based on the

57  experimental data. Researchers demonstrated that the creep property is affected by the microstructure of

58  soft clays (Yin and Chang 2009; Yin et al. 2014b), and physical properties somehow represent the

59  microstructure of clay (Jin et al. 2019). Nakase et al. (1998) proposed a linear formula to describe the

60  relationship between $C_\alpha$ and plasticity index $I_p$, similar relationship was also formulated by (Yin 1999).

61  Zeng et al. (2012) pointed out that the void ratio $e$ and the void ratio at liquid limit $e_L$ are the significant

62  factors for the creep behavior of soft clays, and a $C_\alpha$ prediction model was thus proposed based on these

63  two factors. Yin et al. (2015) formulated a linear expression of $C_\alpha$ with $e$ in a double logarithm plane.

64  More recently, Zhu et al. (2016) further developed a formulation of $C_\alpha$ considering both soil density and

65  soil structure. Nevertheless, these empirical formulae are merely capable of describing few soft clays.

66  Meanwhile, influential factors taken into consideration are limited, one or two in most formulae, although

67  the value of $C_\alpha$ depends on more influential factors. Therefore, $C_\alpha$ calculated by these empirical methods

68  is not enough reliable, and a model with wider adaptability between $C_\alpha$ and physical properties of soft

69  clays needs to be determined.

70      Machine learning (ML) algorithms are characterized by the strong capability of capturing the non-

71  linear relationships among high-dimension variables and worth to try. Various ML algorithms such as

72  back-propagation neural network (BPNN) (Basheer 2000; Habibagahi and Bamdad 2003; He and Li

73  2009; Penumadu and Zhao 1999; Rashidian and Hassanlourad 2014; Turk et al. 2001), evolutionary

74  neural network (ENN) (Johari et al. 2011), recurrent neural network (RNN) (Romo et al. 2001; Zhu et al.

75  1998), support vector machines (SVMs) (Kohestani and Hassanlourad 2016), evolutionary polynomial

76  regression (EPR) (Faramarzi et al. 2012; Javadi and Rezania 2009; Nassr et al. 2018), genetic

77  programming (GP) (Cabalar and Cevik 2011) and Bayesian-related methods (Gamse et al. 2018; Qi and

78  Zhou 2017; Tan et al. 2016; Tan et al. 2018b; Zhou et al. 2012), have been extensively utilized in

79  geotechnical engineering, e.g. the prediction of tunneling-induced settlement (Chen et al. 2019a; Chen et

80  al. 2019b; Hasanipanah et al. 2016), slope displacement and stability (Qi and Tang 2018a; Xu and Niu

81  2018; Yang et al. 2019), pile behaviors (Pooya Nejad and Jaksa 2017), soil physical and mechanical

82  characteristics (Feng et al. 2019; Kirts et al. 2018; Pham et al. 2018; Zhou et al. 2016b), etc. Recently, an

83 advanced ensemble algorithm random forest (RF) has been applied in geotechnical engineering practice.

84 The overfitting issue can be avoided and the importance of variables can be determined internally in the

85 RF algorithm. The superiority of RF algorithm has been proved in other comprehensive studies such as

86 prediction of slope stability (Qi and Tang 2018b), rockburst (Zhou et al. 2016a) and soil temperature

87 (Feng et al. 2019). Furthermore, (Chen et al. 2019b) conducted a comprehensive comparison of different

88 ML algorithms in predicting tunneling-induced settlement, and concluded that RF algorithm obviously

89 outperforms other ML algorithms. However, there is no research to develop $C_\alpha$ prediction model based on

90 ML algorithm. Meanwhile, the hyper-parameters such as the number of hidden layers and neurons in

91 numerous ML prediction models are generally determined by trial and error method and the deterministic

92 algorithms tend to optimize the general parameters (e.g., the weights and bias), which is time-consuming

93 and tends to fall in local optimum resulting in a poor performance of the obtained model.

94 To resolve these deficiencies, this paper proposes a novel RF intelligent model with integrating the

95 particle swarm optimization (PSO) algorithm to predict $C_\alpha$ of reconstituted clays. A database consisting

96 of four input physical properties of remoulded clays and the corresponding $C_\alpha$ is formed first. Thereafter,

97 11 combinations of these four parameters (one combination with four parameters, four combinations with

98 three parameters and six combinations with two parameters) are set as input variables for establishing RF

99 models, thereby the optimum combination with lowest errors for predicting $C_\alpha$ can be determined. Herein,

100 the hyper-parameters in each RF model are identified by PSO. To enhance the robustness of the proposed

101 model, the average prediction error of 10-fold cross-validation sets is set as the fitness function in PSO

102 algorithm. Finally, the relationships between $C_\alpha$ and input variables, and the sensitivity of input variables

103 in the proposed model are investigated comprehensively.

104 **2. Methodology**

105 *2.1 Random forest*

106 Random forest (RF) is an ensemble algorithm consisting of a collection of tree-structured classifiers

107 (Breiman 2001). Bagging (Breiman 1996) and random feature selection (Ho 1998) are integrated in the

108  RF algorithm. Each new bootstrap training set $\mathbf{N}_k$ with replacement samples is from the original training

109  set $\mathbf{N}$. Then a tree that is a sub-predictor is built based on the new bootstrap training set $\mathbf{N}_k$. For each $y$, $\mathbf{x}$

110  in the training set, aggregate the votes only over those sub-predictors for which $y$, $\mathbf{x}$ does not exist in the

111  $\mathbf{N}_k$. These sub-predictors are termed as out-of-bag (OOB) predictors and these datasets without the $\mathbf{N}_k$ are

112  termed as OOB datasets (accounting for one-third of the original training datasets). OOB datasets are

113  employed to evaluate the performance of sub-predictor developed upon the new bootstrap training set $\mathbf{N}_k$.

114  The hyper-parameters in this algorithm are the number of trees and random features at each node, as

115  shown in Table 1 with their ranges. Furthermore, random forest algorithm has been proved not to overfit

116  with the increase in the number of trees (Breiman 2001). They are more robust with respect to noise, and

117  the importance of input variables can be evaluated internally based on the value of Gini index (Breiman

118  2001). The detailed description of Gini index can refer to (Lovatti et al. 2019). Figure 1 presents the

119  process of building a random forest, which is also showing as following:

120  1. Draw $n$ bootstrap training sets from the original training data. Each bootstrap training set has about 2/3

121     of the original training datasets. The features in each bootstrap training set are selected at random.

122  2. Generate a decision tree based on each bootstrap training set. OOB datasets are used to evaluate the

123     performance of a decision tree which ultimately selects the best features/split among the training set.

124     All decision trees form a random forest.

125  3. Predict the new dataset by averaging the predictions of $n$ decision trees. (i.e., average of regression).

126     The output of the RF can be expressed as:

127
$$y = \frac{1}{n}\sum_{i=1}^{n} y_i(\mathbf{x}) \tag{1}$$

128  where, $y_i(\mathbf{x})$ = individual prediction of a tree for an input $\mathbf{x}$; $n$ = a total number of decision trees.

129  *2.2 Particle swarm optimization*

130     Particle swarm optimization (PSO) is a computational method that optimizes a problem by

131  iteratively improving candidate solutions, here termed as particles (Kennedy and Eberhart 1995). Each

132  particle has its position vector, $X_i^k$, velocity vector, $V_i^k$, and fitness value, where $k$ is the current

133 generation and $i$ is the $i$th particle. In the search-space, these particles move toward the global best

134 positions based on local best position and velocity. Herein, lower fitness value donates better position.

135 The global best positions, that is, the optimum hyper-parameters of RF algorithm (see Table 1) can be

136 determined when the fitness reaches the minimum value and keeps a constant. The velocity and position

137 of each particle are updated using the following equations:

$$V_i^{k+1} = \omega V_i^k + c_1 r_1 \left( P_i^k - X_i^k \right) + c_2 r_2 \left( P_g^k - X_i^k \right) \tag{2}$$

$$X_i^{k+1} = X_i^k + V_i^{k+1} \tag{3}$$

140 where, $c_1$ and $c_2$ = acceleration coefficients; $\omega$ = a weight which is called "inertia weight", equal to 1

141 which is a typical value; $r_1$ and $r_2$ = random numbers, distributing in the interval [0, 1]; $P_i$ = current best

142 location of $i$th particles; $P_g$ = the global best among all particles. In this study, both of $c_1$ and $c_2$ are equal

143 to 1.49, which is the typical value used in PSO. Upper and lower bound of $Vk I$ are 1 and −1, respectively.

144 In the PSO, only three parameters (i.e., $\omega$, $c_1$ and $c_2$) need to be set. These values mainly affect the

145 convergence speed (El-Gallad et al. 2002; Zheng et al. 2003) and slightly affect the final optimization

146 results as long as a large number of generations is performed. Therefore, the maximum generation is 100,

147 which is large enough to find the best solution. To enhance the PSO performance through selecting proper

148 parameters, many experiences on the selection of parameters of PSO can be found in previous studies

149 (Shi and Eberhart 1998; Trelea 2003).

150 *2.3 Evaluation indicators*

151    Three common indicators root mean square error (RMSE), mean absolute error (MAE) <u>and</u> mean

152 absolute percentage error (MAPE) are employed to evaluate the performance of RF models. RMSE and

153 MAE can directly reflect the prediction error, but they are the scale-dependent indicators, whose scale is

154 related to the scale of the data. In other words, low values of RMSE and MAE cannot indicate the great

155 performance if the values of model output variables are small. Meanwhile, large values of RMSE and

156 MAE cannot indicate the bad performance if the values of model output variables are large. RMSE has

157 the same scale as the data, but it is more sensitive to outliers than MAE. MAPE is a scale-independent

158 indicator, thereby it is not affected by the scale of data. Nevertheless, the MAPE value is infinite or

159 undefined if the observed values are close to zero. A comprehensive comparison of these indicators can

160 refer to (Hyndman and Koehler 2006). The combination of RMSE, MAE and MAPE can effectively

161 evaluates model performance, and these three indicators are thus used in this study. The expression of

162 these three indicators can be obtained by

163
$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(r_i - p_i)^2}$$
(4)

164
$$MAE = \frac{1}{n}\sum_{i=1}^{n}|r_i - p_i|$$
(5)

165
$$MAPE = \frac{1}{n}\sum_{i=1}^{n}\left|\frac{r_i - p_i}{r_i}\right| \times 100\%$$
(6)

166 where, $r$ = measured output value; $p$ = predicted output value; $n$ = a total number of datasets. Low values

167 of these three indicators indicate a model with great performance.

168 *2.4 K-fold cross validation*

169 The whole process of establishing a ML model includes three phases: training, validation and

170 testing. *K*-fold cross-validation (CV) method has been extensively used to validate model (Stone 1974)

171 for improving the robustness of ML models and avoiding overfitting problem. In this method, the original

172 training set is randomly divided into $k$ sub-datasets. Herein, $k$-1 sub-datasets are used to train model and a

173 remaining dataset is used to validate model. Each sample thus has opportunity to train and test model. In

174 open literatures, the $k$ was recommended to be set as 10 (Kohavi 1995), therefore, 10-fold CV is applied

175 in this study.

176 At each round, RF model with fixed hyper-parameters will be trained ten times by random nine sub-

177 datasets, and the remaining one sub-dataset will be used to validate model. The performance of the RF

178 model with fixed hyper-parameters will be evaluated by the mean prediction error for ten validation sets,

179 that is, the fitness function in the PSO algorithm.

