

Journal Pre-proof

Development of a wide-range soft sensor for predicting wastewater BOD₅ using an eXtreme gradient boosting (XGBoost) machine

P.M.L. Ching, X. Zou, Di Wu, R.H.Y. So, G.H. Chen



PII: S0013-9351(22)00280-8

DOI: <https://doi.org/10.1016/j.envres.2022.112953>

Reference: YENRS 112953

To appear in: *Environmental Research*

Received Date: 29 October 2021

Revised Date: 6 February 2022

Accepted Date: 10 February 2022

Please cite this article as: Ching, P.M.L., Zou, X., Wu, D., So, R.H.Y., Chen, G.H., Development of a wide-range soft sensor for predicting wastewater BOD₅ using an eXtreme gradient boosting (XGBoost) machine, *Environmental Research* (2022), doi: <https://doi.org/10.1016/j.envres.2022.112953>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022 Published by Elsevier Inc.

1 **Development of a wide-range soft sensor for predicting wastewater**2 **BOD₅ using an eXtreme gradient boosting (XGBoost) machine**

P.M.L. Ching ^a, X. Zou ^b, Di Wu ^{b,c,d*}, R. H.Y. So ^e, G.H. Chen ^b

a. Bioengineering Graduate Program, Department of Chemical and Biological Engineering, The Hong Kong University of Science and Technology, Hong Kong SAR, China.

b. Department of Civil and Environmental Engineering, The Hong Kong University of Science and Technology, Hong Kong SAR, China.

3 c. Center for Environmental and Energy Research, Ghent University Global
4 Campus, Republic of Korea.

5 d. Department of Green Chemistry and Technology, Ghent University, Belgium.

6 e. Department of Industrial Engineering and Decision Analytics, The Hong Kong
7 University of Science and Technology, Hong Kong SAR, China.

8

9 * Corresponding author:

10 Di WU, Center for Environmental and Energy Research, , Ghent University Global
11 Campus, Republic of Korea. (Email: di.wu@ghent.ac.kr)

12

13

Abstract

14 In wastewater monitoring, detecting extremely high pollutant concentrations is
15 necessary to properly calibrate the treatment process. However, existing hardware
16 sensors have a limited linear range which may fail to measure extremely high levels
17 of pollutants; and likewise, the conventional “soft” model sensors are not suitable for
18 the highly-skewed data distributions either. This study developed a new soft sensor by
19 using eXtreme Gradient Boosting (XGBoost) machine learning to ‘measure’ the
20 wastewater organics (in terms of 5-day biochemical oxygen demand (BOD₅)). The
21 soft sensor was tested on influent and effluent BOD₅ of two different wastewater
22 treatment plants to validate the results. The model results showed that XGBoost can
23 detect these extreme values better than conventional soft sensors. This new soft sensor
24 can function using a sparse input matrix via XGBoost’s sparsity awareness algorithm
25 - which can address the limitation of the conventional soft sensor with the fallibility of
26 supporting hardware sensors even.

27 **Keywords:** Soft sensor, machine learning, XGBoost, real-time monitoring,
28 biochemical oxygen demand (BOD)

1. Introduction

29 Online monitoring is an important prerequisite for advancements in wastewater
30 treatment. Real-time information allows the plant to implement more cost-efficiently
31 and gives evidence that the quality regulations are consistently being met. A
32 conventional online monitoring system for relevant wastewater parameters (e.g.
33 chemical oxygen demand (COD), ammonia concentration) would be composed of
34 hardware sensors. However, the existing hardware sensors for these parameters have a
35 limited useful lifespan due to the harsh conditions of wastewater. The accumulation of
36 sludge and precipitates on the sensor lowers its accuracy over time, and necessitates
37 frequent maintenance (Haimi et al., 2013). The sensor itself loses its functionality
38 over time, such as the dissolution of Ag/AgCl layers observed in electrode-based
39 sensors (Hill et al., 2020) or the degradation of the microorganism culture used in
40 biosensors (Raud et al., 2012). Besides, there is no mature sensor product for
41 measuring five-day biochemical oxygen demand (BOD₅) to reflect the biodegradable
42 organics content in the wastewater. To solve this problem, one good option is to use a
43 machine learning-based soft sensor model, which estimates parameter values from
44 other hardware sensors using machine learning. In this way, soft sensors can facilitate
45 real-time monitoring by avoiding the delays or missing data resulting from frequent
46 maintenance; manual measurement. However, the accuracy of the soft sensor is still
47 dependent on the (1) choice of hardware sensors used as its basis for estimation; (2)
48 the volume and range of data, and (3) the appropriateness of the machine learning
49 model used in estimation.

