

A Novel Machine Learning-based Energy Consumption Model of Wastewater Treatment Plants

Shike Zhang¹, Hongtao Wang^{1,2*}, Arturo A. Keller³

1 College of Environmental Science and Engineering, Key Laboratory of Yangtze River Water Environment, Ministry of Education, UNEP-Tongji Institute of Environment for Sustainable Development, Tongji University, 1239 Siping Rd, Shanghai, 200092, China

2 Shanghai Institute of Pollution Control and Ecological Security, Tongji University, 1239 Siping Rd, Shanghai, 200092, China.

3 Bren School of Environmental Science and Management, UC Santa Barbara, Santa Barbara, California 93106, United States

***Corresponding author: Email: hongtao@tongji.edu.cn (Hongtao Wang)**

Abstract: Wastewater treatment plants (WWTPs) can account for up to 1% of a country's energy consumption. Meanwhile, WWTPs have high energy-saving potential. To achieve this, it is necessary to establish appropriate energy consumption models for WWTPs. Several recent models have been developed using logarithmic, exponential, or linear functions. However, the behavior of WWTPs is non-linear, and difficult to fit with simple functions particularly for non-numerical variables. Thus, traditional modeling methods cannot effectively describe the relationship between water and energy in WWTPs. Therefore, a machine learning method was adopted in this study to investigate the energy consumption in WWTPs; a novel energy consumption model with a non-numerical variable (discharge standard) for WWTPs was developed using the random forest algorithm. The model can also predict the energy consumption of WWTPs after upgrading discharge standards. We found that the unit electricity consumption of WWTPs exhibited an average increase of 17% after the effluent

29 discharge standard was raised from Class I B to Class I A (per China's classification).

30 The correlation coefficient of the model was 0.702. Thus, the developed model can

31 provide a better understanding of energy efficiency in WWTPs.

32 **Keywords:** Machine learning; random forest; energy efficiency; energy consumption

33 model; wastewater treatment plant

34

35 **1. Introduction**

36 Improving energy efficiency of WWTPs is receiving increasing attention, as saving

37 energy can help reduce economic costs and conserve the resources and environment

38 (E.Açikkalp, 2018). With continuous development and accelerating urbanization in

39 society, wastewater discharge is rapidly increasing (Habib et al. 2020) and water quality

40 requirements are more stringent; therefore, the total energy consumption of WWTPs is

41 also increasing. WWTPs are the primary energy-consuming units of the urban water

42 cycle (Sabia et al. 2020). Thus, the high energy consumption of WWTPs has become a

43 global concern. It has been estimated that in 2018 the energy demand of WWTPs in

44 some European countries accounted for 1% of the energy consumption of the entire

45 country (Sabia et al. 2020). What's more, the U.S. municipal wastewater treatment

46 systems use approximately 30.2 billion kWh per year, which is about 0.8% of the total

47 electricity use in the U.S. (EPRI, 2013). In recent years, several energy evaluation

48 methods have been proposed (e.g., Hernández-Sancho et al. 2009, Mizuta et al. 2010,

49 Molinos-Senante et al. 2014) to investigate the energy consumption of WWTPs. In
50 these methods, the energy consumption of WWTPs is commonly related to factors such
51 as the capacity and influent and effluent concentrations of pollutants (Torregrossa et al.
52 2016).

53 Many stakeholders have been exploring solutions to reduce the energy consumption of
54 WWTPs, such as equipment renewal and maintenance (Daw et al. 2012), energy
55 recovery (Behera et al. 2020), and technical process improvements (Farahbakhsh et al.
56 2020). Previous studies have been mostly focused on technical processes. However,
57 with the use of new equipment, technologies, and new standards, the change in energy
58 consumption has gradually attracted significant attention (Sabia et al. 2020). With
59 increasing urban wastewater discharge and high requirements of clean water, new
60 WWTPs are regularly established, and the discharge standards of the old WWTPs have
61 gradually improved (Smith et al. 2019). However, a larger process capacity and higher
62 standards may lead to higher energy consumption. Therefore, while considering water
63 quality, one should also focus on the energy efficiency of WWTPs.

64 Machine learning is an important and relatively novel method in environmental
65 modeling, particularly with regards to energy efficiency (e.g., Wang et al. 2019) or
66 WWTP operations (e.g. Hernandez-del-Olmo et al. 2019). Machine learning can be
67 utilized in a real-time agent modeling, employing real-time data so that operators can
68 forecast WWTP's future operating status; the model itself can be improved

69 continuously as new data becomes available, with the ability to adopt non-linear
70 relationships. Once a set of inputs and corresponding outputs are presented to the model,
71 it learns the relationship between the inputs and outputs. Accordingly, for a new set of
72 inputs, the trained model can generalize this relationship to produce the corresponding
73 outputs (e.g., Heslot et al. 2014, Song et al. 2017, Xing et al. 2019, Zhu et al. 2020).

74 Random forest is an ensemble learning algorithm used for classification (e.g., Duro et
75 al. 2012), regression (e.g., Wei et al. 2019), and other tasks (e.g., Chen et al. 2018).
76 During training, numerous decision trees are generated to operate and finally obtain
77 prediction results (e.g., Tian et al. 2020); therefore, random forest has a high prediction
78 accuracy, and is not prone to overfitting (e.g., Breiman 2001). Random forest is one of
79 the most popular methods in data mining (e.g., Wylie et al. 2019) and big data fields
80 (e.g., Pamulaparty et al. 2017); it has the advantages of a fast-training speed and is
81 suitable for processing high-dimensional data (e.g., Belgiu and Dragut 2016). The
82 method has been widely used in several other fields, such as medicine (Yeşilkanat 2020),
83 criminal investigation (Tian et al. 2020), and architecture (Cheng et al. 2020). Random
84 forest has also been applied to environmental engineering, such as for mapping canopy
85 nitrogen (Loozen et al. 2020) and in environmental assessment (Paul et al. 2020).
86 However, random forest is rarely used to analyze and predict energy consumption of
87 WWTPs (e.g., Bagherzadeh et al. 2021, Perez et al. 2021). Data Envelopment Analysis
88 (e.g., Yang et al. 2021, Huang et al. 2021) and Multiple Linear Regression (e.g., Xu et

89 al. 2018) are commonly used methods to analyze the energy efficiency of WWTPs. The
90 main principle of Data Envelopment Analysis is using the method of linear
91 programming and Multiple Linear Regression is a kind of generalized linear model.
92 However, Data Envelopment Analysis cannot be used when some data is missing, and
93 it cannot predict the future trend in a statistical way. Multiple Linear Regression cannot
94 achieve the accuracy we need.

95 This study aims to develop an energy consumption model of WWTPs through machine
96 learning, using data from 2,472 WWTPs in China, employing the random forest
97 approach. This model is expected to provide a better understanding of energy
98 consumption in WWTPs.

99

100 **2. Data and Methods**

101 **2.1 Data Sources**

102 The data of this study was collected from the *2015 Urban Drainage Yearbook of China*.

103 A total of 2,472 entries were selected from this yearbook with relatively complete and
104 reliable data. For the energy consumption of WWTPs, the energy consumed by
105 processing 1 m³ wastewater is often used as an evaluation indicator of the energy
106 intensity of WWTPs (Scott et al. 2011). Related studies (e.g., Mjalli et al. 2007,
107 Gazendam et al. 2016, Trenouth et al. 2018, Habib et al. 2020) commonly use electricity
108 intensity (kWh/m³) to indicate the energy consumption of WWTPs. In this study, the

109 following parameters from the Yearbook, which have also been used in related studies
110 (Mjalli et al. 2007, Gazendam et al. 2016, Trenouth et al. 2018, Habib et al. 2020), have
111 been adopted as the primary factors affecting the energy consumption of WWTPs:
112 influent BOD₅ concentration (BOD_i), influent COD concentration (COD_i), influent
113 NH₃-N concentration (NH₃-N_i), effluent BOD₅ concentration (BOD_e), effluent COD
114 concentration (COD_e), effluent NH₃-N concentration (NH₃-N_e), effluent discharge
115 standards, wastewater treatment capacity, annual load rate (actual treatment capacity
116 divided by designed treatment capacity), moisture content of sludge, and dry weight of
117 sludge. The discharge standards primarily include Class I A, Class I B, and Class II,
118 referring to the Chinese National Standard, *Discharge standard of pollutants for*
119 *municipal wastewater treatment plant* (GB 18918-2002).