180
$$Fitness = \frac{1}{10}\sum_{i=1}^{10} \text{MAE}_i$$
(7)

181 where, $\text{MAE}_i$ = prediction error for $i$th validation set.

182     Note that the use of 10-fold cross validation can also reduce the effect of the selection of different

183     amounts of data on the model performance.

184     *2.5 Grey relational grade*

185     Grey relational grade (GRG) has been extensively employed to evaluate uncertain correlations

186     among variables (Jiang and He 2012; Li and Chen 2019). The geometric similarity of the time series of

187     the two variables is taken into consideration in this method. Given a reference sequence $x_r = x_r (x_r(1),$

188     $x_r(2), …, x_r(n))$ and a compared sequence $x_i = x_i (x_i(1), x_i(2), …, x_i(n))$. The grey relational coefficient

189     between two sequences at $j$th ($j = 1, 2, …, n$) criterion is defined as following

190
$$\gamma\left(x_r\left(j\right),x_i\left(j\right)\right)=\frac{\min\limits_{i}\min\limits_{j}\left|x_r\left(j\right)-x_i\left(j\right)\right|+\delta\max\limits_{i}\max\limits_{j}\left|x_r\left(j\right)-x_i\left(j\right)\right|}{\left|x_r\left(j\right)-x_i\left(j\right)\right|+\delta\max\limits_{i}\max\limits_{j}\left|x_r\left(j\right)-x_i\left(j\right)\right|}$$
(8)

191     where, $\delta$ = resolving coefficient, in the range of [0, 1], which is usually considered as 0.5. Thereafter, the

192     GRG between sequences $x_r$ and $x_i$ can be obtained by

193
$$\gamma\left(x_r,x_i\right)=\frac{1}{n}\sum_{j=1}^{n}\gamma\left(x_r\left(j\right),x_i\left(j\right)\right)$$
(9)

194     where, large GRG value means the high correlations between sequences $x_r$ and $x_i$.

195     **3.  Proposed intelligent model**

196     *3.1 Model framework*

197     Figure 2 presents the process of establishing the proposed creep index $C_\alpha$ prediction model. The

198     whole process can be categorized into three phases: data preprocessing, training and testing RF prediction

199     models. At the first phase, main influential factors of $C_\alpha$ need to be determined and collected, a database

200     consisting of input and output parameters are then established. Herein, 80% of data are randomly selected

201     for training the model while the remaining are used to test the model. The selection of input variables are

202     vitally important to the model performance. According to previous investigations, the liquid limit $w_L$,

203     plasticity index $I_p$ and void ratio $e$ have been used to form an empirical equation to predict $C_\alpha$ (Nakase et

204     al. 1998; Zeng et al. 2012; Zhu et al. 2016). In addition, the clay content (*CI*) also has a significant

205 influence on predicting $C_\alpha$ (Jin et al. 2019). Therefore, these four parameters $w_L$, $I_p$, $e$ and $CI$ are

206 preferably considered as the model input variables in this study. The correlation of selected parameters to

207 $C_\alpha$ is examined by GRG method, as shown in Fig 4.

208 Feature selection method has been successfully conducted to process high-dimensional data and

209 select the most relevant factors to the outputs of model (Gao et al. 2018; Lu et al. 2018). In order to

210 determine the optimum combination of input variables for predicting $C_\alpha$, a total number of 11

211 combinations of input variables are used to train the prediction model respectively, which can be divided

212 into three groups: 6 combinations of two variables, 4 combinations of three variables, and 1 combination

213 of four variables. The one variable as input parameter is not taken into consideration in this research,

214 because the prediction model trained with only one variable suffers from underfitting, losing

215 generalization capability. Therefore, a total number of 11 $C_\alpha$ prediction models with different

216 combinations of input variables need to be trained. The model with the optimum performance will be

217 recommended to predict $C_\alpha$ in practice engineering.

218 The objective at the process of training model is to identify the optimum hyper-parameters in 11

219 prediction models. PSO algorithm is employed to search for the optimum hyper-parameters in the RF

220 algorithm. At each round, the proposed model starts from randomly assigning hyper-parameters to RF.

221 Training set is then randomly divided into ten parts using 10-fold CV method. If the termination criterion

222 is satisfied, that is, whether or not reaches maximum 100 generations and fitness value converges, the

223 optimum hyper-parameters in one model can be determined. Otherwise, RF will be assigned a new set of

224 hyper-parameters by the PSO algorithm. In this way, the hyper-parameters in 11 prediction models can be

225 determined ultimately.

226 At the last phase, the 11 prediction models with optimum hyper-parameters will be evaluated by the

227 test set. The model with the lowest error will be selected as the optimum $C_\alpha$ prediction model.

228 Meanwhile, the optimum model in each group can also be determined. Therefore, engineers and

229 researchers can select the most appropriate prediction model based on their existing experiment data.

*3.2 Data source*

231    The data used in this research were collected from published research works and consist of various

232    remoulded clays in the world (Li et al. 2012; Yin 1999; Yin et al. 2015; Zeng et al. 2012; Zhu et al.

233    2016). A total number of 151 datasets were ultimately collected. The various remoulded clays in this

234    database facilitate the development of a uniform $C_\alpha$ prediction model for all remoulded clays. The

235    histogram of all variables in the database is presented in the diagonal line of Fig. 3. Scatter plots of

236    pairwise variables are also plotted in this figure. It can be observed that all variables cover a wide range

237    of values, which sufficiently extends the applicability of the proposed model.

238    To eliminate the effect of different magnitudes of input variables on the model's performance and

239    also to reduce the computational cost, all datasets have been normalized into the range of (−1, 1) using the

240    following expression:

241
$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}} \left( \overline{x}_{max} - \overline{x}_{min} \right) + \overline{x}_{min} \tag{10}$$

242    where $x$ = actual value of input variables, $x_{min}$ = minimum value of input variables and $x_{max}$ = maximum

243    value of input variables. $\overline{x}_{min} = -1$; $\overline{x}_{max} = 1$.

244    Figure 4 presents the GRG values among input and output variables. The number in each panel

245    represents the GRG value between pairwise variables. It can be observed that the largest GRG value

246    reaches 0.81, showing high correlation between input $e$ and output $C_\alpha$, whereas the GRG values of

247    remaining three input variables are roughly identical. Overall, GRG values of four input parameters

248    exceed 0.65, suggesting the selected influential factors are appropriate to predict $C_\alpha$.

249    **4.  Results**

250    *4.1 Determination of hyper-parameters*

251    PSO algorithm is utilized for tuning two hyper-parameters in the RF model. Figure 5 shows the

252    evolution of the fitness value within 100 generations in 11 $C_\alpha$ prediction models. In regard to the CV sets,

253    the evolution of fitness is different and ultimately converges at various values among 11 models. RF

254      model with the combination of $I_p$–$e$ outperforms remaining five models trained by two input variables.

255      The ultimate fitness values is only 0.00408 and it actually maintains steadily from the initial generation.

256      The fitness value of this model is even less than fitness values produced by all RF models with three input

257      variables, whereas the fitness values of remaining five RF models are larger than that. RF model with the

258      combination of $w_L$–$I_p$–$e$ yields the lowest fitness value of 0.0043, compared with remaining three models

259      trained by three input variables. RF model with the combination of four variables shows the best

260      performance with the fitness value of 0.00406. It is noteworthy that no further decrease in the fitness

261      value is observed after the generation exceeds 16 in 11 RF prediction models, which indicates that 100

262      generations are large enough to determine the optimum hyper-parameters in RF models. Meanwhile, the

263      fitness values of RF models with two input variables range from 0.00408 to 0.00591, whereas this value

264      ranges from 0.0043 to 0.00504 in RF models with three input variables. It actually indicates that the

265      performance of $C_\alpha$ prediction models is increasingly steady with the increase in the number of input

266      variables. Overall, the fitness value in RF prediction model with two input variables is largest, and the

267      lowest value appears in the RF prediction model with four input variables.

268      *4.2 Prediction of $C_\alpha$ for the validation sets*

269      In order to reveal the reason behind the difference in the converged fitness value in 11 RF prediction

270      models, Figure 6 presents the values of three indicators in each CV set. It can be observed that the biggest

271      difference of RMSE and MAE values appears in the second CV set, where the RF models with four

272      combinations, that are, $w_L$–$e$, $I_p$–$e$, $w_L$–$I_p$–$e$, and $CI$–$w_L$–$I_p$–$e$, produces much less RMSE and MAE values

273      than the remaining seven RF models. Meanwhile, the variation of RMSE and MAE values in these four

274      models is slight, compared with the remaining RF models. Therefore, the fitness values produced by these

275      four RF models are lowest in each group. Aside from the second CV set, the evolution of RMSE and

276      MAE values is roughly identical in each RF model. It can be observed from Fig. 6(c) that the evolution of

277      MAPE differs from RMSE and MAE due to its scale-independent characteristic. The performance of RF

278      models in each CV set presents obvious difference, especially in the second, seventh and ninth CV sets,

279      but the four RF models as mentioned above still show the lower MAPE value at the second CV set.

280    The distribution of RMSE, MAE and MAPE values in each RF model is presented by boxplot, as

281    shown in Fig. 7. It is clear that the distribution of RMSE and MAE are roughly identical in each RF

282    model. In each group, the RF models with better performance produce lower mean RMSE and MAE

283    values. Meanwhile, the ranges of RMSE and MAE are much smaller. From the perspective of MAPE, RF

284    model with the combination of $I_p$–$e$ shows the best performance among all prediction models with the

285    minimum mean prediction error of 20.8%. RF model with the combination of $CI$–$w_L$–$e$ outperforms the

286    remaining RF models trained by three parameters. These two characteristics are different from RMSE and

287    MAE results. It will affect the performance of RF models for the test set, which will be revealed at the

288    next section. Overall, for the CV sets, RF model with the combination of two parameters produces

289    maximum prediction error as well as a wider range of errors.

290    *4.3 Prediction of $C_\alpha$ for the test set*

291    On the basis of hyper-parameters determined in the former section, 11 optimum RF models can be

292    established. Figure 8 presents the scatter plot of predicted $C_\alpha$ for the training and test sets using optimum

293    models. The predicted results of RF models with two input variables are illustrated in Figs. 8(a)-(f). The

294    combinations of $CI$–$e$, $w_L$–$e$, and $I_p$–$e$ clearly outperform the remaining three RF models. The predict $C_\alpha$

295    for the training set exist perfect agreement with the measured $C_\alpha$, and the predict $C_\alpha$ for the test set are

296    also close to the P = M line. In contrast, predicted results using the remaining three RF models widely

297    scatter, and the disagreement of the predicted and measured $C_\alpha$ is frequently observed. Figures 8(g)-(j)

298    present the predicted results using RF models with three input variables. The range of predicted $C_\alpha$ using

299    RF model with the combination of $CI$–$w_L$–$I_p$ is much smaller than the measured $C_\alpha$, losing fidelity at most

300    points. The predicted results using the remaining three models show great agreement with the measured

301    results. The predicted $C_\alpha$ for both training and test sets using the RF model with four input variables are

302    closer to the line with slope of 1 than remaining 10 models (see Fig. 8(k)).