50 In choosing the hardware basis for the machine learning-based soft sensor, the
51 ideal choice is to choose simple and stable sensors (e.g. pH, conductivity). Yet, the

52 majority of soft-sensor studies add complex sensors (e.g. chemical oxygen demand
53 (COD), NH_4) to enhance accuracy. When there is a large quantity of potential soft
54 sensors, parameter selection techniques can be employed to reduce the number of
55 model inputs. These techniques aim to identify the input parameters sharing the
56 strongest relationship with the output parameter(s) (Zhu et al., 2017). In addition, the
57 performance of the soft sensor may be improved by the removal of some inputs, as
58 collinearity between the input variables may promote overfitting (Asante-Okyere et
59 al., 2020).

60 It should note that datasets used in soft sensor development vary in size (Ye et
61 al., 2020). While there is no defined minimum for the size of the dataset, a larger
62 dataset is preferred for higher generalizability. The volume and range of wastewater
63 datasets are limited by sensor degradation and infrequent sampling, resulting in
64 missing sensor readings in the dataset. These missing values can be filled in using a
65 statistic (e.g. mean, median), or using a statistical method to impute the missing
66 values (Wu et al., 2008). While these methods can produce additional samples to the
67 dataset, samples with missing parameters may increase the uncertainty in the model,
68 and skew the estimations of the soft sensor (Li et al., 2020).

69 Although any mathematical model can be applied in soft sensor development,
70 machine learning approaches are preferred in recent studies. One reason is that these
71 utilize the existing wastewater treatment databases, and produce new insights without
72 additional experimentation (Asami et al., 2021; Qiu et al., 2021). Using machine
73 learning, mathematical relationships are automatically ‘learned’ instead of manually
74 developed based on theoretical knowledge, and this may be more efficient in some
75 cases. Some examples include applications in predicting the concentration of novel
76 pollutants and pathogens of interest (Abdeldayem et al., 2022). It can also capture a

77 broad range of operating conditions, whereas traditional mechanistic modelling is
78 typically limited to steady state analysis (Wang et al., 2021).

79 Currently, the most popular machine learning models applied in wastewater
80 treatment are artificial neural networks (ANN) and support vector machines (SVM)
81 (Ye et al., 2020). The ANN model is composed of several layers of node equations,
82 which form a highly nonlinear relationship. Its primary advantage is its ability to
83 present complex underlying relationships between variables, and has improved the
84 accuracy of predicting several key wastewater parameters (Matheri et al., 2021).
85 However, the disadvantage of this complex nonlinear structure is that ANN models
86 have a tendency to overfit to the dataset used for training, and thus require a large
87 number of samples in order for the trained model to be generalizable (Ye et al., 2020).
88 Some modifications of the classical neural network have been proposed: Zhu et al.
89 (2017) integrated the radial basis function in an ANN model for predicting total
90 phosphorus (TP), as this function is associated with enhanced generalizability even
91 with smaller datasets. Cong and Yu (2018) used wavelet transforms in an ANN model,
92 to prevent it from overfitting to noise in the training set.

93 On the other hand, the advantage of SVM is its generalizability. Specifically,
94 the objective function used in determining the optimal parameters of an SVM model
95 seeks to maximize generalizability (Liu & Xie, 2020; Jiang et al., 2020). Because of
96 this, SVM can be used even with relatively small datasets, which can be important
97 when analysing novel processes and technologies (Hosseinzadeh et al., 2022; Moufid
98 et al., 2021). The disadvantage of SVM is that its generalizability objective may lead
99 the model to overfit to the dominant condition in the dataset (Jaramillo et al., 2018). A
100 soft sensor based on SVM may thus fail in accurately measuring extreme values in the
101 statistical distribution of a parameter, or in differentiating between normal and
102 abnormal operating conditions.

103 It should also be noted that, aside from the recurring problems in terms of
104 missing sensor readings and noise, data on water treatment is characterized by skewed
105 and non-normal distributions. This may render approaches that emphasize
106 generalizability unsuitable for modeling. Ensemble models are a non-parametric
107 modeling approach that makes estimations using the average of a large number of
108 simple models (Sharafati et al., 2020). Each model within the ensemble may represent
109 a characteristic of the distribution of the predicted parameter. This enhances the
110 robustness of the model while allowing it to model non-normal variables.

111 In this study, extreme gradient boosting (XGBoost), a new ensemble method,
112 is proposed in soft sensor development for BOD₅ analysis. This method was selected
113 because of its robustness and ability to model non-normal variables. In addition,
114 XGBoost includes a sparsity-awareness algorithm that allows it to train using samples
115 with missing sensor readings. Operating as a soft sensor, XGBoost can also make
116 inferences from inputs with missing parameters, which is faster compared to using a
117 separate model to estimate the missing values. This study used two case studies of
118 wastewater treatment plants to identify the dataset characteristics. Finally, the
119 comparisons with other popular machine learning techniques were drawn to verify the
120 merits of XGBoost machine learning.