120 **2.2 Data Cleaning**

121 To import data to develop the model, all numerical data (including int64, float64) was
122 converted to float64, all non-numerical data (including string and object, such as the
123 discharge standard) was transferred into the object, and all the default data was
124 converted to NaN (Not a Number).

125 Since a few of the WWTPs did not have the “unit electricity consumption” (UEC)
126 parameter in the Yearbook, to ensure the reliability of the model, we removed the data
127 for those 85 WWTPs, so that the number of the remaining WWTPs was 2,387.

128 Although there were some outliers, the data base is large and Random Forest model is

129 good at dealing with this situation, so there was no need to eliminate them. Essentially,
130 the basic learner of random forest is robust to outliers, which makes the random forest
131 algorithm robust to outliers. Unlike linear regression, the entire space in linear
132 regression has the same equation, so a very simple model can be locally fitted to each
133 subspace.

134 In the case of regression, it is usually a very low-order regression model. Therefore, for
135 regression, extreme values do not affect the entire model because they are averaged
136 locally.

137

138 **2.3 Preprocessing of Regression Variables**

139 The numerical variables can be directly applied to the regression. For the object type
140 variables, such as discharge standards, their classification scheme was transformed into
141 a matrix with 0 and 1 values, such that, the row of the matrix represents the different
142 WWTPs, and the column represents the different discharge standards. A value of 1
143 indicated that the WWTP represented by this row used the discharge standard of this
144 column. Otherwise, a value 0 was assigned. For example ([Fig. 1](#)), the discharge
145 standards of WWTP A, WWTP B, and WWTP C are Class I A, Class II, and Class I B,
146 respectively. Therefore, the sum of each row in the matrix is 1, and the sum of each
147 column equals the total number of WWTPs using the discharge standard of this column.
148 In this study, numerical values (e.g., 1, 2, 3, etc.) were not used to represent the different

149 discharge standards. This is because the values themselves include the potential
 150 relationship of size or numerical operation, and substituting a relationship which is
 151 unrelated to the statistical content into the regression model will lead to model deviation.

$$\begin{array}{r}
 \text{WWTP A} \\
 \text{WWTP B} \\
 \text{WWTP C}
 \end{array}
 \begin{array}{c}
 \text{Class I A} \\
 \text{Class I B} \\
 \text{Class II}
 \end{array}
 \begin{pmatrix}
 1 & 0 & 0 \\
 0 & 0 & 1 \\
 0 & 1 & 0
 \end{pmatrix}$$

152

153

Fig. 1 Matrix of discharge standard in WWTPs

154

155 2.4 Random Forest

156 A preliminary evaluation of the relationship between UEC (kWh/m³) and other
 157 parameters was performed using Python, to conduct Multiple Linear Regression
 158 analysis between UEC and BOD_i, COD_i, NH₃-N_i, BOD_e, COD_e, NH₃-N_e, wastewater
 159 treatment capacity, annual load rate, moisture content of sludge, and dry weight of
 160 sludge. It was found that $R^2 \leq 0.2$, which is too small, thus the regression equation was
 161 not sufficiently reliable. Moreover, the discharge standard is a character-type variable
 162 that cannot be included in the statistics and effectively predict its influence by Multiple
 163 Linear Regression. Owing to the high correlation between the parameters, a large fitting
 164 deviation occurred when using the Multiple Linear Regression method to obtain the
 165 relationship between UEC and dependent variables. In addition, the use of a Multiple
 166 Linear Regression model is limited in this case because of the existence of non-
 167 numerical variables, such as the discharge standard.

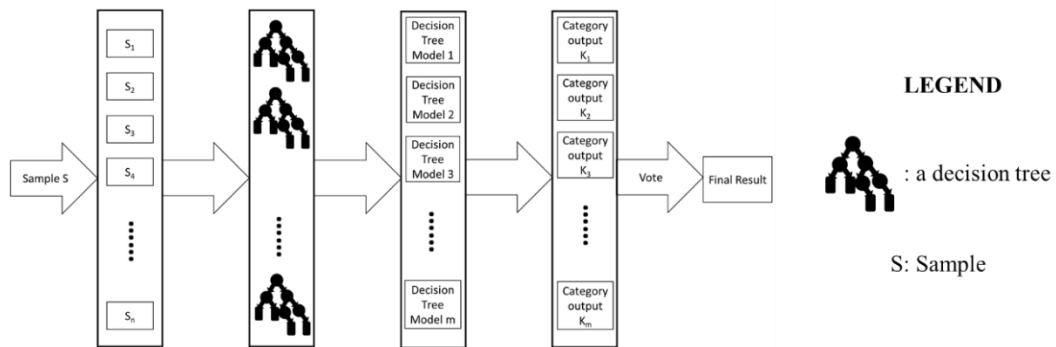
168 Machine learning algorithms like TensorFlow or Keras require very large databases,
169 which are not available for this study. While there are machine learning algorithms like
170 Lars, Lasso or Support Vector Machine (SVM) which only need small databases, they
171 cannot reach the accuracy needed for this study. Therefore, we considered a subset of
172 machine learning algorithms, including Random Forest (RF), Boosting Tree, Gradient
173 Boosting Decision Tree (GBDT) and XGBoost. These algorithms are actually
174 combination of different algorithmic frameworks and decision trees (Detail see [Table](#)
175 [S1](#)), so they perform quite similarly. We selected Random Forest because it is the only
176 algorithm that can show us the importance of each variable, which is very valuable for
177 the subsequent analysis.

178 Therefore, a random forest algorithm was introduced to extract the relationship between
179 UEC and the different variables, including non-numerical ones. Simultaneously, the
180 factors indicating the influence of each variable on UEC were calculated, and then the
181 factors that significantly affected UEC were selected for further analysis and to develop
182 a model for evaluating the UEC. Finally, the change in UEC was calculated using the
183 model after a simulated improvement of the discharge quality to meet a higher standard,
184 which can help in future management of WWTPs.

185 The steps conducted for the random forest approach were showed in [Fig. 2](#).

186 From a mathematical perspective, a complex functional relationship exists between
187 independent variables and dependent variables, which is composed of the basic

188 operations of independent variables. Random forest approximates the coefficients
 189 before each dependent variable by learning from a large amount of data. All the models
 190 in this study were coded in Python 3.7.3, and the prediction curves were plotted from
 191 the Python data using MATLAB R2018a.



192
 193 Fig. 2 Process flow of the random forest method

194
 195 **2.5 Model Validation**

196 Random forest uses a bootstrapping algorithm for sampling. As the bootstrapping
 197 algorithm returns samples after sampling, some data are not extracted. By calculating
 198 the limit, it was observed that approximately 1/3 of the data were not extracted.

199 Because out-of-bag (OOB) data were not been used, random forest can use these data
 200 for model validation. Moreover, as each sample obtained by bootstrap trains a small
 201 model S_n , the OOB data can be tested for each model of the sample.