303    Table 2 summarizes the values of indicators in 11 optimum RF models. The values of indicator are

304    roughly consistent with the model performance presented in Fig. 8, that is, lower values of indicators

305    harmonize with more reasonable distribution of the predicted $C_\alpha$. Meanwhile, models with great

306    performance in the CV sets also exhibit better performance in the test set. For instance, RF models with

307    combination of $I_P$–$e$ outperform the remaining models with two input variables in both CV and test sets,

308    and RF model with combination of $CI$–$w_L$–$I_p$–$e$ always presents the best performance among all models

309    throughout the analysis. However, a special case is observed in the RF models with three input variables,

310    where values of indicators of RF model with combination of $CI$–$w_L$–$I_p$ are lowest for the test set, but Fig.

311    8(g) presents the distribution of predicted $C_\alpha$ using this model is not acceptable. In reality, the indicators

312    values for the training set are larger than the remaining three RF models with three input variables. Figure

313    8 (g) also presents that the distribution of predicted $C_\alpha$ using this model is much more concentrated, the

314    slight variation of predicted $C_\alpha$ thus leads to the lower values of indicators. Similar conditions also appear

315    in the combinations of $CI$ –$w_L$, $CI$–$I_P$, and $w_L$–$I_P$. However, the smaller range of predicted results indicates

316    that the prediction applicability of these models is limited, that is, weak generalization capability. To sum

317    up, the combinations of $I_P$–$e$, $CI$–$I_P$–$e$, and $CI$–$w_L$–$I_p$–$e$ are optimum RF model in respective group, and

318    these models are recommended to predict $C_\alpha$ in engineering practice.

319    **5.  Discussions**

320    *5.1 Comparison with empirical formula*

321    In order to evaluate the predictive ability of the proposed model, five commonly used empirical

322    methods are used for comparison (see Table 3). Several input parameters of the proposed model are also

323    used in these empirical methods. Figure 9 presents scatterplot of the predicted $C_\alpha$ using five empirical

324    methods. The predicted $C_\alpha$ values in the Figs. 9(a) and (b) are roughly identical, because $I_p$ is the only

325    parameter in these two methods. It leads to predicted $C_\alpha$ with the identical value and the small range. In

326    Fig. 9(c), the values of predicted $C_\alpha$ vary dramatically, losing fidelity at most points. The predicted $C_\alpha$

327    using the empirical method proposed by (Zhu et al. 2016) exhibit great agreement with measured $C_\alpha$ (see

328    Figs. 9(d) and (e)).

329    Figure 10 presents the comparison between the empirical and proposed models in predicting $C_\alpha$,

330    where the results of three optimum models in each group are plotted. It is clear that three proposed

331  models obviously outperform the empirical methods with lower values of RMSE, MAE and MAPE. The

332  mean RMSE and MAE values produced by the proposed models decrease by 0.0012 and 0.001,

333  respectively, compared with the empirical methods. Further, the mean MAPE decreases from 29.33% to

334  18.09%.

335  *5.2 Parametric investigation*

336  A robust ML model exhibits smooth functions with respect to the input and output variables and

337  exhibits physical explanation for the behaviors (Shahin et al. 2005). Therefore, the correlations between

338  four input variables and the $C_\alpha$ at three typical points in the test set are investigated. Note that the $C_\alpha$

339  presented in this section are predicted by RF model with four input variables for comprehensively

340  investigating the effects of all variables on the $C_\alpha$. At each round, values of three input variables are fixed

341  in the RF model, whereas the value of a remaining variable increases from $0.2v$ to $2v$ with an interval of

342  $0.2v$ ($v$ donates the original value of the investigated parameter).

343  Figure 11 presents the correlations between input and output variables in the RF model with four

344  input parameters. The fixed values of three variables are also plotted in figures. Note that a perfect

345  smooth correlation between input and output variables in RF model is hardly obtained, because RF model

346  is developed based on measured data, and a smooth correlation is not observed in the measured data as

347  shown in Fig. 3. Therefore, the correlations between input and output variables in the RF model merely

348  reflect a general trend. It can be seen from Fig. 11(a) that the predicted $C_\alpha$ decreases initially with the

349  increase in the *CI*, $C_\alpha$ then starts to increase after reaches minimum value. $C_\alpha$ ultimately holds steadily

350  with the continuously increase in the *CI*. The correlation between $w_L$ and $C_\alpha$ shows a similar condition.

351  An opposite trend is observed between $w_L$ and $C_\alpha$. With an increase in the $w_L$, $C_\alpha$ increases initially, then

352  starts to decreases after reaches maximum value, and $C_\alpha$ ultimately maintain a constant value. In regard to

353  $e$, $C_\alpha$ increases monotonically with the increase in the value of $e$. After reaching maximum value, $C_\alpha$ is

354  constant. The evolution of predicted $C_\alpha$ with the change in the four input variables is similar in three

355  points, but the magnitude of $C_\alpha$ is different. Note that the predicted $C_\alpha$ using RF model always holds

356  steadily when the input variables exceeds a certain range, which more or less violates laboratory test

357 results. This is a limitation for all data-drive model, since the prediction capability of this kind of model

358 will be useless when values of new input variables exceed the range of original database. Overall, the

359 correlations presented in Fig. 11 are consistent with physical explanation, indicating robustness and

360 reasonability of the proposed RF models.

361     To investigate the generalization ability of the proposed model, a database including a total number

362 of 10000 random samples is established. Each sample has four input variables ($CI$, $w_L$, $I_p$, $e$), and it

363 assumes that each variable complies with lognormal distributions (Cao and Wang 2014; Zhang et al.

364 2009; Zhang et al. 2018). Herein, the values of mean and standard deviation for each variable are

365 consistent with the results presented in Fig. 3. Thereafter, the $C_\alpha$ is predicted by the RF model with four

366 input variables.

367     Figure 12 shows the distribution of predicted $C_\alpha$. It can be observed that the $C_\alpha$ roughly meet the

368 lognormal distribution with the coefficient of determination of 0.91. The mean and standard deviation

369 values of predicted $C_\alpha$, 0.020 and 0.008, respectively, show great agreement with measured values, 0.019

370 and 0.011, respectively. Note that the range of predicted $C_\alpha$ are perfectly consistent with measured $C_\alpha$,

371 because the prediction ability of RF algorithm is useless when the datasets exceed the range of the

372 original database. Overall, the performance of proposed RF model is absolutely reliable for the unseen

373 datasets within the range of original database.

374 *5.3 Sensitivity of variables*

375     Variable importance measure (VIM) provides a basis for understanding the contributions of different

376 input variables to the model output (Hapfelmeier et al. 2012). As mentioned in the Random Forest

377 section, variable importance can be measured internally in RF algorithm, which is achieved by

378 investigating the influence of the variation of input variables on the Gini index. Input variable that causes

379 larger variation in Gini index is more significant to the model predictive capability (Breiman et al. 1984).

380     The mean decrease in Gini index caused by the change in each variable is shown in Fig. 13. It can be

381 observed that model output $C_\alpha$ is much more sensitive to $e$ than the remaining three variables, followed by

382 $CI$, $I_p$ and $w_L$, but the difference among these three variables can be negligible. It explains the reason that

383     $e$ is a common input variable in three optimum models as mentioned above, meanwhile $w_L$ as an input

384     variable merely appears in the RF model with four input variables.

385     **6. Conclusions**

386     A new hybrid surrogate intelligent model based on particle swarm optimization (PSO) and random

387     forest (RF) algorithms was proposed in this study for predicting $C_\alpha$. A database with four input variables

388     liquid limit $w_L$, plasticity index $I_p$, void ratio $e$, clay content $CI$ and one output variable $C_\alpha$ was first

389     established. 80% of data was used to train model, and the remaining 20% of data was used to test model.

390     High values of grey relation grade (>0.65) between four input variables and one output variable indicated

391     that the selected influential factors are appropriate to predict $C_\alpha$.

392     To search the optimum combination of input variables with respect to predicting $C_\alpha$, a total number

393     of 11 combinations of input variables were used to train prediction model respectively, which can be

394     divided into three groups: 6 combinations of two variables, 4 combinations of three variables, and 1

395     combination of four variables. Therefore, 11 RF models with different combinations of input variables

396     were established.

397     To determine the optimum hyper-parameters in the RF algorithm, meta-heuristic algorithm PSO was

398     integrated with RF algorithm. The fitness function in the PSO algorithm was defined as the mean

399     prediction error for ten cross-validation sets, enhancing the robustness of RF models. The predicted

400     results for the training and testing sets indicate that the combinations of $I_P$–$e$, $CI$–$I_P$–$e$, and $CI$–$w_L$–$I_p$–$e$ are

401     optimum RF models in respective group. Therefore, these models are recommended to predict $C_\alpha$ in

402     engineering practice.

403     Compared with the empirical methods of predicting $C_\alpha$, three proposed models obviously outperform

404     the empirical methods with lower values of RMSE, MAE and MAPE. The mean RMSE and MAE values

405     produced by the proposed models decrease by 0.0012 and 0.001, respectively, and the mean MAPE

406     decreases from 29.33% in empirical methods to 18.09% in the proposed models.

407     Parametric investigation indicates that the relationships between the $C_\alpha$ and four input variables in

408 the proposed RF models harmonize with the physical explanation, verifying the robustness and

409 reasonability of the proposed models. Gini index in the RF algorithm was employed to evaluate the

410 sensitivity of input variables. The results indicate that the model performance is much more sensitive to $e$

411 than other three variables, followed by $CI$, $I_p$ and $w_L$, but the difference among these three variables can

412 be negligible.

413     As mentioned above, the performance of the RF model depends significantly on the datasets.

414 Although the database used in this study includes numerous remoulded clays and the range of variables

415 are large, it should be further expanded in the future with more data for facilitating the application of

416 proposed models.

417     In order to allow readers to quickly perform the training and get results, the used datasets and the

418 MATLAB source code for the proposed hybrid RF and PSO on predicting the $C_\alpha$ are provided and can be

419 downloaded from the website of

420 https://www.researchgate.net/publication/334450481_Matlab_code_for_predicting_creep_index_using_h

421 ybrid_Random_Forest_and_Particle_Swarm_Optimization_algorithms.

426 **References**

427 Basheer, I.A. 2000. Selection of methodology for neural network modeling of constitutive hystereses
428     behavior of soils. Computer-Aided Civil and Infrastructure Engineering, **15**, 440–458.

429 Breiman, L. 1996. Bagging Predictors. Machine Learning, **24**, 123-140, doi: 10.1007/bf00058655.

430 Breiman, L. 2001. Random Forests. Machine Learning, **45**, 5–32.

431 Breiman, L.I., Friedman, J.H., Olshen, R.A. & Stone, C. 1984. Classification and Regression Trees
432     (CART). Encyclopedia of Ecology, **57**, 582-588.

433 Cabalar, A.F. & Cevik, A. 2011. Triaxial behavior of sand–mica mixtures using genetic programming.

434    Expert Systems with Applications, **38**, 10358-10367.

435  Cao, Z.J. & Wang, Y. 2014. Bayesian model comparison and selection of spatial correlation functions for

436    soil parameters. Structural Safety, **49**, 10-17, doi: 10.1016/j.strusafe.2013.06.003.

437  Chen, R.P., Zhang, P., Kang, X., Zhong, Z.Q., Liu, Y. & Wu, H.N. 2019a. Prediction of maximum

438    surface settlement caused by EPB shield tunneling with ANN methods. Soils and Foundations, in

439    press, doi: 10.1016/j.sandf.2018.11.005.