121

2. Materials and Methods

122 The framework of developing a new machine learning-based soft sensor is illustrated
123 in Figure 1. The details were described as 1) the data source (two case studies for
124 BOD₅ soft sensor development); 2) the general steps involved in the development of
125 the soft sensor; 3) the method of developing Modified Partial Least Squares used in
126 selecting supporting sensors for the soft sensor; 4) the methods for missing sensor

127 reading in the dataset; 5) the development approach for the proposed XGBoost soft
128 sensor and other potential soft sensor development methods for comparison.

2.1. Data Source

129 The proposed soft sensor development approach was demonstrated through two case
130 studies: Case 1, the public wastewater treatment dataset published by the UCI
131 Machine Learning Repository (Dua & Graff, 2019); and Case 2, a dataset collected
132 from a wastewater treatment plant in Hong Kong (see supplementary information
133 Figure S1). The data used in this study came from manual measurements to allow
134 easier comparability of results, avoid variance resulting from the choice of sensor
135 and sensor performance. Thus, we can assume that all model input data is accurate.
136 Although in the context of real operation, input data collected from sensors would
137 suffer from noise and interference, there is already existing work on mathematical
138 models that address these problems (see Ba-Alawi et al., 2021; Fan et al., 2020; Wang
139 et al., 2020).

140 The case study based on the dataset of the UCI Machine Learning Repository
141 (Case 1) describes the treatment of urban wastewater in an unnamed plant. It contains
142 527 daily readings, with missing data found in 84 samples in the influent and 72
143 samples in the effluent. The Hong Kong dataset (Case-2) was collected from January
144 2013 to December 2018. It contains 2,189 daily samples, with missing data in 1,576
145 samples in the influent and 1,575 samples in the effluent. The supplementary
146 information for this study contains more details on the statistical properties of this
147 dataset, namely its range, skewness, and the number of missing readings for each
148 parameter in the dataset (Table S1 and Table S2). Generally speaking, both datasets
149 are highly skewed, with a higher level of skewness among effluent parameters.
150 Skewness measures the tendency of samples in the dataset to cluster towards lower
151 (positive skew) or higher values (negative skew). High levels of skewness indicate

152 that the dataset is not normally distributed, which is a key assumption in most data-
153 driven models. It is also notable that the distribution of effluent BOD (BOD_{eff}) is more
154 skewed compared to influent BOD (BOD_{inf}), while BOD_{inf} has a higher variance
155 compared to BOD_{eff} .

2.2. General Soft Sensor Development

156 This study developed soft sensors for BOD_{inf} and effluent BOD_{eff} for the two cases
157 described in the previous section. For Case 1 (using data from the UCI repository),
158 BOD_{inf} was modeled using other influent parameters as supporting sensors; and
159 likewise, BOD_{eff} was modeled only using other effluent parameters as supporting
160 sensors. Case 2 differs slightly as it includes ambient temperature (represented by
161 temperature measured at the reactor, $Temp_{Reac}$) as a potential supporting sensor. This
162 was included as a supporting sensor for BOD_{inf} , representing the potential for organic
163 degradability before the treatment process.

164 There are multiple potential supporting sensors, and some information is
165 redundant across the different hardware sensors (e.g., NH_3-N and NO_2-N). To identify
166 the best-supporting sensors to use as the basis for the soft sensor, the study
167 incrementally added supporting sensors as inputs to the soft sensor model and
168 evaluated the change in performance as the output. The order of adding supporting
169 sensors to the model was based on a modified Partial Least Squares approach (the
170 details see next section). Limiting the number of input variables also limits the
171 potential for the soft sensor to fail; its inputs are other supporting sensors, therefore its
172 performance is dependent on its supporting sensors.

173 Given that a significant portion of the samples in both cases include missing
174 parameters, the study considered three methods of handling the missing values: (i)
175 removing the samples with missing values; (ii) using k -nearest neighbors (kNN) to fill
176 in the missing values; and (iii) using the sparsity-awareness algorithm of XGBoost to

177 train a model using samples with missing parameters. The disadvantage of removing
178 the samples with missing values is that it significantly reduces the size of the dataset.
179 Depending on the distribution and noisiness of the data, a smaller dataset could
180 prevent the machine learning models from representing the complete and general
181 behavior of BOD₅. Conversely, using a model to impute the missing values could also
182 worsen soft sensor performance through the errors in the imputed values. For most
183 estimation models (e.g. ANN and SVM), it is necessary to have a separate method
184 such as *k*-nearest neighbors for handling the missing values. But XGBoost differs
185 from these methods as it has a built-in algorithm to incorporate the samples with
186 missing values in the training process. This is one of the key advantages of XGBoost
187 and will be described further in the following section. It will be compared with the
188 aforementioned methods (i) and (ii) of handling missing values.

189 Performance analysis was based on root mean square error (RMSE, see Eq. 1)
190 in units of mg/L. This reflects the actual deviation of the soft sensor reading from the
191 ‘real’ value, based on laboratory tests. It also reflects the effect of differences in
192 dataset characteristics such as the minimum, maximum and kurtosis on the magnitude
193 of error.