202 The self-detection of the model uses the mean squared error (MSE), average absolute
 203 percentage (MAPE), root mean square error (RSME), mean absolute error (MAE),
 204 median absolute error (MedAE) and mean squared logarithmic error (MSLE) (Zhong

205 et al. 2021, Gupta et al. 2021) as follows:

$$206 \quad MSE = \frac{1}{n} \sum_{t=1}^n (actual(t) - predicted(t))^2 \quad (1)$$

$$207 \quad MAPE = \frac{1}{n} \sum_{t=1}^n \frac{|actual(t) - predicted(t)|}{actual(t)} \times 100\% \quad (2)$$

$$208 \quad RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n (actual(t) - predicted(t))^2} = \sqrt{MSE} \quad (3)$$

$$209 \quad MAE = \frac{1}{n} \sum_{t=1}^n |actual(t) - predicted(t)| \quad (4)$$

$$210 \quad MedAE = median(|actual(t_1) - predicted(t_1)|, \dots, |actual(t_n) - predicted(t_n)|) \quad (5)$$

$$211 \quad MSLE = \frac{1}{n} \sum_{t=1}^n (\log_e(1 + actual(t)) - \log_e(1 + predicted(t)))^2 \quad (6)$$

212 where n is the number of decision tree model, $actual(t)$ is the actual UEC of a

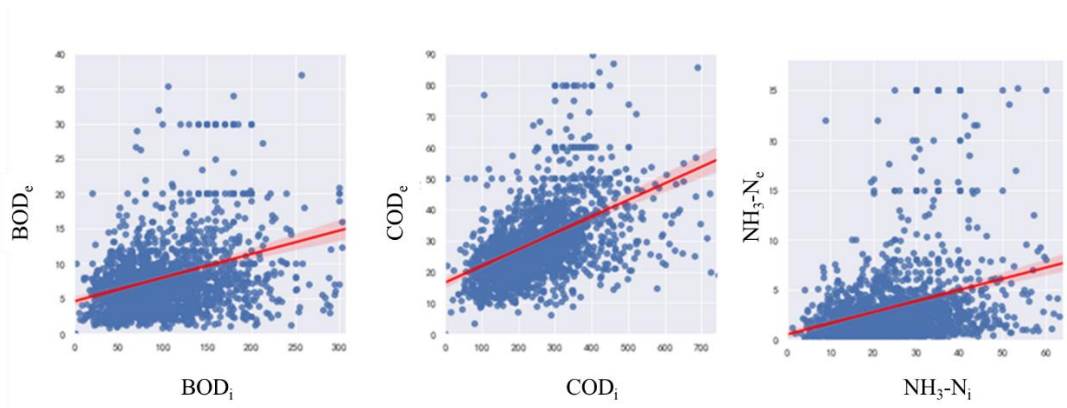
213 WWTP, and $predicted(t)$ is the predicted UEC of a WWTP.

214

215 2.6 Data Preprocessing

216 Certain evident linear relationships exist between some variables in the Yearbook,

217 which were removed before modeling to obtain a model with improved accuracy.



218

219 Fig. 3 Scatter plot and linear regression curve of

220 (a) BOD_i/BOD_e , (b) COD_i/COD_e , and (c) NH_3-N_i/NH_3-N_e .

221 In Fig. 3, each point in the figure represents a WWTP dataset, and it may be noted that
222 the BOD_i and BOD_e , COD_i and COD_e , NH_3-N_i and NH_3-N_e , for most WWTPs are
223 evenly distributed along a straight line. The light-red area around the regression curve
224 represents the confidence interval. Therefore, we performed a linear regression between
225 BOD_i and BOD_e , COD_i and COD_e , NH_3-N_i and NH_3-N_e , which showed that the linear
226 correlation between BOD_i and BOD_e , COD_i and COD_e , NH_3-N_i and NH_3-N_e was high.
227 At the same time, the correlation between each variable was reduced to the lowest
228 possible limit to improve the accuracy of machine learning. In this study, the removal
229 ratios BOD_i/BOD_e , COD_i/COD_e , and NH_3-N_i/NH_3-N_e were used instead of a single
230 variable for analysis, which practically represent the reduction multiple of BOD, COD,
231 and NH_3-N_e of treated wastewater.

232

233 **2.7 The Importance of Features**

234 The importance of a feature X in a random forest was calculated as follows:

235 A. For each decision tree in a random forest, the corresponding OOB data were used

236 to calculate the OOB data error, which was recorded as err_{OOB1} .

237 B. Random noise interference was added to the characteristic X of all samples of the

238 OOB data, and the OOB data error was calculated again, which was recorded

239 as err_{OOB2} .

240 C. If there are N trees in the random forest, then the importance of the feature is

241 given as

$$242 X_{Importance} = \sum (\text{errOOB2} - \text{errOOB1}) / N \quad (7)$$

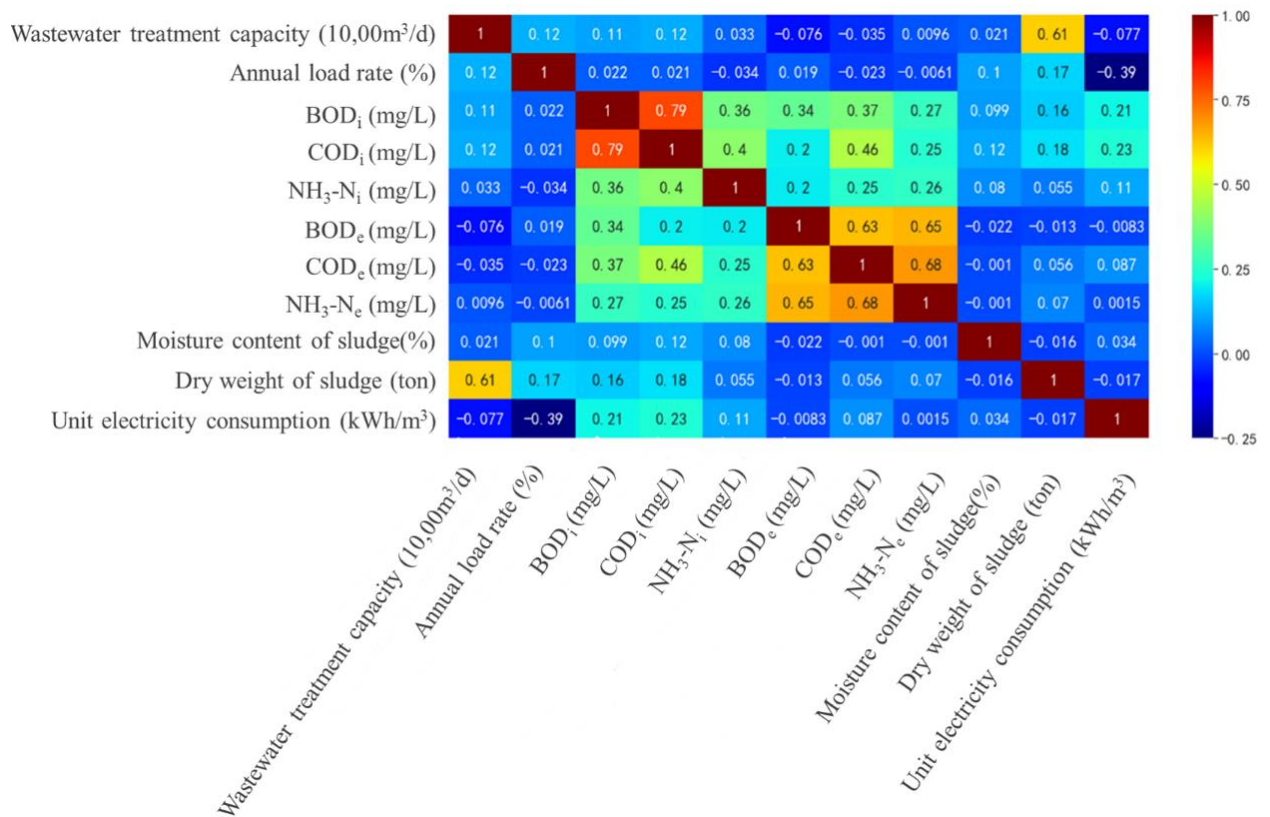
243 This expression can be used as a measure of the importance of corresponding features
 244 because if a feature was randomly added with noise, the accuracy rate outside the bag
 245 was highly reduced, which indicated that this feature had a high influence on the
 246 classification results of samples; in other words, it was of high importance.