440  Chen, R.P., Zhang, P., Wu, H.N., Wang, Z.T. & Zhong, Z.Q. 2019b. Prediction of shield tunneling-

441    induced ground settlement using machine learning techniques. Frontiers of Structural and Civil

442    Engineering, in press.

443  El-Gallad, A., El-Hawary, M., Sallam, A. & Kalas, A. 2002. Enhancing the particle swarm optimizer via

444    proper parameters selection.    *IEEE CCECE2002. Canadian Conference on Electrical and*

445    *Computer Engineering. Conference Proceedings (Cat. No. 02CH37373)*. IEEE, 792-797.

446  Faramarzi, A., Javadi, A.A. & Alani, A.M. 2012. EPR-based material modelling of soils considering

447    volume changes. Computers & Geosciences, **48**, 73-85, doi: 10.1016/j.cageo.2012.05.015.

448  Feng, Y., Cui, N., Hao, W., Gao, L. & Gong, D. 2019. Estimation of soil temperature from

449    meteorological data using different machine learning models. Geoderma, **338**, 67-77, doi:

450    10.1016/j.geoderma.2018.11.044.

451  Gamse, S., Zhou, W.-H., Tan, F., Yuen, K.-V. & Oberguggenberger, M. 2018. Hydrostatic-season-time

452    model updating using Bayesian model class selection. Reliability Engineering & System Safety, **169**,

453    40-50.

454  Gao, W., Hu, L., Zhang, P. & He, J. 2018. Feature selection considering the composition of feature

455    relevancy. Pattern Recognition Letters, **112**, 70-74, doi: 10.1016/j.patrec.2018.06.005.

456  Habibagahi, G. & Bamdad, A. 2003. A neural network framework for mechanical behavior of unsaturated

457    soils. Canadian Geotechnical Journal, **40**, 684-693, doi: 10.1139/t03-004.

458  Hapfelmeier, A., Hothorn, T., Ulm, K. & Strobl, C. 2012. A new variable importance measure for random

459    forests with missing data. Statistics and Computing, **24**, 21-34, doi: 10.1007/s11222-012-9349-1.

460  Hasanipanah, M., Noorian-Bidgoli, M., Jahed Armaghani, D. & Khamesi, H. 2016. Feasibility of PSO-

461    ANN model for predicting surface settlement caused by tunneling. Engineering with Computers, **32**,

462    705-715, doi: 10.1007/s00366-016-0447-0.

463  He, S. & Li, J. 2009. Modeling nonlinear elastic behavior of reinforced soil using artificial neural

464    networks. Applied Soft Computing, **9**, 954-961, doi: 10.1016/j.asoc.2008.11.013.

465  Ho, T.K. 1998. The random subspace method for constructing decision forests. IEEE Transactions on

466    Pattern Analysis & Machine Intelligence, **20**, 832-844, doi: 10.1109/34.709601.

467  Hyndman, R.J. & Koehler, A.B. 2006. Another look at measures of forecast accuracy. International

468        Journal of Forecasting, **22**, 679-688, doi: 10.1016/j.ijforecast.2006.03.001.

469    Javadi, A.A. & Rezania, M. 2009. Applications of artificial intelligence and data mining techniques in
470        soil modeling. Geomechanics and Engineering, **1**, 53-74, doi: 10.12989/gae.2009.1.1.053.

471    Jiang, H. & He, W. 2012. Grey relational grade in local support vector regression for financial time series
472        prediction. Expert Systems with Applications, **39**, 2256-2262, doi: 10.1016/j.eswa.2011.07.100.

473    Jin, F., Zhang, C.H., Wang, G. & Wang, G.L. 2003. Creep modeling in excavation analysis of a high rock
474        slope. Journal of Geotechnical and Geoenvironmental Engineering, **129**, 849-857, doi:
475        10.1061/(ASCE)1090-0241(2003)129:9(849).

476    Jin, Y.F., Yin, Z.Y., Zhou, W.H., Yin, J.H. & Shao, J.F. 2019. A single-objective EPR based model for
477        creep index of soft clays considering L2 regularization. Engineering Geology, **248**, 242-255, doi:
478        10.1016/j.enggeo.2018.12.006.

479    Johari, A., Javadi, A.A. & Habibagahi, G. 2011. Modelling the mechanical behaviour of unsaturated soils
480        using a genetic algorithm-based neural network. Computers and Geotechnics, **38**, 2-13, doi:
481        10.1016/j.compgeo.2010.08.011.

482    Karstunen, M. & Yin, Z.Y. 2010. Modelling time-dependent behaviour of Murro test embankment.
483        Géotechnique, **60**, 735-749, doi: 10.1680/geot.8.P.027.

484    Kennedy, J. & Eberhart, R. 1995. Particle swarm optimization. *1995 IEEE International Conference on*
485        *Neural Networks*, Perth, Australia, 1942-1948.

486    Kirts, S., Panagopoulos, O.P., Xanthopoulos, P. & Nam, B.H. 2018. Soil-Compressibility Prediction
487        Models Using Machine Learning. Journal of Computing in Civil Engineering, **32**, doi:
488        10.1061/(asce)cp.1943-5487.0000713.

489    Kohavi, R. 1995. A study of Cross-Validation and bootstrap for accuracy estimation and model selection.
490        *International joint conference on artificial intelligence*. Morgan Kaufmann Publishers Inc., 1137-
491        1143.

492    Kohestani, V.R. & Hassanlourad, M. 2016. Modeling the mechanical behavior of carbonate sands using
493        artificial neural networks and support vector machines. International Journal of Geomechanics, **16**,
494        04015038, doi: 10.1061/(ASCE).

495    Li, Q., Ng, C.W.W. & Liu, G.B. 2012. Low secondary compressibility and shear strength of Shanghai
496        Clay. Journal of Central South University, **19**, 2323-2332, doi: 10.1007/s11771-012-1278-9.

497    Li, Z. & Chen, L. 2019. A novel evidential FMEA method by integrating fuzzy belief structure and grey
498        relational projection method. Engineering Applications of Artificial Intelligence, **77**, 136-147, doi:
499        10.1016/j.engappai.2018.10.005.

500    Lovatti, B.P.O., Nascimento, M.H.C., Neto, Á.C., Castro, E.V.R. & Filgueiras, P.R. 2019. Use of
501        Random forest in the identification of important variables. Microchemical Journal, **145**, 1129-1134,

502    doi: 10.1016/j.microc.2018.12.028.

503    Lu, Q., Li, X. & Dong, Y. 2018. Structure preserving unsupervised feature selection. Neurocomputing,
504        **301**, 36-45, doi: 10.1016/j.neucom.2018.04.001.

505    Meng, F.Y., Chen, R.P. & Xin, K. 2018. Effects of tunneling-induced soil disturbance on the post-
506        construction settlement in structured soft soils. Tunnelling and Underground Space Technology, **80**,
507        53–63, doi: 10.1016/j.tust.2018.06.007.

508    Nakase, A., Kamei, T. & Kusakabe, O. 1998. Constitutive parameters estimated by plasticity index
509        Journal of Geotechnical Engineering, **114**, 844-858.

510    Nassr, A., Esmaeili-Falak, M., Katebi, H. & Javadi, A. 2018. A new approach to modeling the behavior
511        of frozen soils. Engineering Geology, **246**, 82-90.

512    Penumadu, D. & Zhao, R.D. 1999. Triaxial compression behavior of sand and gravel using artificial
513        neural networks (ANN). Computers and Geotechnics, **24**, 207-230.

514    Pham, B.T., Son, L.H., Hoang, T.-A., Nguyen, D.-M. & Tien Bui, D. 2018. Prediction of shear strength
515        of soft soil using machine learning methods. Catena, **166**, 181-191, doi:
516        10.1016/j.catena.2018.04.004.

517    Pooya Nejad, F. & Jaksa, M.B. 2017. Load-settlement behavior modeling of single piles using artificial
518        neural networks and CPT data. Computers and Geotechnics, **89**, 9-21, doi:
519        10.1016/j.compgeo.2017.04.003.

520    Qi, C.C. & Tang, X.L. 2018a. A hybrid ensemble method for improved prediction of slope stability.
521        International Journal for Numerical and Analytical Methods in Geomechanics, **42**, 1823-1839, doi:
522        10.1002/nag.2834.

523    Qi, C.C. & Tang, X.L. 2018b. Slope stability prediction using integrated metaheuristic and machine
524        learning approaches: A comparative study. Computers & Industrial Engineering, **118**, 112-122, doi:
525        10.1016/j.cie.2018.02.028.

526    Qi, X.-H. & Zhou, W.-H. 2017. An efficient probabilistic back-analysis method for braced excavations
527        using wall deflection data at multiple points. Computers and Geotechnics, **85**, 186-198.

528    Rashidian, V. & Hassanlourad, M. 2014. Application of an artificial neural network for modeling the
529        mechanical behavior of carbonate soils. International Journal of Geomechanics, **14**, 142-150, doi:
530        10.1061/(asce)gm.1943-5622.0000299.

531    Romo, M.P., García, S.R., Mendoza, M.J. & Taboada‐Urtuzuástegui, V. 2001. Recurrent and
532        Constructive‐Algorithm Networks For Sand Behavior Modeling. International Journal of
533        Geomechanics, **1**, 371-387, doi: 10.1061/(asce)1532-3641(2001)1:4(371).

534    Shahin, M.A., Maier, H.R. & Jaksa, M.B. 2005. Investigation into the robustness of artificial neural
535        networks for a case study in civil engineering. In: Zerger A, Argent RM, editors. MODSIM 2005

536     International Congress on Modelling and Simulation. Modelling and Simulation Society of Australia
537     and New Zealand, pp.79-83.

538 Shen, S.L., Wu, H.N., Cui, Y.J. & Yin, Z.Y. 2014. Long-term settlement behaviour of metro tunnels in
539     the soft deposits of Shanghai. Tunnelling and Underground Space Technology, **40**, 309-323, doi:
540     10.1016/j.tust.2013.10.013.

541 Shi, Y. & Eberhart, R.C. 1998. Parameter selection in particle swarm optimization. *International*
542     *conference on evolutionary programming*. Springer, 591-600.

543 Stone, M. 1974. Cross-validatory choice and assessment of statistical predictions. Journal of the Royal
544     Statistical Society, **36**, 111-147, doi: 10.2307/2344741.

545 Tan, F., Zhou, W.-H. & Yuen, K.-V. 2016. Modeling the soil water retention properties of same-textured
546     soils with different initial void ratios. Journal of Hydrology, **542**, 731-743.

547 Tan, F., Zhou, W.-H. & Yuen, K.-V. 2018a. Effect of loading duration on uncertainty in creep analysis of
548     clay. International Journal for Numerical and Analytical Methods in Geomechanics, **42**, 1235-1254,
549     doi: 10.1002/nag.2788.

550 Tan, F., Zhou, W.H. & Yuen, K.V. 2018b. Effect of loading duration on uncertainty in creep analysis of
551     clay. International Journal for Numerical and Analytical Methods in Geomechanics, **42**, 1235-1254.

552 Trelea, I.C. 2003. The particle swarm optimization algorithm: convergence analysis and parameter
553     selection. Information processing letters, **85**, 317-325.