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}} \quad \text{Eq. 1}$$

194 The model results were validated using 10-fold cross-validation. This
195 approach divides the training set into 10 sets with no overlap. Each of the ten sets
196 represents a test set for the model, where it will be trained using all the other samples
197 not included in the test set (as shown in Figure 1b). The purpose of this method is to
198 determine the general performance of each method using different datasets. This also
199 allows for a comparison of the consistency of model performance.

2.3. Modified Partial Least Squares for Supporting Sensor Selection

200 PLS is a form of linear regression that maximizes the covariance between model
201 inputs and the predicted output. It has been used in related studies for soft sensor
202 development because it simultaneously maximizes the variance in the inputs, and the
203 correlation between the model's inputs and outputs (Zhu et al., 2017; Qin et al.,
204 2012). This means that the strongest supporting sensors with the least redundancy will
205 be selected as inputs. Although there are various ways of interpreting the results of
206 PLS for input variable selection, one of the most reliable and straightforward ways is
207 by measuring the absolute value of the PLS regression coefficients (Mehmood et al.,
208 2020). The greater the value of the regression coefficient, the more significant its
209 corresponding input variable is based on PLS regression.

210 However, in the context of wastewater treatment, the effectiveness of the
211 supporting sensor has to be weighed with respect to the practicality of selecting this
212 particular soft sensor. Simpler sensors (i.e., pH, conductivity, temperature, flow rate)
213 may be easier to maintain or replace. Using these sensors as supporting sensors would
214 make the proposed soft sensor more reliable, although these variables may not have
215 the strongest correlation with BOD₅. This study applied a modified PLS approach in
216 selecting the supporting sensors. Several versions of the soft sensor were built using
217 different sets of supporting sensors as inputs to the model. There lessen the number of
218 combinations that would have to be tested, the study used the modified PLS approach
219 to guide the selection process. This approach prioritizes the simpler sensors as inputs
220 for the initial model. Then, sensors are added incrementally in the order of their PLS
221 regression coefficients. The optimal soft sensor design was selected based on the
222 model which resulted in the lowest and most consistent RMSE.

2.4. Methods for Missing Sensor Readings in the Dataset

223 In general, having a larger dataset is preferred as it should help enhance the
224 generalizability of the model. Several studies have attempted to fill in missing values
225 in a dataset to enhance model performance. Among these studies, *k*-Nearest neighbors
226 (kNN) has emerged as a standard for determining the missing values. To fill in the
227 missing parameters of a sample, kNN uses the weighted average of samples with the
228 highest similarity based on the available parameters for that sample (Qi et al., 2021).
229 In this case, the similarity is based on a distance measure such as Euclidean distance
230 (Alfeilat et al., 2019).

231 XGBoost has its algorithm for addressing the missing values. This algorithm
232 is known as a sparsity awareness split-finding algorithm, referring to the dataset
233 splitting involved in determining the optimal structure for the XGBoost model. The
234 sparsity-awareness algorithm applies for any commonly recurring value (e.g., NaN,
235 0). In the context of wastewater treatment, this can apply to missing values and very
236 low levels of effluent pollutants below the threshold for recording.

237 The sparsity awareness algorithm differs from kNN, as the former is a method
238 that is integrated in model training, while the latter is completely independent of soft
239 sensor development. This study sought to identify the best method for handling the
240 missing sensor readings in Cases 1 and 2, according to the characteristics of these
241 respective datasets. The study compared three methods of handling the missing values
242 by developing models using (1) a dataset containing no samples with missing
243 readings; (2) a dataset where the missing sensor readings were filled in using kNN;
244 and (3) the sparsity-awareness split-finding algorithm to train using a dataset with
245 missing values.

2.5. XGBoost and Comparison with other Soft Sensor Models

246 XGBoost is an ensemble method, meaning that it is a collection of weaker models, as
 247 opposed to being a single, highly complex model (i.e., ANN, SVM) (Chen &
 248 Guestrin, 2016). Specifically, it is composed of regression trees (f_k) (Eq. 2). The
 249 structure of the regression tree is represented by its leaves, which correspond to a
 250 numerical weight (w). Each sample is assigned to a set of leaves based on the values
 251 of its input variables. The model's estimated output for that sample is obtained by
 252 adding the sum of the leaves assigned to that sample for each regression tree
 253 (visualized in Figure 2-b).

$$\hat{y}_i = \sum_{k=1}^K f_k(x_i) \quad \text{Eq. 2}$$

254 These regression trees are introduced additively to the ensemble (as f_i for
 255 iteration t), such that each new regression tree minimizes the learning objective (eq.
 256 3). This is different from singular models, which tend to have a pre-defined structure
 257 and are optimized in a Euclidean space.

$$\mathcal{L}^t = \sum_{i=1}^n l(y_i, \hat{y}_i^{t-1}) + f_t(x_i) + \Omega(f_t) \quad \text{Eq. 3}$$

258

259 For benchmarking, XGBoost was compared with an ANN model and an SVM
 260 model. The ANN model was based on Zhu et al. (2017), which is composed of one
 261 hidden layer with 10 neurons, using the radial basis function as its activation function.
 262 The SVM model was based on Zaghoul et al. (2020), which used a Gaussian kernel
 263 function.