247

248 3. Results and Discussion

249 3.1 Correlation between Variables

250 To accurately analyze the relationship between UEC and the different variables, we
 251 calculated the correlation between these variables, as shown in Fig. 4.



252

253 Fig. 4 Correlation thermodynamic diagram of variables

254 In Fig. 4, the correlation of UEC, wastewater treatment capacity, annual load rate,
255 moisture content of sludge, dry weight of sludge, BOD_i, BOD_e, COD_i, COD_e, NH₃-N_i,
256 and NH₃-N_e are described by the thermodynamic diagram. The number on the color
257 block represents the correlation between the corresponding variables of the abscissa
258 and ordinate. A darker red implies a higher correlation, and a darker blue indicates a
259 lower correlation. The UEC is highly correlated with BOD, COD, and NH₃-N of
260 influent and effluent (Fig. 4).

262 3.2 Regression

263 The analysis of the importance of the independent variables is presented in Table 1 and
264 Fig. S1. The regression model had an $R^2 = 0.702$, which was significantly higher than
265 that of the Multiple Linear Regression, implying higher accuracy.

266 As shown in Table 1 and Fig. S1, the most important variable was the wastewater
267 treatment capacity, which is expected since this determines the sizing of pumps, air
268 blowers and other equipment that consumes electricity (Torregrossa et al. 2018). This
269 is followed by the annual load rate, which is also expected to be a major factor
270 (Torregrossa et al. 2018). Wastewater treatment capacity and annual load rate can reflect
271 the influence of the design and practical operation of WWTPs on energy consumption,
272 with a total importance of 0.38. The high importance indicated that the design of a

273 WWTP was very important, so a clear treatment target would significantly affect the
 274 energy consumption of WWTPs.

275

276 Table 1. Variables and their importance

Variable	Importance
Wastewater treatment capacity (m³/d)	0.2130
Annual load rate (%)	0.1758
COD_i /COD_e	0.1655
BOD_i /BOD_e	0.1170
Moisture content of sludge (%)	0.1134
NH₃-Ni/NH₃-N_e	0.0846
Dry weight of sludge (ton)	0.0747
Discharge standard	0.0560

277 The removal efficiency of COD and BOD had a significant impact on the energy
 278 consumption of WWTPs, which also verified that the level of removal of COD and
 279 BOD highly affected the energy consumption of WWTPs (Longo et al. 2016). This is
 280 consistent with other studies, and the pollution load is consistent with the energy
 281 consumption load of WWTPs (Torregrossa et al. 2018). COD_i/COD_e is significantly
 282 more important than BOD_i/BOD_e since the pollutants measured by BOD are subset of
 283 pollutants measured by COD, so COD contains some pollutants that do not belong to

284 BOD. Second, the model may divide the pollutants that both belong to BOD and COD
285 into the importance of BOD_i/BOD_e and COD_i/COD_e . These two factors may contribute
286 to the finding that the importance of COD_i/COD_e is significantly higher than
287 BOD_i/BOD_e .

288 However, the removal efficiency of NH_3-N , one of the primary pollutants in sewage, is
289 of relatively low importance to UEC that is because the primary function of most
290 WWTPs in China is to remove organic matter rather than denitrification, which leads
291 to the lower importance of NH_3-N_i/NH_3-N_e . Some studies have shown that COD, BOD,
292 and NH_3-N are correlated (Luo et al. 2019) and the energy to power blower fans are
293 actually the main factor in electricity consumption for the removal of COD, BOD, and
294 NH_3-N (Piotrowski and Ujazdowski 2020). Considering the fact that the smaller the
295 number of highly correlated variables, the more convenient is the practical application
296 of the model, we assigned the importance of the overlap between variables to the high
297 correlation variable, which led to the low importance of NH_3-N_i/NH_3-N_e . Depending
298 on the request to the accuracy of the model, NH_3-N_i/NH_3-N_e can be neglected during
299 practical usage, but to analyze the model more clearly and completely we will still take
300 NH_3-N_i/NH_3-N_e into consideration in the following discussion.

301 There are limits to the moisture content of sludge, so there will be energy consumption
302 to separate water from sludge. However, from Fig. 4 it appears that moisture content of
303 sludge has low correlation with other variables, so its importance will be higher. During

304 sludge conditioning, drying, and incineration, a large amount of energy is required;
 305 however, in the current statistical yearbook of WWTPs, there is no data on energy
 306 consumed for sludge disposal. Therefore, we could not further analyze the importance
 307 of sludge treatment.

308 The results in Table 1 show that compared with the data type variables, the importance
 309 of the discharge standard (Table 2) was low because BOD_i, COD_i, and NH₃-N_i are
 310 limited by discharge standard, so discharge standard is highly correlated to them. Since
 311 the current model is built to minimize the influence of this correlation, so the
 312 importance of discharge standard is low. Table 1 indicates that the discharge standard
 313 is low in importance; hence, in the following analysis, the discharge standard was not
 314 used to forecast the UEC of WWTPs.

315 Table 2. Discharge standard of COD, BOD, and NH₃-N

Parameter	Class I A	Class I B
COD _e (mg/L)	50	60
BOD _e (mg/L)	10	20
NH ₃ -N _e (mg/L)	5(8)	8(15)

316 *Note: The value in the bracket means standard at temperature ≤ 12 °C, which was*
 317 *not modeled in this study.*

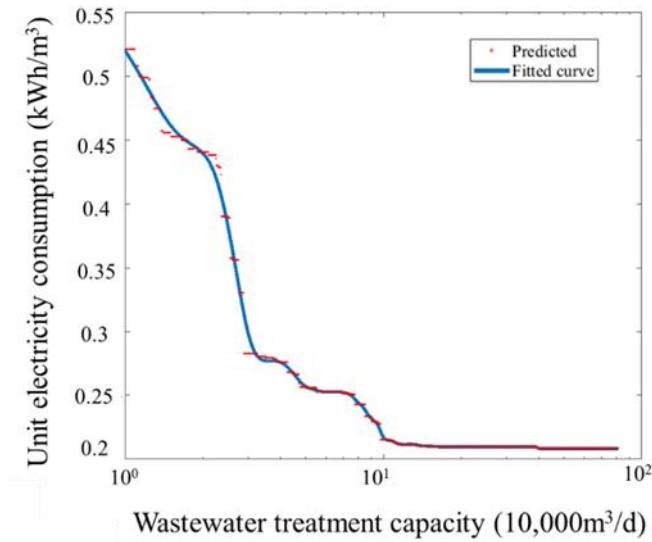
318 3.3 Model Demonstration and Prediction of Energy Consumption

319 An energy consumption model for WWTPs was established through training using a

320 large amount of data. By changing the input variables, we can predict the change in the
321 energy consumption of a WWTP. We selected the following variables with high
322 importance: design treatment capacity, annual average load rate, and removal ratios
323 (BOD_i/BOD_e , COD_i/COD_e , and NH_3-N_i/NH_3-N_e) to obtain the prediction function. The
324 model can be directly presented through this curve and it can be applied to the
325 management of energy efficiency of real WWTPs.

326 **3.3.1 Wastewater Treatment Capacity**

327 From the predictive model shown in [Fig. 5](#), it is evident that the wastewater treatment
328 capacity is negatively related to UEC, and for wastewater treatment capacities from
329 10,000 m³/d to 100,000 m³/d, the UEC decreases rapidly with an increase in the design
330 treatment capacity. Above 100,000 m³/d there is minimal decrease in UEC, which is a
331 consideration for the design of WWTPs. The overall trend is consistent with the finding
332 of previous studies (e.g., Yang et al. 2021, Huang et al. 2021). What's more, this finding
333 also follows the scale economy of WWTPs (Hernández-Chover et al. 2018).