554 Turk, G., Logar, J. & Majes, B. 2001. Modelling soil behaviour in uniaxial strain conditions by neural
555     networks. Advances in Engineering Software, **32**, 805-812.

556 Xu, S.L. & Niu, R.Q. 2018. Displacement prediction of Baijiabao landslide based on empirical mode
557     decomposition and long short-term memory neural network in Three Gorges area, China. Computers
558     & Geosciences, **111**, 87-96, doi: 10.1016/j.cageo.2017.10.013.

559 Yang, B., Yin, K., Lacasse, S. & Liu, Z. 2019. Time series analysis and long short-term memory neural
560     network to predict landslide displacement. Landslides, doi: 10.1007/s10346-018-01127-x.

561 Yao, Y.-P., Qi, S.-J., Che, L.-W., Chen, J., Han, L.-M. & Ma, X.-Y. 2018. Postconstruction Settlement
562     Prediction of High Embankment of Silty Clay at Chengde Airport Based on One-Dimensional Creep
563     Analytical Method: Case Study. International Journal of Geomechanics, **18**, doi:
564     10.1061/(asce)gm.1943-5622.0001191.

565 Yin, J.H. 1999. Properties and behaviour of Hong Kong marine deposits with different clay contents.
566     Canadian Geotechnical Journal, **36**, 1085-1095.

567 Yin, J.H. & Graha, J. 1989. Viscous-elastic-plastic modelling of one-dimensional time-dependent
568     behaviour. Canadian Geotechnical Journal, **26**, 199-209.

569 Yin, Z.-Y. & Chang, C.S. 2009. Microstructural modelling of stress-dependent behaviour of clay.

570      International Journal of Solids and Structures, **46**, 1373-1388, doi: 10.1016/j.ijsolstr.2008.11.006.

571 Yin, Z.-Y., Chang, C.S., Karstunen, M. & Hicher, P.-Y. 2010a. An anisotropic elastic–viscoplastic model
572      for soft clays. International Journal of Solids and Structures, **47**, 665-677, doi:
573      10.1016/j.ijsolstr.2009.11.004.

574 Yin, Z.-Y., Yin, J.-H. & Huang, H.-W. 2014a. Rate-Dependent and Long-Term Yield Stress and Strength
575      of Soft Wenzhou Marine Clay: Experiments and Modeling. Marine Georesources & Geotechnology,
576      **33**, 79-91, doi: 10.1080/1064119x.2013.797060.

577 Yin, Z.-Y., Zhu, Q.-Y. & Zhang, D.-M. 2017. Comparison of two creep degradation modeling approaches
578      for soft structured soils. Acta Geotechnica, **12**, 1395-1413, doi: 10.1007/s11440-017-0556-y.

579 Yin, Z.Y., Karstunen, M., Chang, C.S., Koskinen, M. & Lojander, M. 2011. Modeling Time-Dependent
580      Behavior of Soft Sensitive Clay. Journal of Geotechnical and Geoenvironmental Engineering, **137**,
581      1103-1113, doi: 10.1061/(asce)gt.1943-5606.0000527.

582 Yin, Z.Y., Karstunen, M. & Hicher, P.Y. 2010b. Evaluation of the influence of elasto-viscoplastic scaling
583      functions on modelling time-dependent behaviour of natural clays. Soils and Foundations, **50**, 203-
584      214.

585 Yin, Z.Y., Xu, Q. & Yu, C. 2015. Elastic-Viscoplastic Modeling for Natural Soft Clays Considering
586      Nonlinear Creep. International Journal of Geomechanics, **15**, doi: 10.1061/(asce)gm.1943-
587      5622.0000284.

588 Yin, Z.Y., Zhu, Q.Y., Yin, J.H. & Ni, Q. 2014b. Stress relaxation coefficient and formulation for soft
589      soils. Géotechnique Letters, **4**, 45-51, doi: 10.1680/geolett.13.00070.

590 Zeng, L.L., Hong, Z.S., Liu, S.Y. & Chen, F.Q. 2012. Variation law and quantitative evaluation of
591      secondary consolidation behavior for remolded clays. Chinese Journal of Geotechnical Engineering,
592      **34**, 1496-1500.

593 Zhang, J., Zhang, L.M. & Tang, W.H. 2009. Bayesian framework for characterizing geotechnical model
594      uncertainty. Journal of Geotechnical and Geoenvironmental Engineering, **135**, 932-940, doi:
595      10.1061/共 ASCE 天 GT.1943-5606.0000018.

596 Zhang, L., Li, D.-Q., Tang, X.-S., Cao, Z.-J. & Phoon, K.-K. 2018. Bayesian model comparison and
597      characterization of bivariate distribution for shear strength parameters of soil. Computers and
598      Geotechnics, **95**, 110-118, doi: 10.1016/j.compgeo.2017.10.003.

599 Zheng, Y.-L., Ma, L.-H., Zhang, L.-Y. & Qian, J.-X. 2003. On the convergence analysis and parameter
600      selection in particle swarm optimization. *Proceedings of the 2003 International Conference on*
601      *Machine Learning and Cybernetics (IEEE Cat. No. 03EX693)*. IEEE, 1802-1807.

602 Zhou, J., Li, X. & Mitri, H.S. 2016a. Classification of Rockburst in Underground Projects: Comparison of
603      Ten Supervised Learning Methods. Journal of Computing in Civil Engineering, **30**, 04016003, doi:

604        10.1061/(asce)cp.1943-5487.0000553.

605    Zhou, W.-H., Garg, A. & Garg, A. 2016b. Study of the volumetric water content based on density,

606        suction    and    initial    water    content.    Measurement,    **94**,    531-537,    doi:

607        10.1016/j.measurement.2016.08.034.

608    Zhou, W.-H., Yuen, K.-V. & Tan, F. 2012. Estimation of maximum pullout shear stress of grouted soil

609        nails using Bayesian probabilistic approach. International Journal of Geomechanics, **13**, 659-664.

610    Zhou, W.H., Tan, F. & Yuen, K.V. 2018. Model updating and uncertainty analysis for creep behavior of

611        soft soil. Computers and Geotechnics, **100**, 135-143, doi: 10.1016/j.compgeo.2018.04.006.

612    Zhu, J.H., Zaman, M.M. & Anderson, S.A. 1998. Modelling of shearing behaviour of a residual soil with

613        Recurrent Neural Network. International Journal for Numerical and Analytical Methods in

614        Geomechanics, **22**, 671-687.

615    Zhu, Q.Y., Jin, Y.F. & Yin, Z.Y. 2019. Modeling of embankment beneath marine deposited soft sensitive

616        clays considering straightforward creep degradation. Marine Georesources & Geotechnology, in

617        press, doi: 10.1080/1064119X.2019.1603254.

618    Zhu, Q.Y., Yin, Z.Y., Hicher, P.Y. & Shen, S.L. 2016. Nonlinearity of one-dimensional creep

619        characteristics of soft clays. Acta Geotechnica, **11**, 887-900, doi: 10.1007/s11440-015-0411-y.