264

3. Results

3.1. Selection of Supporting Sensors

265 A modified PLS approach was used to identify the best-supporting sensors for the
 266 proposed BOD₅ soft sensors. In typical implementations of PLS for parameter
 267 selection, the PLS regression coefficients are used as the basis for selection. The
 268 regression coefficients obtained from building models for BOD₅ using data from
 269 Cases 1 and 2 are presented in Figure 2. In all cases, the supporting sensor with the
 270 highest PLS regression coefficient was a complex sensor, i.e. COD and total
 271 suspended solids (TSS). On the other hand, simpler sensors, i.e. Temp_{Reac}, flow rate
 272 (Q), conductivity (Cond), and pH, ranked lower in the order of recommended
 273 supporting sensors.

274 Through the modified PLS approach, the potential of using these simpler
 275 sensors to build the soft sensor was explored. The study compared the difference in
 276 RMSE resulting from using different sets of supporting sensors (see Figure 3). The
 277 initial soft sensor model for each case was built using only simple sensors.
 278 Specifically, for Case 1, these simple sensors were pH_{inf} , flowrate (Q_{inf}) and
 279 conductivity ($Cond_{inf}$) for BOD_{inf} , and pH_{eff} and $Cond_{eff}$ for BOD_{eff} . For Case 2, simple
 280 sensors refer to Q_{inf} , $Temp_{Reac}$ and pH_{inf} for BOD_{inf} , and Q_{eff} and pH_{eff} for BOD_{eff} .
 281 Supporting sensors were incrementally added based on the order of their PLS
 282 regression coefficients until all potential supporting sensors were exhausted. The
 283 sensitivity analysis showed that using a large number of complex supporting sensors
 284 did not improve accuracy. Based on these results, we found that: A soft sensor for
 285 BOD₅ could be built using simple sensors and one complex sensor.

286 The results for Case 1 showed that a soft sensor for BOD_{inf} could be developed
 287 based on simple sensors and COD_{inf} . There was a significant decrease in RMSE from
 288 the model using only simple sensors, to the model using simple sensors and COD_{inf} .

289 However, the improvement in RMSE became minimal for additional supporting
290 sensors. As such, the proposed soft sensor for BOD_{inf} in Case 1 is based on the simple
291 sensors and COD_{inf} . The soft sensor for BOD_{eff} was the only case where COD did not
292 have the highest PLS regression coefficient. In this specific case, the coefficient for
293 COD_{eff} is lower than TSS_{eff} and $Sediments_{eff}$, although it is notable that there was a
294 slight decrease in RMSE when COD_{eff} is added as a soft sensor along with TSS,
295 sediments and the simple sensors. It is also notable that soft sensor performance
296 worsens both in terms of average performance and consistency between the soft
297 sensor with 3 supporting sensors (i.e., simple sensors and TSS_{eff}), and that with 4
298 supporting sensors (i.e., simple sensors, TSS_{eff} and $Sediments_{eff}$). This suggests that
299 having more supporting sensors may generally even worsen performance, potentially
300 due to noise or multicollinearity. Thus, the proposed BOD_{eff} sensor for Case 1 takes
301 the version of the model using 3 supporting sensors. Based on the same reasoning,
302 the proposed supporting sensors for BOD_{inf} are the simple sensors and COD_{inf} . For
303 BOD_{eff} , the proposed supporting sensors are the simple sensors, COD_{eff} , and
304 orthophosphates ($OP-P_{eff}$). This is the only case where further improvement was
305 observed from adding more than one complex sensor as an input. This shows that a
306 soft sensor can be developed with relatively few and accessible supporting sensors.

3.2. Comparing Methods for Missing Sensor Readings

307 The wastewater datasets contain a significant number of samples with missing
308 readings, owing to sensor failure or manual measurement of the parameters. This is a
309 common problem in data-driven modeling, particularly in water treatment (Ma et al.,
310 2020). This study compared different approaches for the missing sensor readings in
311 the dataset. Specifically, the study compared (i) the case where only samples without
312 missing readings were used in training, (ii) the case where kNN was used to fill in the
313 missing readings, and (iii) the case where a dataset with missing values was used to

314 train the XGBoost model, to be processed by its sparsity awareness split-finding
315 algorithm. For Case 1, including the missing sensor readings in training generally
316 improved performance (see Table 1). The sparsity awareness split-finding algorithm
317 of XGBoost resulted in the highest accuracy (i.e., lowest RMSE) for both BOD_{inf} and
318 BOD_{eff} , although the XGBoost model using missing values was obtained from kNN
319 had the highest consistency. Meanwhile, for Case 2, the results of the models were in
320 favor of removing the samples with missing readings from the dataset, for both
321 BOD_{inf} and BOD_{eff} .