334

335 Fig. 5 Predictive model of UEC as a function of wastewater treatment capacity with
 336 other variables constant.

337

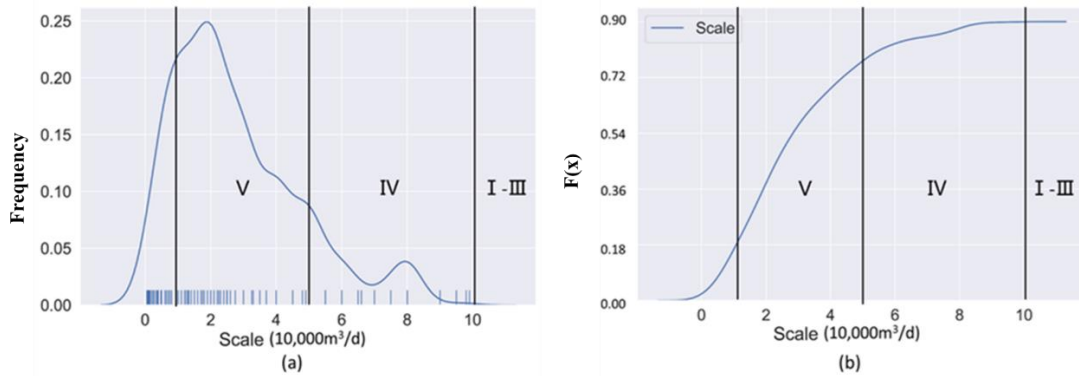
338 The construction scale of WWTPs in China can be divided into five categories
 339 (Ministry of Construction of China, 2001), as shown in Table 3:

Table 3 Standard of construction scale of WWTPs in China

Category	Construction Scale
<i>I</i>	500,000~1,000,000 m³/d
<i>II</i>	200,000~500,000 m³/d
<i>III</i>	100,000~200,000 m³/d
<i>IV</i>	50,000~100,000 m³/d
<i>V</i>	10,000~50,000 m³/d

340 From the predicted data (Fig. 5), we found that the UEC of WWTPs with a scale of I,

341 II, and III was relatively low. Therefore, we can conclude that the WWTPs larger than
 342 100,000 m³/d have effectively reduced energy consumption, and there are 245 WWTPs
 343 in this range in the database, which is 9.91% of the total WWTPs considered in this
 344 study (Fig. 6).



345
 346 Fig. 6 (a) Kernel frequency distribution (b) Probability distribution of wastewater
 347 treatment capacity of WWTPs in the model (to make the figure clearer, all the WWTPs
 348 with a wastewater treatment capacity above 100,000 m³/d were not counted in the
 349 figure).

350 The ordinate of a point on the curve in Fig.6 (a) shows the proportion of the scale that
 351 is reflected by abscissa of total WWTPs in China. The ordinate of a point on the curve
 352 in Fig.6 (b) shows the accumulated proportion of the scale smaller than abscissa of total
 353 WWTPs in China and the slope of the point shows the frequency. In Fig.6 (a) the curve
 354 peaks in Category V which means the scale of WWTPs concentrated in Category V and,
 355 after the peak, the number of WWTPs decreases with the scale in an overall trend. In
 356 Fig.6 (b), the tangent slope of the curve also shows that Category V includes most of
 357 the WWTPs in China. In short, Fig.6 shows the scale distribution of WWTPs in China

358 and we can find that most WWTPs in China are small-scale.

359

360 **3.3.2 Annual Load Rate**

361 Annual Load Rate means the percentage usage of wastewater treatment capacity over

362 the year, it reflects the divergence of design and actual usage of WWTPs. As shown in

363 [Fig. 7](#), the annual load rate has a significant impact on the UEC. The UEC remained

364 high when the annual load rate was less than 40%, but the UEC decreased significantly

365 when the annual load rate was between 40% and 100%; meanwhile, the UEC remained

366 stable in a low range after the annual load rate was more than 100%. (i.e., overload).

367 However, as overload may damage the instruments and equipment, a load rate between

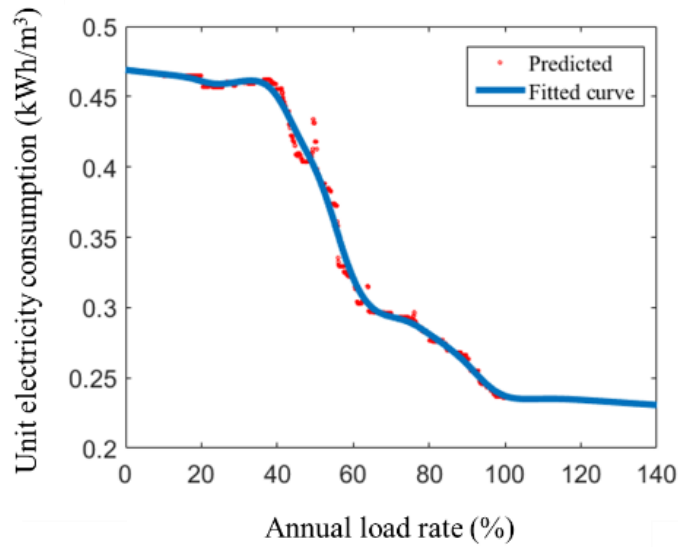
368 60% and 100% should be maintained in the design and operation of WWTPs. The trend

369 in this study corresponds well the results of Huang et al. 2021 and it also follows the

370 rule of extensive models of different factory managements (Gerami et al. 2021). This

371 finding indicates that it will be better to do more study on the amount of wastewater

372 needed to be treated in one area before designing the treatment capacity of the WWTP.



373

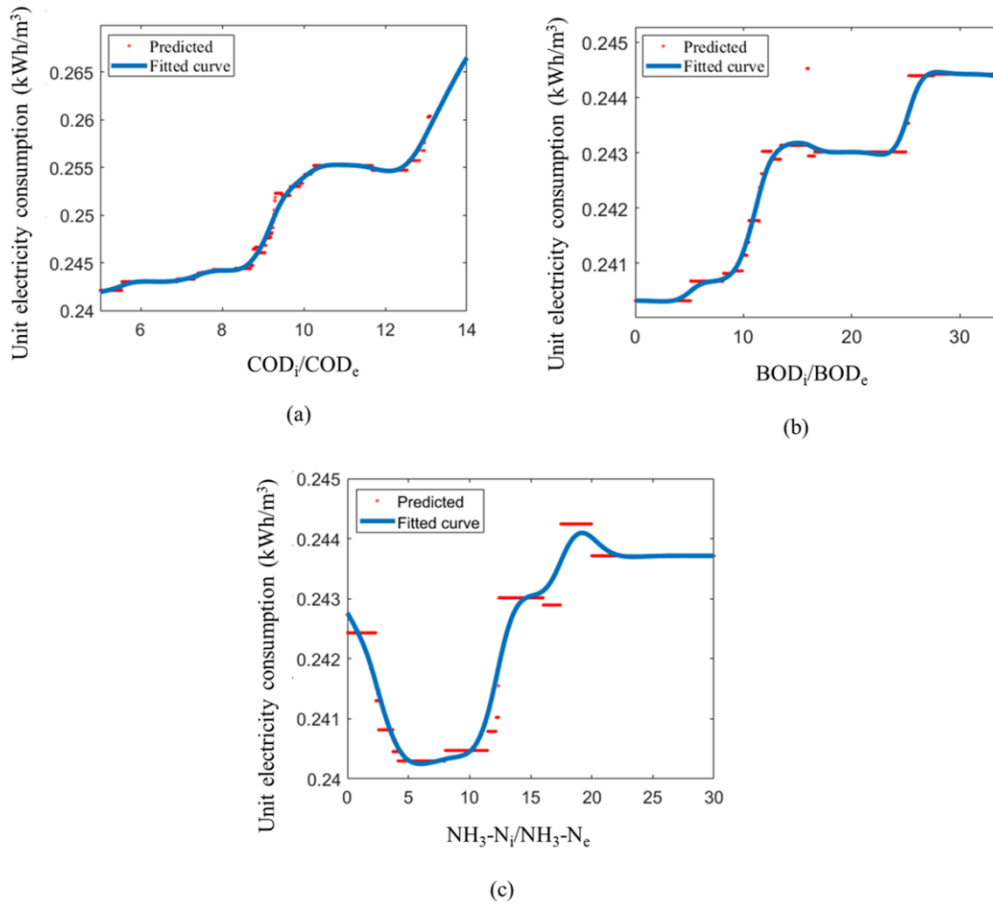
374 Fig. 7 Predictive model of UEC as a function of annual load rate with other variables
 375 constant.