620

621


622

**Captions of figures**

**Fig. 1** Flowchart of building a random forest

**Fig. 2** Flow chart of the proposed $C_\alpha$ prediction model

**Fig. 3** Distributions of all variables in the database

**Fig. 4** GRG values among variables

**Fig. 5** Evolution of fitness value in all prediction models

**Fig. 6** Values of indicators in ten CV sets: (a) RMSE; (b) MAE; (c) MAPE

**Fig. 7** Distribution of indicators values in ten CV sets: (a) RMSE; (b) MAE; (c) MAPE

**Fig. 8** Predicted $C_\alpha$ for training and test sets by all prediction models

**Fig. 9** Predicted $C_\alpha$ for the test set by all published methods

**Fig. 10** Comparison of empirical methods and proposed models in predicting $C_\alpha$

**Fig. 11** Predicted $C_\alpha$ using RF model against (a) $CI$; (b) $w_L$; (c) $I_p$; (d) $e$

**Fig. 11** Predicted $C_\alpha$ using RF model against (a) $CI$; (b) $w_L$; (c) $I_p$; (d) $e$

**Fig. 12** Distribution of predicted $C_\alpha$

**Fig. 13** Mean decrease in Gini index of four input variables


**Captions of tables**

**Table 1** Hyper-parameters in random forest

**Table 2** Performance for all prediction models
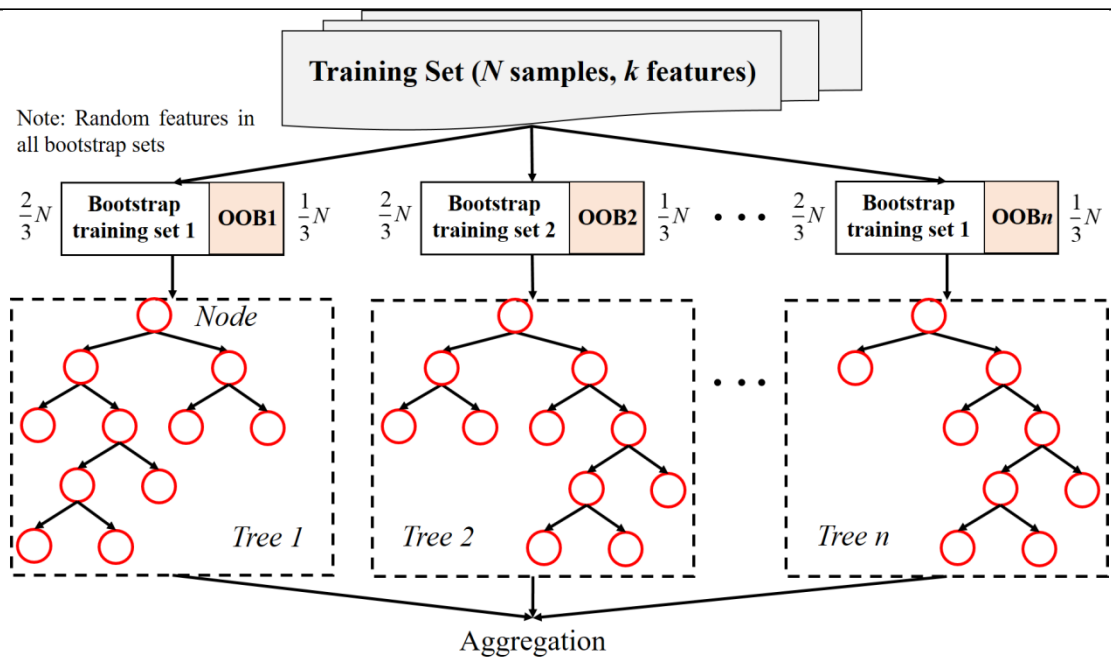
**Table 3** Performance for all prediction models
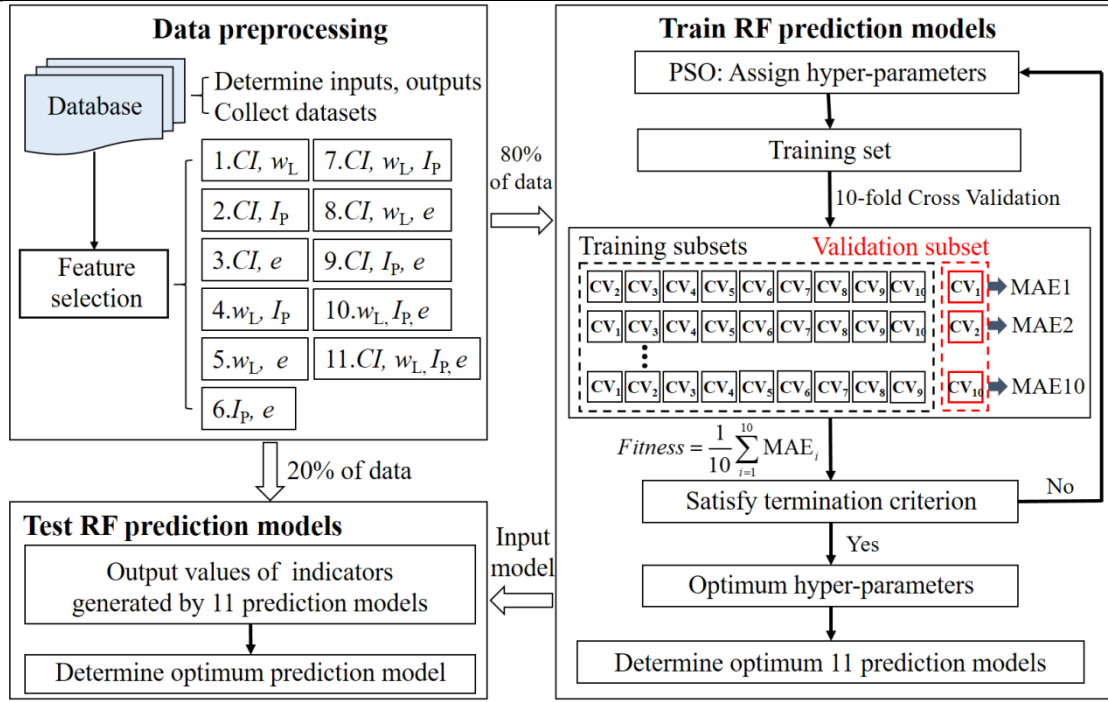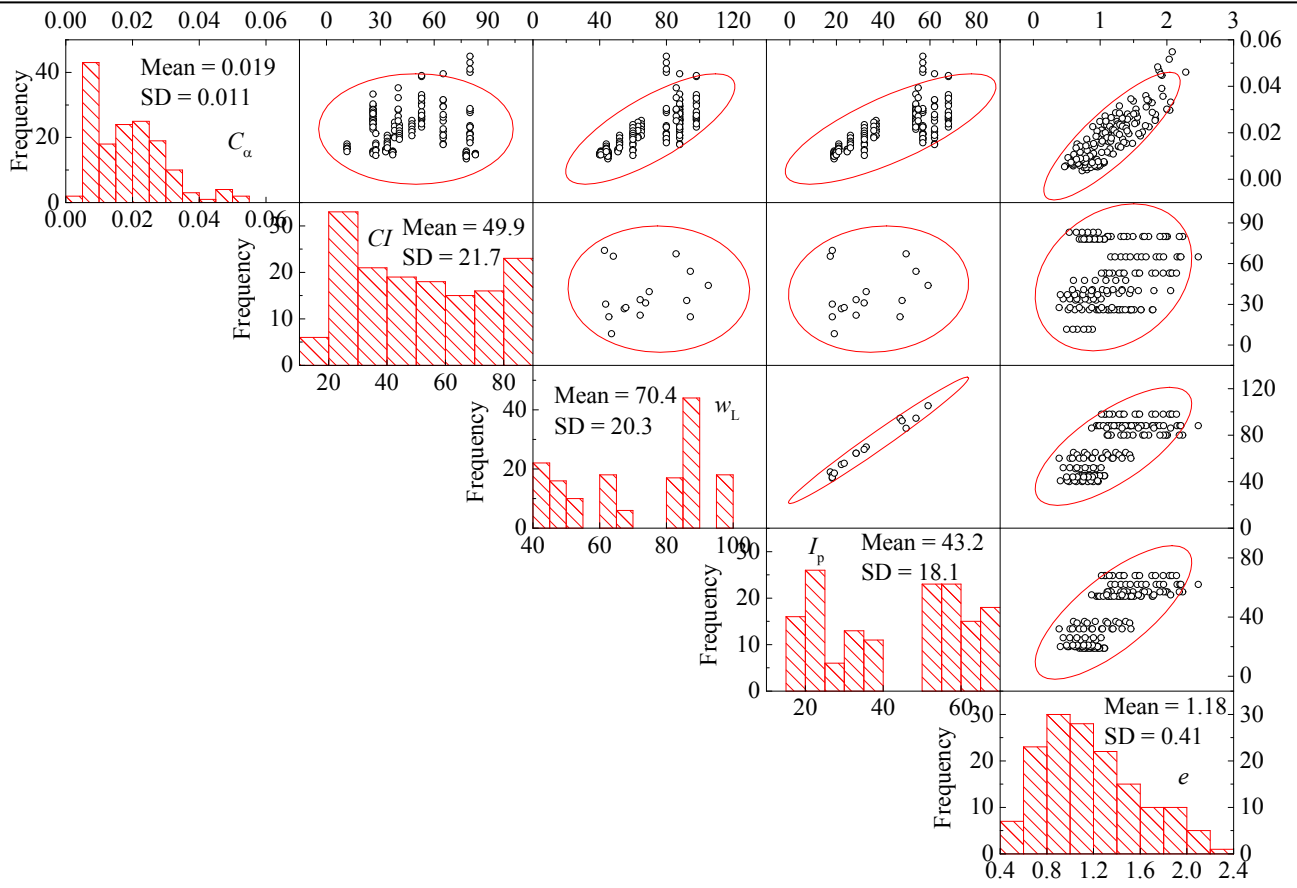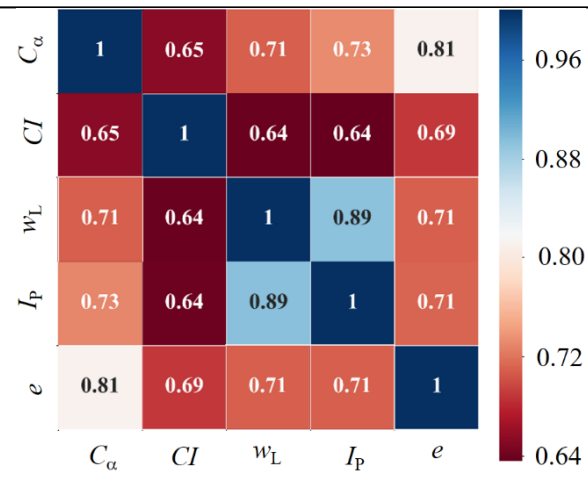
**Fig. 1** Flowchart of building a random forest

**Data preprocessing**

Database

Determine inputs, outputs
Collect datasets

| | |
|---|---|
| 1.$CI$, $w_L$ | 7.$CI$, $w_L$, $I_P$ |
| 2.$CI$, $I_P$ | 8.$CI$, $w_L$, $e$ |
| 3.$CI$, $e$ | 9.$CI$, $I_P$, $e$ |
| 4.$w_L$, $I_P$ | 10.$w_L$, $I_P$, $e$ |
| 5.$w_L$, $e$ | 11.$CI$, $w_L$, $I_P$, $e$ |
| 6.$I_P$, $e$ | |

Feature selection

80% of data

**Train RF prediction models**

PSO: Assign hyper-parameters

Training set

10-fold Cross Validation

Training subsets        Validation subset

| $CV_2$ | $CV_3$ | $CV_4$ | $CV_5$ | $CV_6$ | $CV_7$ | $CV_8$ | $CV_9$ | $CV_{10}$ | $CV_1$ | → MAE1 |
| $CV_1$ | $CV_3$ | $CV_4$ | $CV_5$ | $CV_6$ | $CV_7$ | $CV_8$ | $CV_9$ | $CV_{10}$ | $CV_2$ | → MAE2 |
| $CV_1$ | $CV_2$ | $CV_3$ | $CV_4$ | $CV_5$ | $CV_6$ | $CV_7$ | $CV_8$ | $CV_9$ | $CV_{10}$ | → MAE10 |

$$Fitness = \frac{1}{10}\sum_{i=1}^{10} \text{MAE}_i$$

Satisfy termination criterion      No

Yes

Optimum hyper-parameters

Determine optimum 11 prediction models

20% of data

Input model

**Test RF prediction models**

Output values of indicators generated by 11 prediction models

Determine optimum prediction model

**Fig. 2** Flow chart of the proposed $C_\alpha$ prediction model

**Fig. 3** Distributions of all variables in the database

**Fig. 4** GRG values among variables

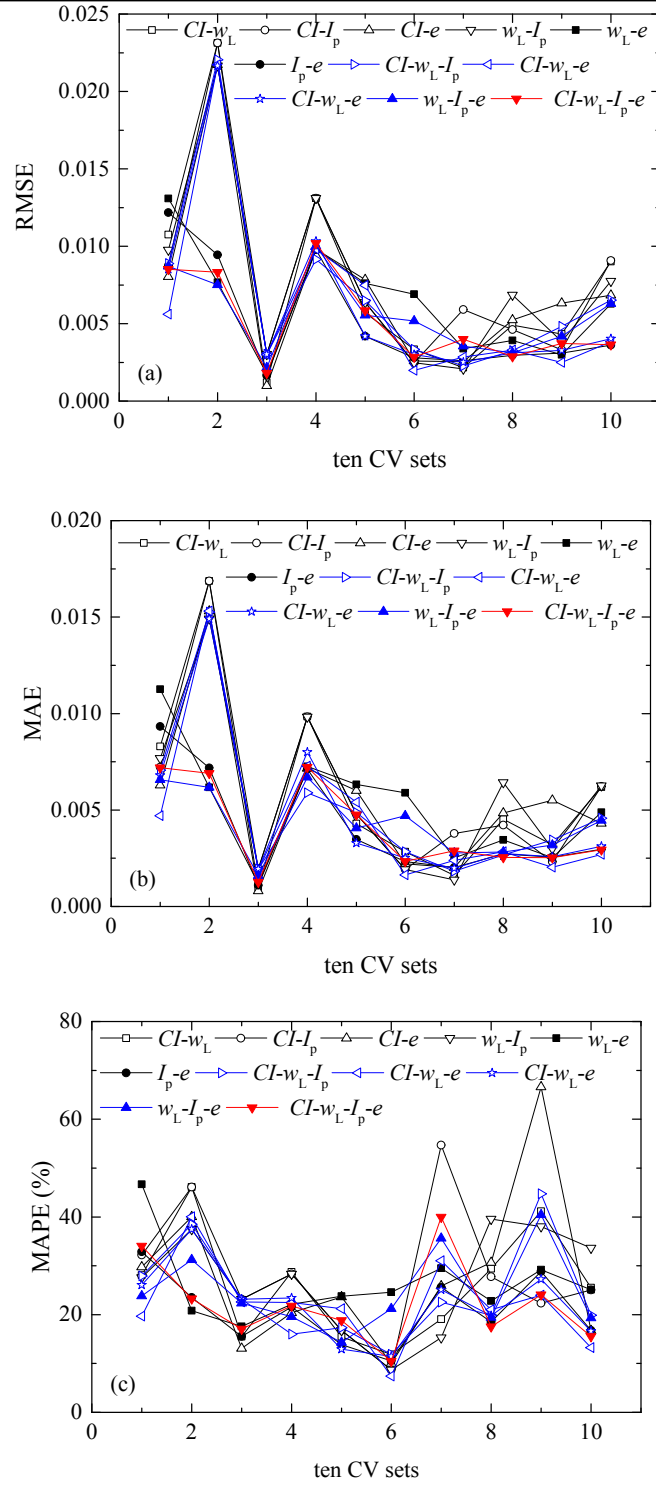**Fig. 5** Evolution of fitness value in all prediction models