322 The difference between the performance of the methods in Cases 1 and 2 was
323 attributed to the volume of missing values in each case. Specifically, only 15.9% of
324 influent samples and 13.7% of effluent samples contained missing sensor readings.
325 This is small compared to Case 2, with 72.0% of influent samples and 71.9% of
326 effluent samples containing missing values. In addition, it was notable that data in
327 Case 1 tended to contain fewer parameters with missing readings in each sample. In
328 comparison, there samples in Case 2 with missing readings tended to contain several
329 missing sensor readings (see supplementary information Table S3). Because of this,
330 kNN and the sparsity awareness algorithm had less inputs for handling the missing
331 values, resulting in poorer estimations.

332 These characteristics of Case 1 make it more viable to include the samples
333 with missing readings in training. This illustrates that there is a threshold for
334 uncertainty in the samples included in the training set. While including some of these
335 samples with missing readings can improve performance, adding a large number of
336 the samples, or using samples with too many missing parameter values, worsen
337 performance. Related studies concerning unlabelled datasets have also encountered
338 this problem, necessitating the selective inclusion of samples for model development
339 (Li et al., 2020).

340 As a method for handling missing values, the results demonstrated that the
341 sparsity awareness algorithm of XGBoost was at least equal to kNN. This makes the
342 estimation process of the soft sensor model more efficient, as the algorithm can
343 directly process the samples with missing readings, whereas kNN results in a two-step
344 approach of imputation and estimation. The significance of the sparsity awareness
345 algorithm method is that it assigns a direction for any sparsely occurring value,
346 whereas other regression tree ensembles would either not be able to use a missing
347 value as an input, or would treat the recurring value as any continuous value. This
348 method is helpful both in training and operating the soft sensor, as the algorithm may
349 allow the soft sensor to continue functioning even if some of the supporting sensors
350 fail.

3.3. Comparison of Soft Sensor Models

351 XGBoost differs from other implementations of regression tree ensembles as its
352 learning objective is penalized with the term $\Omega(f_k)$. This limits the complexity of the
353 regression trees, preventing overfitting. The learning objective is used to determine
354 the optimal structure of regression trees, the assignment of leaves for each sample,
355 and the weighted value of the leaves. The performance of XGBoost was compared
356 with more popular methods in soft sensor development, i.e. ANN and SVM. First, a
357 comparison of observed (laboratory-tested) and estimated (soft sensor) values was
358 conducted to identify the source of error in the models in relation to RMSE. Results to
359 demonstrate this analysis in Case 1 is shown in Figure 4. For BOD_{inf} , the RMSE of
360 XGBoost was inferior to both ANN and SVM; and for BOD_{eff} , the RMSE of XGBoost
361 was superior to both models. The stark difference in performance indicates the dataset
362 characteristics where each model would be more appropriate. Specifically, a
363 continuous regression approach seems to be more effective for the high-variance case

364 of BOD_{inf} , while the ensemble learning approach is compatible with the high
365 skewness BOD_{eff} .

366 Figure 5 shows the results for Case 2. In this case, XGBoost ranks second to
367 ANN in terms of performance for BOD_{inf} . This supports the notion that continuous
368 regression is more appropriate for BOD_{inf} . However, more cases would be needed to
369 identify the difference between Cases 1 and 2 that allowed XGBoost to have an
370 advantage over SVM. On the other hand, the results for BOD_{eff} show that XGBoost
371 had the lowest RMSE in this case, which supports the conclusions drawn from Case 1
372 on the effectiveness of XGBoost on skewed and non-normal distributions. In spite of
373 this, it was found that all three models were challenged when it came to estimating
374 extremely high and extremely low values. In particular, given the high skewness of
375 BOD_{eff} , there were significantly fewer samples to represent extremely high values of
376 BOD_{eff} in the dataset, which can account for poor performance. The visual comparison
377 of observed and estimated values shows that XGBoost is superior in estimating some
378 of these low-frequency cases.

379 The findings based on Cases 1 and 2 analysis were validated with 10-fold
380 cross-validation. This means that each model was tested using 10 different test sets, in
381 cases where these samples were not included in training the model. The results of
382 cross-validation for Case 1 are presented in Table 2a. For BOD_{inf} , the results
383 confirmed that continuous regression was superior for this case; and likewise, the
384 results for BOD_{eff} confirmed that XGBoost was advantageous for skewed
385 distributions. In addition, it can also be observed that while XGBoost did not always
386 have the lowest RMSE, it consistently had the lowest standard deviation, which
387 supports the notion that the residual errors were not higher for extreme values.