376

377 3.3.3 Reduction Ratios

378 In this section, the effects of COD_i/COD_e , BOD_i/BOD_e , and NH_3-N_i/NH_3-N_e on UEC
 379 are analyzed. As shown in Table 1, the importance of the UEC of COD removal was
 380 significantly greater than that of BOD and NH_3-N , and UEC was primarily affected by
 381 COD removal. To achieve a comprehensive study, the influence of BOD removal and
 382 NH_3-N removal on UEC is also discussed herein. However, the low importance is
 383 reflected on the ordinate of BOD and NH_3-N . In general, the trends of COD and BOD
 384 is generally in line with Huang et al. 2021, while the trend of NH_3-N is a little conflict
 385 with the common sense, we are going to explain it in the following discussion.

386



387

388 Fig. 8 Predictive model of UEC as a function (a) COD_i/COD_e , (b) BOD_i/BOD_e , and

389

(c) NH_3-N_i/NH_3-N_e with other variables constant.

390

When the independent variable was too large, a flat response occurred in this region of

391

the predictive model due to the lack of data, in other words, when BOD_i/BOD_e or NH_3-

392

N_i/NH_3-N_e was too large there will be insufficient data in the database to train the model

393

at this level and the prediction will be the average of these data; therefore, these regions

394

are not discussed in the following analysis. But with the construction of huge WWTPs

395

in China, there will be more data in the future, reducing this issue.

396

From the predictive model shown in Fig. 8 (a), it is seen that COD_i/COD_e and UEC are

397

positively correlated; that is, the higher the COD reduction ratio, the higher the energy

398 consumption. As expected, a WWTP that seeks to have higher removal efficiency
399 requires more energy. In the predictive model, a relatively flat response occurs when
400 the reduction multiple is less than 9. From the available data, the average COD_i/COD_e
401 was 9.23, near the edge of the region with a minimal slope, indicating that WWTPs had
402 a high energy efficiency in the removal of COD.

403 As shown in the predictive model in Fig. 8 (b), BOD_i/BOD_e and UEC are positively
404 correlated for BOD_i/BOD_e in the range of 0–15 and 25–30, with a minimal slope
405 (except for the platform) at 0–10, and a region with a slope almost 0 in the 15–25 range.
406 From the available data, the average BOD_i/BOD_e was 18.59, in the middle of the region
407 with a slope almost 0, implying that WWTPs had a high energy efficiency in the
408 treatment of BOD, but requires further improvement.

409 As shown in Fig. 8 (c), the overall trend of the predictive model of NH_3-N_i/NH_3-N_e
410 indicates that the UEC decreases monotonically when NH_3-N_i/NH_3-N_e is lower than 5,
411 increases monotonically after NH_3-N_i/NH_3-N_e is greater than 10, and finally tends to a
412 constant value. A minimum slope is observed between 5 and 10. Therefore, the optimal
413 value of ammonia nitrogen reduction should be between 5 and 10. When NH_3-N_i/NH_3-
414 N_e is less than 10, it is negatively correlated with UEC, which is contrary to the common
415 understanding that a larger reduction multiple, leads to a higher energy consumption.
416 The specific reasons for this require further analysis. However, the possible reasons are
417 as follows: 1. As the importance of NH_3-N_i/NH_3-N_e in UEC is low, which causes the

418 difference between the maximum and minimum values of the final prediction result to
 419 be ≤ 0.04 kWh, the measuring instrument may not be highly accurate. 2. When $\text{NH}_3\text{-}$
 420 $\text{N}_i/\text{NH}_3\text{-N}_e$ is 5–10, the reduction multiple is easily achieved. For a lower value, energy
 421 consumption may be required to limit the reduction multiple. This study aimed to
 422 predict the UEC of WWTPs if the plant upgrades to a higher standard, thus improving
 423 the removal ratios ($\text{COD}_i/\text{COD}_e$, $\text{BOD}_i/\text{BOD}_e$, and $\text{NH}_3\text{-N}_i/\text{NH}_3\text{-N}_e$). The number of
 424 WWTPs applied Class I A were 1,041 and 1,184 in Class I B, with only 162 in Class II.
 425 Therefore, we primarily considered the improvement of the discharge standard from
 426 Class I B to Class I A. The specific discharge standards are listed in [Table 2](#).

427 The ratios $\text{COD}_i/\text{COD}_e$, $\text{BOD}_i/\text{BOD}_e$, and $\text{NH}_3\text{-N}_i/\text{NH}_3\text{-N}_e$ were used as variables in
 428 the model, since the discharge standard restricts COD_e , BOD_e , $\text{NH}_3\text{-N}_e$. Therefore, the
 429 following modifications were adopted in this study:

$$430 \quad y_{\text{COD}} = \frac{\text{COD}_i/\text{COD}_e}{x_{\text{COD}}} \times \text{COD}_e \quad (8)$$

$$431 \quad y_{\text{BOD}} = \frac{\text{BOD}_i/\text{BOD}_e}{x_{\text{BOD}}} \times \text{BOD}_e \quad (9)$$

$$432 \quad y_{\text{N}} = \frac{\text{NH}_3\text{-N}_i/\text{NH}_3\text{-N}_e}{x_{\text{N}}} \times \text{NH}_3 - \text{N}_e \quad (10)$$

433 where y_{COD} : $\text{COD}_I/\text{COD}_E$ after upgrading to a higher class, x_{COD} : COD_E of Class I
 434 A, y_{BOD} : $\text{BOD}_I/\text{BOD}_E$ after upgrading to a higher class, x_{BOD} : BOD_E of Class I A,
 435 y_{N} : $(\text{NH}_3 - \text{N}_I/\text{NH}_3 - \text{N}_E)$ after upgrading to a higher class, and x_{N} : $(\text{NH}_3 - \text{N}_E)$ of
 436 Class I A.

437 From equation 4-6, a new value of $\text{COD}_i/\text{COD}_e$, $\text{BOD}_i/\text{BOD}_e$, $\text{NH}_3 - \text{N}_i/\text{NH}_3 - \text{N}_e$

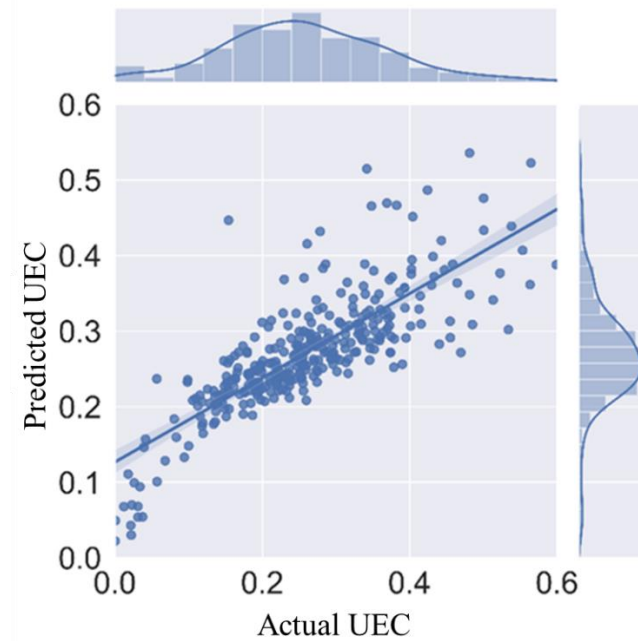
438 after the improvement of the discharge standard from Class I B to Class I A was
439 calculated. The model predicted the UEC of WWTPs with the new data.