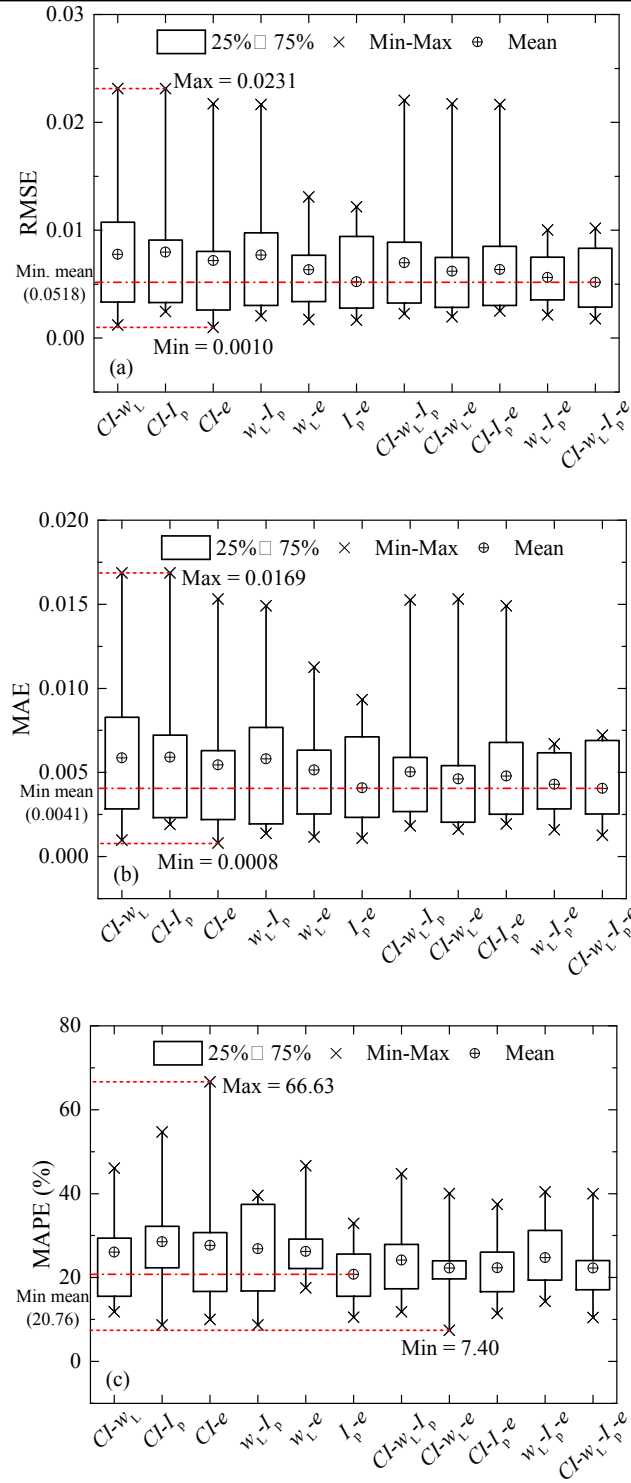**Fig. 6** Values of indicators in ten CV sets: (a) RMSE; (b) MAE; (c) MAPE

**Fig. 7** Distribution of indicators values in ten CV sets: (a) RMSE; (b) MAE; (c) MAPE
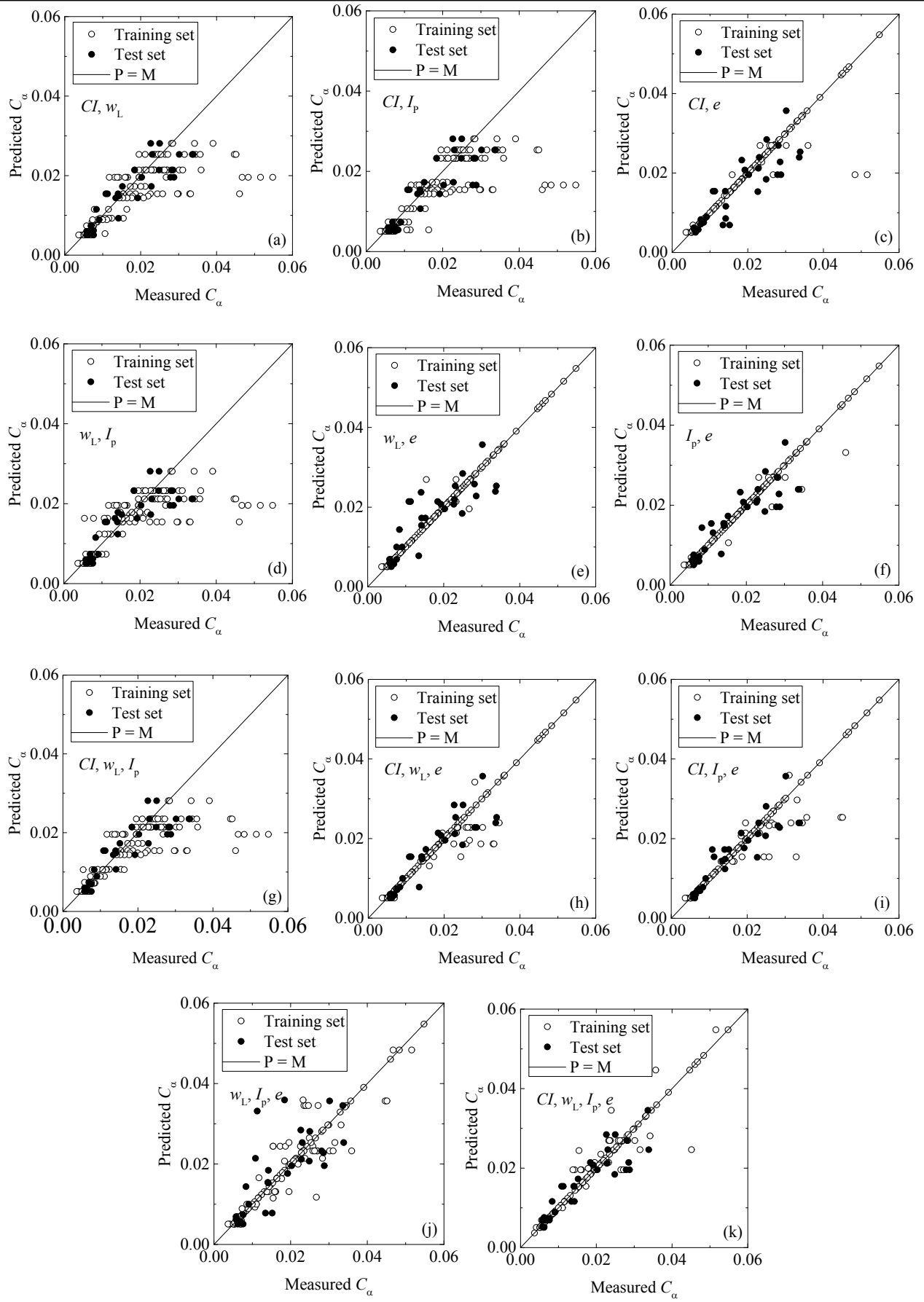
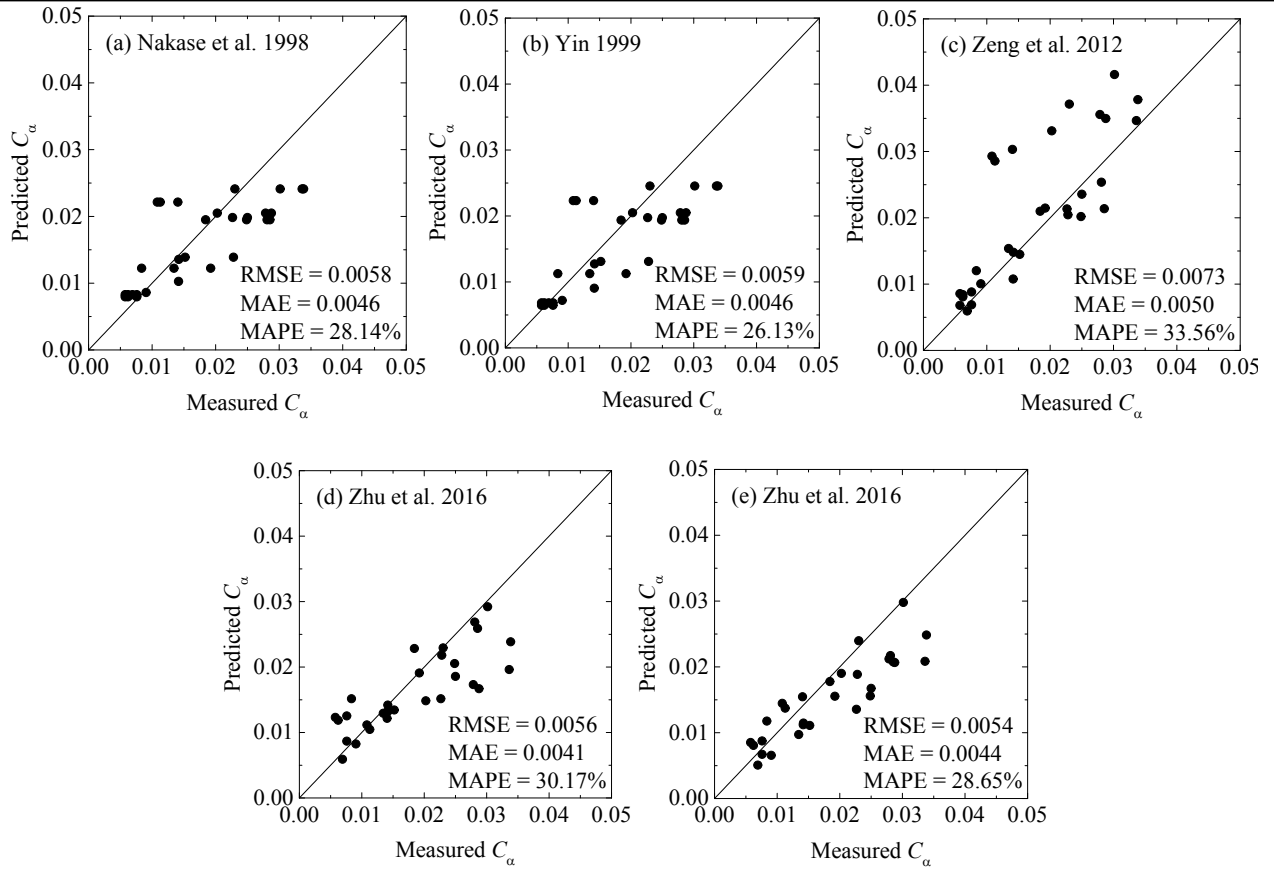**Fig. 8** Predicted $C_\alpha$ for training and test sets by all prediction models