388 The results of cross-validation for Case 2 (shown in Table 2b) confirm that
389 both ANN and SVM are superior to XGBoost for BOD_{inf} . Previously, the results of a

390 single test set presented in Figure 5 showed that XGBoost was more accurate than
391 SVM for at least one case. Although the average RMSE from cross validation seemed
392 to converge (between 67.49 – 67.79 mg/L), some variation on a case-to-case basis can
393 be expected given that the SVM model had the highest standard deviation based on
394 cross validation. A model can achieve the highest accuracy for a certain fold if it is
395 the most suitable model for the characteristics of the data in that fold. This was
396 demonstrated by the results of BOD_{eff} , which supported the appropriateness of
397 XGBoost for skewed datasets. Specifically, XGBoost had the highest accuracy and
398 consistency among the three models.

399 In general, XGBoost has some advantages over singular models in terms of
400 robustness and scalability. The characteristic of being an ensemble of models is
401 intended to allow each model to capture some aspect of the data structure. Being
402 composed of several weaker (less complex) models prevents the likelihood of
403 overfitting, even for smaller datasets. Together, the ensemble characteristic and the
404 additive model development process prevent convergence to local minima, a tendency
405 of singular models. These characteristics can also help XGBoost to cover a larger
406 space of potential solutions, resulting in a higher potential for good potential.

407

4. Discussion

408 This study developed soft sensors for BOD_5 for two different wastewater treatment
409 plants. In both cases, the supporting sensors used were a combination of relatively
410 simple sensors (e.g., pH, temperature, flow rate) and minimal complex sensors (i.e.,
411 COD and/or nutrients). This is a common approach in most soft sensor development
412 studies, as simpler sensors may be more stable or easily replaced, while the complex
413 sensors may share stronger correlations with BOD_5 . In a literature study, Xiao et al.
414 (2019) predicted effluent BOD_5 from sensors for pH, effluent ammonia, influent TSS,

415 and influent COD using multivariate regression models. The soft sensor designed by
416 Ebrahimi et al. (2017) predicted effluent BOD₅ from influent TSS, influent total
417 phosphorus (TP) and influent total nitrogen (TN), specifically using the interactions
418 between these parameters in the soft sensor model. Similarly, the supporting sensors
419 for the soft sensor developed by Liu (2017) include influent TSS, effluent ammonia,
420 and simpler sensors such as dissolved oxygen, oxidation-reduction potential, and flow
421 rate.

422 Notably, most studies used sensors for nutrients (e.g., TN, TP, ammonia),
423 COD and TSS as supporting sensors. In this study, both cases showed significant
424 accuracy improvement when COD was included as a supporting sensor. Case 2 also
425 demonstrated the potential improvement from using sensors for nutrients (i.e., OP-P)
426 in predicting effluent BOD₅. However, this study was able to keep the complex
427 sensors to a minimum by using the modified PLS approach for prioritization and
428 performing a sensitivity analysis of soft sensor performance using different supporting
429 sensors.

430 The two cases used in this study varied in terms of statistical properties (see
431 supplementary information Table S4). This affected the model's performance based
432 on RMSE, where data with a higher range (Case 1) also resulted in higher RMSE.
433 Because of this, it is difficult to compare the reported performance of soft sensors
434 developed using different datasets. It should note that most studies use a private
435 dataset, which further limits the potential for comparison. These datasets may have
436 unique characteristics which will influence the conclusions of the study. The size of
437 the dataset alone is an influential factor, affecting the generalizability of the soft
438 sensor. Mjalli et al. (2007) used a relatively small dataset of 73 samples from Doha
439 West Wastewater Treatment Plant. In comparison, the dataset used by Ebrahimi et al.

440 (2017) was composed of 9,180 samples from Floyds Forks Water Quality Treatment
441 Center.

442 This study aimed to make a comprehensive summary of the characteristics of
443 the two datasets used in its analysis. This was intended to allow for comparison
444 between the results of this study on XGBoost, as well as past and future efforts in soft
445 sensor development for wastewater parameters. Aside from summarizing the
446 characteristics of the datasets used, the study used a public dataset in Case 1 (Dua &
447 Graff, 2019), allowing future studies to have the opportunity to make a direct
448 comparison using the same dataset.

449 It was also observed that the majority of studies tended to focus on effluent
450 prediction. In most cases, effluent parameters were predicted using influent
451 parameters. The availability of influent parameters as supporting sensors may be one
452 reason for the majority of soft sensor studies being concerned with the effluent.
453 Previously, some studies were cited which used measures such as influent TSS and
454 influent COD to predict BOD₅. Aside from this, influent parameters such as ammonia
455 and flow rate have been used to predict effluent COD (Cong & Yu, 2018; Grieu et al.,
456 2005). Effluent TP has been predicted using TP and TSS in the influent (Wang et al.,
457 2021; Bagheri et al., 2015). Conversely, it makes no logical reason to predict the
458 influent parameters using effluent data, which may be one reason that there are
459 significantly more soft sensors that have been developed for the effluent, compared to
460 the influent (Ye et al., 2020). In comparison, relatively few soft sensors have been
461 developed for influent parameters. These include models for influent COD and TP
462 developed by Wang et al. (2019); and the model for influent TP of Zhu et al. (2017).
463 So far, this study is one of the only a few studies to develop a soft sensor for the
464 influent BOD₅; the XGBoost-based machine learning model provided good
465 opportunity for achieving this objective.