440 The results showed that when the discharge quality of WWTPs was upgraded from
441 Class I B to Class I A, the increase in the UEC of WWTPs varied due to the various
442 effluent qualities. The UEC of WWTPs had an average increase of 17%, obtained from
443 Equations 4-6.

444

445 **3.3.4 Model Validation**

446 As previously mentioned, the R^2 of the model was 0.702. $MSE = 0.00662 \text{ (kWh/m}^3\text{)}^2$,
447 $MAPE = 5.74\%$, $RSME = 0.106 \text{ kWh/m}^3$, $MAE = 0.0416 \text{ kWh/m}^3$, $MedAE = 0.0416$
448 kWh/m^3 and $MSLE = 0.00327$ (obtained from equations (1) to (6)), which are very low
449 (Weber et al. 2020). These low evaluation metrics indicate that the model for UEC of
450 WWTPs developed in this study was quite accurate. As shown in [Fig. 9](#) and [Fig. S2](#),
451 the actual and predicted UEC exhibit the same trend when the UEC is not too high or
452 too low. The predicted UEC was not accurate when the corresponding actual value was
453 too high or too low because of insufficient data for fully developing the model. In fact,
454 in practice, there are not many cases of too large or too small WWTPs, so the effect of
455 these WWTPs are not significant.



456

457 Fig. 9 QQPlot and simplified kernel frequency distribution representing actual and
 458 predicted UEC for 347 WWTPs

459 3.4 Comparison with other approaches

460 Compared to the Data Envelopment Analysis, Random Forest is more stable when some
 461 input data is missing, which means a unified model can be made without considering
 462 special cases that one or more variables are missing. Considering the fact that it's hard
 463 to set up a monitoring system that would include thousands of WWTPs across China
 464 with exactly the same variables, it would be impossible to build a normalized model
 465 using Data Envelopment Analysis.

466 As mentioned in Section 3.1, the WWTP variables are correlated, so the accuracy of a
 467 Multiple Linear Regression model compared to a Random Forest model is relatively
 468 low. In this study, Multiple Linear Regression model was considered, but the R^2 (0.147)

469 was too low. In comparison, the random forest model can achieve a much higher R^2
470 (0.702). Therefore, Random Forest is more suitable to build the model than Data
471 Envelopment Analysis or Multiple Linear Regression.

472

473 **4. Conclusion**

474 In this study, an energy consumption model for WWTPs was developed using machine
475 learning. The UEC of a WWTP can be predicted with a few key parameters by the
476 model using the random forest algorithm. It can also predict the UEC of a WWTP for
477 policy formulation and improvement of sewage treatment standards. This model can be
478 a useful tool for investigating the water-energy nexus in WWTPs. Although the
479 particular model in this study is based on data from Chinese WWTPs, it can be easily
480 applied to WWTPs worldwide by changing the input data. In this study, we didn't
481 investigate the influence of local climate and treatment technologies due to insufficient
482 data, which are also very important and deserve further research in the future.

483

484 **References**

485 Bagherzadeh, F., M. J. Mehrani, M. Basirifard and J. Roostaei (2021). "Comparative study on total
486 nitrogen prediction in wastewater treatment plant and effect of various feature selection
487 methods on machine learning algorithms performance." Journal of Water Process
488 Engineering **41**: 102033.
489 Behera, C. R., R. Al, K. V. Gernaey and G. Sin (2020). "A process synthesis tool for WWTP – An
490 application to design sustainable energy recovery facilities." Chemical Engineering Research
491 and Design **156**: 353-370.
492 Belgiu, M. and L. Dragut (2016). "Random forest in remote sensing: A review of applications and
493 future directions." Isprs Journal of Photogrammetry and Remote Sensing **114**: 24-31.

494 Breiman, L. (2001). "Random Forests." Machine Learning **45**(1): 5-32.

495 Chen, W., S. Zhang, R. W. Li and H. Shahabi (2018). "Performance evaluation of the GIS-based data
496 mining techniques of best-first decision tree, random forest, and naive Bayes tree for landslide
497 susceptibility modeling." Science of the Total Environment **644**: 1006-1018.

498 Cheng, L., J. De Vos, P. Zhao, M. Yang and F. Witlox (2020). "Examining non-linear built
499 environment effects on elderly's walking: A random forest approach." Transportation
500 Research Part D: Transport and Environment **88**: 102-552.

501 Construction Standard of Urban Sewage Treatment Project: Ministry of Construction of the
502 People's Republic of China; 2001.

503 Daw, J., K. Hallett, J. Dewolfe and I. Venner (2012). "Energy Efficiency Strategies for Municipal
504 Wastewater Treatment Facilities." Office of Scientific & Technical Information Technical
505 Reports.

506 Duro, D. C., S. E. Franklin and M. G. Dube (2012). "A comparison of pixel-based and object-based
507 image analysis with selected machine learning algorithms for the classification of agricultural
508 landscapes using SPOT-5 HRG imagery." Remote Sensing of Environment **118**: 259-272.

509 E. Açıkkalp ve S. Yerel Kandemir (2018). "Optimum insulation thickness of the piping system with
510 combined economic and environmental method." Energy Sources Part A-Recovery Utilization
511 and Environmental Effects, **40** (23): 2876-2885.

512 EPRI, Electricity use and management in the municipal water supply and wastewater industries,
513 2013, 5-16.

514 Gazendam, E., B. Gharabaghi, J. D. Ackerman and H. Whiteley (2016). "Integrative neural networks
515 models for stream assessment in restoration projects." Journal of Hydrology: 339-350.

516 Habib, R. Z., T. Thiemann and R. A. Kendi (2020). "Microplastics and Wastewater Treatment
517 Plants—A Review." Journal of Water Resource and Protection **12**(1): 1-35.

518 Hernandez-del-Olmo, F., E. Gaudioso, N. Duro and R. Dormido (2019). "Machine Learning
519 Weather Soft-Sensor for Advanced Control of Wastewater Treatment Plants." Sensors **19**(14).

520 Heslot, N., D. Akdemir, M. E. Sorrells and J. L. Jannink (2014). "Integrating environmental covariates
521 and crop modeling into the genomic selection framework to predict genotype by
522 environment interactions." Theoretical and Applied Genetics **127**(2): 463-480.

523 Hernández-Sancho, F. and R. Sala-Garrido (2009). "Technical efficiency and cost analysis in
524 wastewater treatment processes: A DEA approach." Desalination **249**(1): 230-234.

525 Li, W., X. Shi, S. Zhang and G. Qi (2019). "Modelling of ammonia recovery from wastewater by air
526 stripping in rotating packed beds." Science of The Total Environment **702**: 134-971.

527 Longo, S., B. M. d'Antoni, M. Bongards, A. Chaparro, A. Cronrath, F. Fatone, J. M. Lema, M.
528 Mauricio-Iglesias, A. Soares and A. Hospido (2016). "Monitoring and diagnosis of energy
529 consumption in wastewater treatment plants. A state of the art and proposals for
530 improvement." Applied Energy **179**: 1251-1268.

531 Loozen, Y., K. T. Rebel, S. M. de Jong, M. Lu, S. V. Ollinger, M. J. Wassen and D. Karssenberg (2020).
532 "Mapping canopy nitrogen in European forests using remote sensing and environmental
533 variables with the random forests method." Remote Sensing of Environment **247**: 111933.

534 Luo, L., M. Dzakpasu, B. Yang, W. Zhang, Y. Yang and X. C. Wang (2019). "A novel index of total

535 oxygen demand for the comprehensive evaluation of energy consumption for urban
536 wastewater treatment." Applied Energy **236**: 253-261.

537 Ministry of Construction of China (2001). Construction standard of urban sewage treatment project.
538 No.77.

539 Mizuta, K. and M. Shimada (2010). "Benchmarking energy consumption in municipal wastewater
540 treatment plants in Japan." Water Science & Technology A Journal of the International
541 Association on Water Pollution Research **62**(10): 2256-2262.