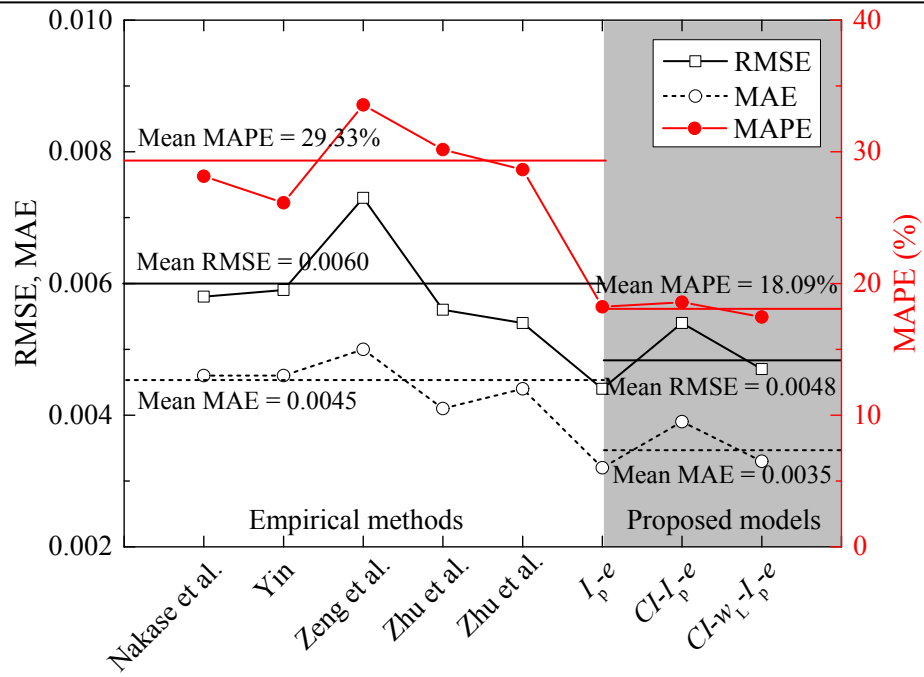**Fig. 9** Predicted $C_\alpha$ for the test set by all published methods

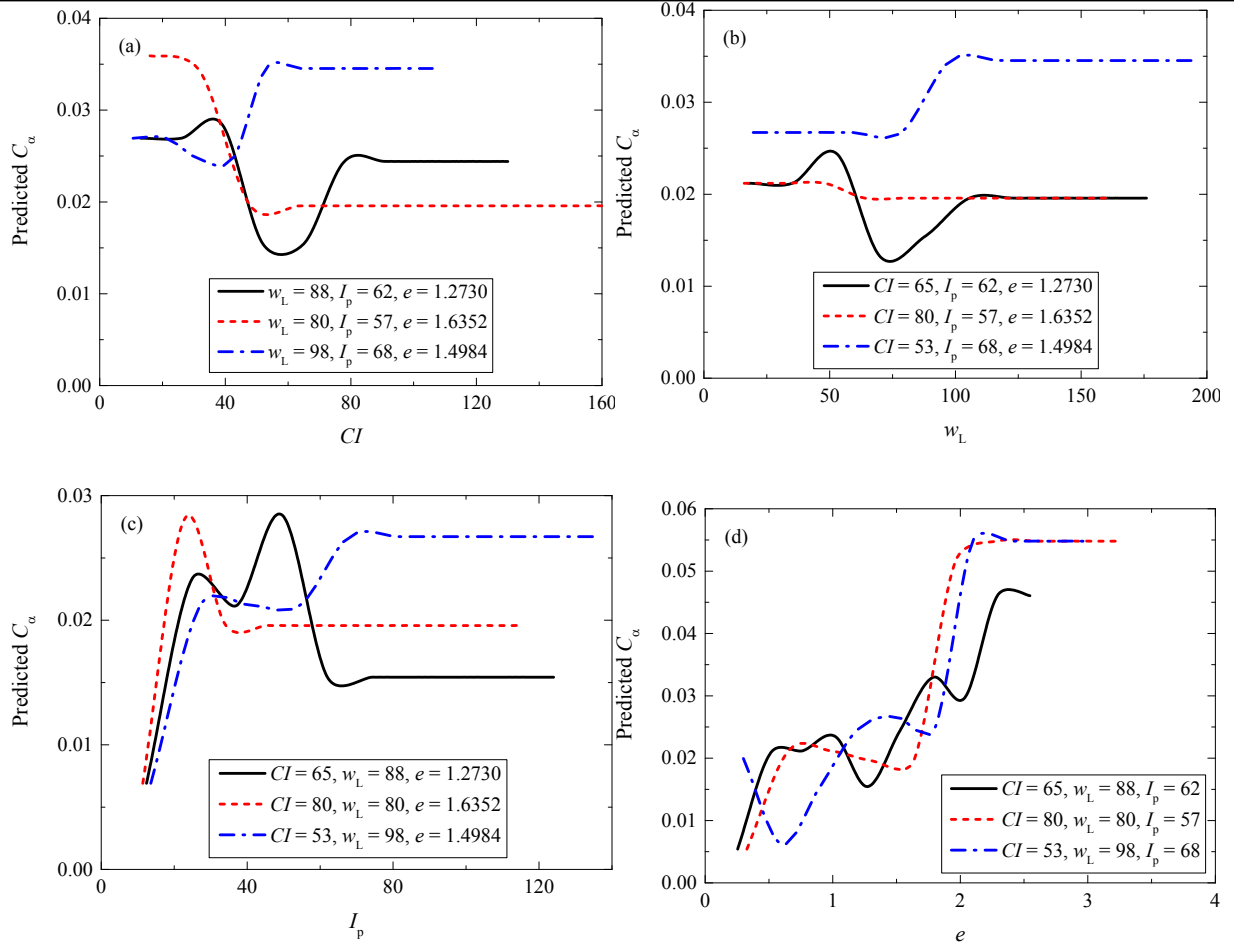**Fig. 10** Comparison of empirical methods and proposed models in predicting $C_\alpha$

**Fig. 11** Predicted $C_\alpha$ using RF model against (a) $CI$; (b) $w_L$; (c) $I_p$; (d) $e$
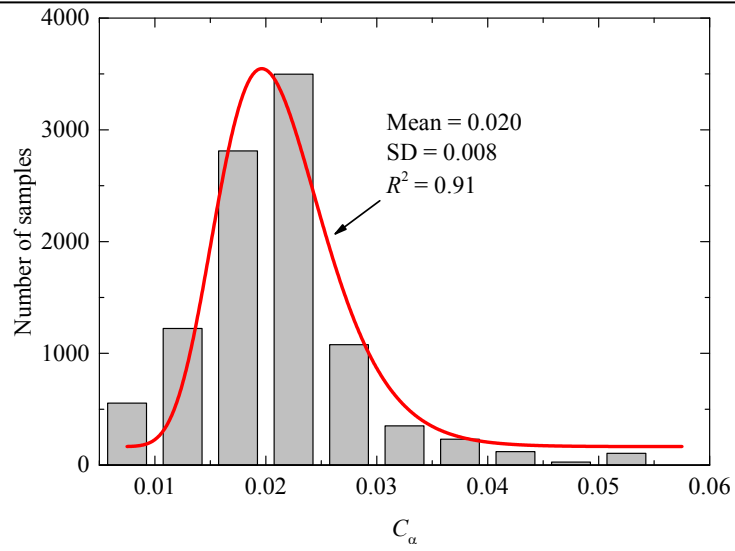
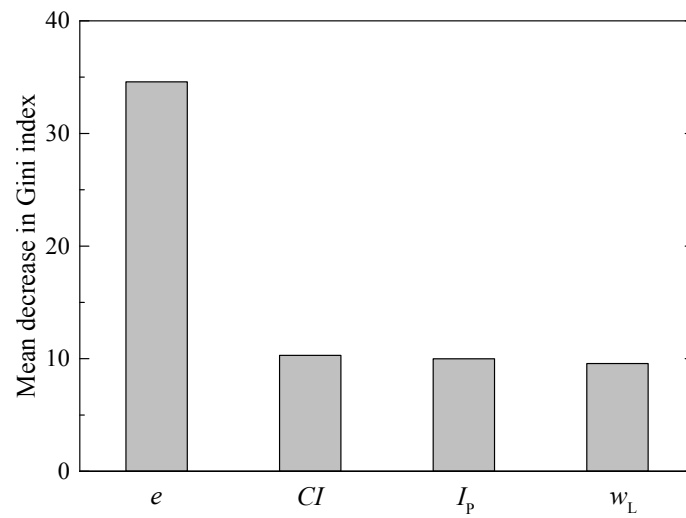**Fig. 12** Distribution of predicted $C_\alpha$



**Fig. 13** Mean decrease in Gini index of four input variables

**Table 1** Hyper-parameters in random forest

| Hyper-parameters | Description | Range |
|---|---|---|
| *ntree* | Number of trees grown | 0~300 |
| *mtry* | Number of predictors sampled for spliting at each node | 0~300 |

**Table 2** Performance for all prediction models

| Variables combination | Training set | | | Test set | | |
|---|---|---|---|---|---|---|
| | RMSE | MAE | MAPE | RMSE | MAE | MAPE |
| $CI, w_L$ | 0.0084 | 0.0046 | 18.75% | 0.0044 | 0.0035 | 19.05% |
| $CI, I_P$ | 0.0087 | 0.0046 | 18.36% | 0.0045 | 0.0037 | 20.02% |
| $CI, e$ | 0.0041 | 0.0008 | 2.83% | 0.0048 | 0.0037 | 19.57% |
| $w_L, I_P$ | 0.0087 | 0.0048 | 21.26% | 0.0052 | 0.0039 | 20.56% |
| $w_L, e$ | 0.0013 | 0.0002 | 1.75% | 0.0055 | 0.0048 | 27.77% |
| **$I_P, e$** | **0.0018** | **0.0004** | **2.35%** | **0.0044** | **0.0032** | **18.23%** |
| $CI, w_L, I_P$ | 0.0085 | 0.0047 | 19.01% | 0.0047 | 0.0036 | 18.37% |
| $CI, w_L, e$ | 0.0030 | 0.0011 | 4.82% | 0.0055 | 0.0040 | 19.42% |
| **$CI, I_P, e$** | **0.0037** | **0.0012** | **4.79%** | **0.0054** | **0.0039** | **18.58%** |
| $w_L, I_P, e$ | 0.0039 | 0.0020 | 9.84% | 0.0055 | 0.0051 | 31.72% |
| **$CI, w_L, I_P, e$** | **0.0030** | **0.0011** | **5.05%** | **0.0047** | **0.0033** | **17.45%** |

Note: bold format donates the optimum combination in each group

**Table 3** Empirical correlations for $C_\alpha$

| References | Empirical formula |
|---|---|
| Nakase et al. (1998) | $C_\alpha = 0.00168 + 0.00033 I_p$ |
| Yin (1999) | $C_\alpha = 0.000369 I_p - 0.00055$ |
| Zeng et al. (2012) | $C_\alpha = \left(-0.0067 + 0.0115 e_L - 0.0016 (e_L)^2\right)(1+e)$ |
| Zhu et al. (2016) | $C_\alpha = \left(-0.0274 + 0.0011 w_L - 0.00048 I_p\right)\left(\dfrac{w}{w_L}\right)^{0.7872 - 0.0369 w_L + 0.0619 I_p}$ |
| Zhu et al. (2016) | $C_\alpha = \left(0.0007 w_L - 0.0223\right)\left(\dfrac{w}{w_L}\right)^{0.014978 w_L - 0.23031}$ |