542 Mjalli, F. S., S. Al-Asheh and H. E. Alfadala (2007). "Use of artificial neural network black-box
543 modeling for the prediction of wastewater treatment plants performance." Journal of
544 Environmental Management **83**(3): 329-338.

545 Molinos-Senante, M., F. Hernandez-Sancho and R. Sala-Garrido (2014). "Benchmarking in
546 wastewater treatment plants: a tool to save operational costs." Clean Technologies and
547 Environmental Policy **16**(1): 149-161.

548 Pamulaparty, L., C. V. G. Rao and M. S. Rao (2017). "Critical review of various near-duplicate
549 detection methods in web crawl and their prospective application in drug discovery."
550 International Journal of Biomedical Engineering and Technology **25**(2-4): 212-226.

551 Paul, G. C., S. Saha and K. G. Ghosh (2020). "Assessing the soil quality of Bansloi river basin, eastern
552 India using soil-quality indices (SQIs) and Random Forest machine learning technique."
553 Ecological Indicators **118**: 106-804.

554 Perez, V. M., J. M. M. Fernandez, J. V. Balsera and C. A. Alvarez (2021). "A Random Forest Model
555 for the Prediction of FOG Content in Inlet Wastewater from Urban WWTPs." Water **13**(9): 1237.

556 Piotrowski, R. and T. Ujazdowski (2020). "Designing Control Strategies of Aeration System in
557 Biological WWTP." Energies **13**(14): 1-17.

558 Sabia, G., L. Petta, F. Avolio and E. Caporossi (2020). "Energy saving in wastewater treatment plants:
559 A methodology based on common key performance indicators for the evaluation of plant
560 energy performance, classification and benchmarking." Energy Conversion and Management
561 **220**: 113-067.

562 Scott, C. A., S. A. Pierce, M. J. Pasqualetti, A. L. Jones, B. E. Montz and J. H. Hoover (2011). "Policy
563 and institutional dimensions of the water–energy nexus." Energy Policy **39**(10): 6622-6630.

564 Siddiqi, A. and L. D. Anadon (2011). "The water–energy nexus in Middle East and North Africa."
565 Energy Policy **39**(8): 4529-4540.

566 Smith, K., S. Guo, Q. Zhu, X. Dong and S. Liu (2019). "An evaluation of the environmental benefit
567 and energy footprint of China's stricter wastewater standards: Can benefit be increased?"
568 Journal of Cleaner Production **219**: 723-733.

569 Song, R. S., A. A. Keller and S. Suh (2017). "Rapid Life-Cycle Impact Screening Using Artificial Neural
570 Networks." Environmental Science & Technology **51**(18): 10777-10785.

571 Tamjidi Farahbakhsh, M. and M. Chahartaghi (2020). "Performance analysis and economic
572 assessment of a combined cooling heating and power (CCHP) system in wastewater
573 treatment plants (WWTPs)." Energy Conversion and Management **224**: 113-351.

574 Tian, H., P. Bai, Y. Tan, Z. Li, D. Peng, X. Xiao, H. Zhao, Y. Zhou, W. Liang and L. Zhang (2020). "A
575 new method to detect methylation profiles for forensic body fluid identification combining

576 ARMS-PCR technique and random forest model." Forensic Science International: Genetics **49**:
577 102-371.

578 Torregrossa, D., U. Leopold, F. Hernandez-Sancho and J. Hansen (2018). "Machine learning for
579 energy cost modelling in wastewater treatment plants." Journal of Environmental
580 Management **223**(OCT.1): 1061-1067.

581 Torregrossa, D., G. Schutz, A. Cornelissen, F. Hernández-Sancho and J. Hansen (2016). "Energy
582 saving in WWTP: Daily benchmarking under uncertainty and data availability limitations."
583 Environmental Research **148**: 330-337.

584 Trenouth, W. R., B. Gharabaghi and H. Farghaly (2018). "Enhanced roadside drainage system for
585 environmentally sensitive areas." Science of The Total Environment **610-611**: 613-622.

586 Wang, H. Z., Z. X. Lei, X. Zhang, B. Zhou and J. C. Peng (2019). "A review of deep learning for
587 renewable energy forecasting." Energy Conversion and Management **198**: 111799.

588 Weber, V. A. M., F. d. L. Weber, A. d. S. Oliveira, G. Astolfi, G. V. Menezes, J. V. de Andrade Porto,
589 F. P. C. Rezende, P. H. d. Moraes, E. T. Matsubara, R. G. Mateus, T. L. A. C. de Araújo, L. O. C.
590 da Silva, E. Q. A. de Queiroz, U. G. P. de Abreu, R. da Costa Gomes and H. Pistori (2020).
591 "Cattle weight estimation using active contour models and regression trees Bagging."
592 Computers and Electronics in Agriculture **179**: 105-804.

593 Wei, J., W. Huang, Z. Q. Li, W. H. Xue, Y. R. Peng, L. Sun and M. Cribb (2019). "Estimating 1-km-
594 resolution PM2.5 concentrations across China using the space-time random forest approach."
595 Remote Sensing of Environment **231**: 111221.

596 Wylie, B. K., N. J. Pastick, J. J. Picotte and C. A. Deering (2019). "Geospatial data mining for digital
597 raster mapping." Giscience & Remote Sensing **56**(3): 406-429.

598 Xing, L., M. G. Xue and M. S. Hu (2019). "Dynamic simulation and assessment of the coupling
599 coordination degree of the economy-resource-environment system: Case of Wuhan City in
600 China." Journal of Environmental Management **230**: 474-487.

601 Yeşilkanat, C. M. (2020). "Spatio-temporal estimation of the daily cases of COVID-19 in worldwide
602 using random forest machine learning algorithm." Chaos, Solitons & Fractals **140**: 110-210.

603 Zhu, Y. A., W. Y. Xu, G. L. Luo, H. L. Wang, J. J. Yang and W. Lu (2020). "Random Forest enhancement
604 using improved Artificial Fish Swarm for the medial knee contact force prediction." Artificial
605 Intelligence in Medicine **103**: 101811.

606 Yang, J., B. Chen (2021). "Energy efficiency evaluation of wastewater treatment plants (WWTPs)
607 based on data envelopment analysis" Applied Energy **289**: 116680

608 Huang, R., Z. Shen, H. Wang, J. Xu, Z. Ai, H. Zheng, R. Liu (2021). "Evaluating the energy efficiency
609 of wastewater treatment plants in the Yangtze River Delta: Perspectives on regional
610 discrepancies" Applied Energy **297**: 117087

611 Xu, J., P. Luo, B. Lu, H. Wang, X. Wang, J. Wu, J. Yan (2018). "Energy-water nexus analysis of
612 wastewater treatment plants (WWTPs) in China based on statistical methodologies" Energy
613 Procedia **152**: 259-264.

614 Hernández-Chover, V., Á. Bellver-Domingo and F. Hernández-Sancho (2018). "Efficiency of
615 wastewater treatment facilities: The influence of scale economies." Journal of Environmental
616 Management **228**: 77-84.

617 Gerami, N., A. Ghasemi, A. Lotfi, L. G. Kaigutha and M. Marzband (2021). "Energy consumption
618 modeling of production process for industrial factories in a day ahead scheduling with
619 demand response." Sustainable Energy, Grids and Networks **25**: 100420.

620 Zhong, S., K. Zhang, M. Bagheri, J. G. Burken, A. Gu, B. Li, X. Ma, B. L. Marrone, Z. J. Ren, J. Schrier,
621 W. Shi, H. Tan, T. Wang, X. Wang, B. M. Wong, X. Xiao, X. Yu, J.-J. Zhu and H. Zhang (2021).
622 "Machine Learning: New Ideas and Tools in Environmental Science and Engineering."
623 Environmental Science & Technology **55**(19): 12741-12754.

624 Gupta, S., D. Aqa, A. Pruden, L. Zhang and P. Vikesland (2021). "Data Analytics for Environmental
625 Science and Engineering Research." Environmental Science & Technology **55**(16): 10895-
626 10907.