5. Conclusion

467 This study developed soft sensors for predicting BOD₅ using XGBoost machine
468 learning. This new method was applied to two cases to evaluate its robustness. In both
469 cases, XGBoost estimated a wide range of BOD₅ values, showing consistent
470 performance across different test sets. Although the average performance of machine
471 learning models tended to converge, XGBoost has an innate method of handling
472 missing values; is less prone to overfitting; and was observed to be more effective in
473 measuring higher values of pollutant concentration. XGBoost was particularly
474 effective in estimating effluent BOD₅ which is characterized by important outliers, as
475 cases of high pollutant concentration rarely occur. The soft sensor developed in this
476 study was validated through 10-fold cross validation; however, in future work, we
477 expect to validate the soft sensor in lab-scale or full-scale operation.

478

Acknowledgement

479 The authors give their warmest thanks to the Drainage Services Department of Hong
480 Kong for supporting the research and providing the dataset used in this study. This
481 study was also partially supported by the Hong Kong Innovation and Technology
482 Commission (grant no ITC- CNERC14EG03), Hong Kong Research Grant Council
483 (grant no T21-604/19-R), and Ghent University (BOF/STA/202109/022), Belgium.
484

References

- 485 1. Abdeldayem, O. M., Dabbish, A. M., Habashy, M. M., Mostafa, M. K., Elhefnawy, M.,
486 Amin, L., ... Rene, E. R., 2022. Viral outbreaks detection and surveillance using
487 wastewater-based epidemiology, viral air sampling, and machine learning techniques: A
488 comprehensive review and outlook. *Science of The Total Environment* 803, 149834.
- 489 2. Abu Alfeilat, H. A., Hassanat, A. B., Lasasmeh, O., Tarawneh, A. S., Alhasanat, M. B.,
490 Eyal Salman, H. S., Prasath, V. S., 2019. Effects of distance measure choice on k-nearest
491 neighbor classifier performance: A review. *Big Data* 7, 221-248.
- 492 3. Asami, H., Golabi, M., Albaji, M., 2021. Simulation of the biochemical and chemical
493 oxygen demand and total suspended solids in wastewater treatment plants: Data-mining
494 approach. *Journal of Cleaner Production* 296, 126533.
- 495 4. Asante-Okyere, S., Shen, C., Ziggah, Y. Y., Rulegeya, M. M., Zhu, X., 2020. Principal
496 component analysis (PCA) based hybrid models for the accurate estimation of reservoir
497 water saturation. *Computers & Geosciences* 145, 104555.

- 498 5. Ba-Alawi, A. H., Vilela, P., Loy-Benitez, J., Heo, S., Yoo, C., 2021. Intelligent sensor
499 validation for sustainable influent quality monitoring in wastewater treatment plants using
500 stacked denoising autoencoders. *Journal of Water Process Engineering* 43, 102206.
- 501 6. Bagheri, M., Mirbagheri, S. A., Ehteshami, M., Bagheri, Z., 2015. Modeling of a
502 sequencing batch reactor treating municipal wastewater using multi-layer perceptron and
503 radial basis function artificial neural networks. *Process Saf. Environ.* 93, 111-123.
- 504 7. Chen, T., Guestrin, C., 2016. XGBoost: A scalable tree boosting system. *Proceedings of
505 the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data
506 Mining (KDD '16)*, 784-794.
- 507 8. Cong, Q., Yu, W., 2018. Integrated soft sensor with wavelet neural network and adaptive
508 weighted fusion for water quality estimation in wastewater treatment process.
509 *Measurement* 124, 436-446.
- 510 9. [dataset] Dua, D., Graff, C., 2019. Water treatment dataset. UCI Machine Learning
511 Repository. <http://archive.ics.uci.edu/ml>
- 512 10. Ebrahimi, M., Gerber, E. L., Rockaway, T. D., 2017. Temporal performance assessment
513 of wastewater treatment plants by using multivariate statistical analysis. *J. Environ.
514 Manag.* 193, 234-246.
- 515 11. Fan, Y., Xu, Z., Huang, Y., Wang, T., Zheng, S., DePasquale, A., ... Li, B., 2020. Long-
516 term continuous and real-time in situ monitoring of Pb (II) toxic contaminants in
517 wastewater using solid-state ion selective membrane (S-ISM) Pb and pH auto-correction
518 assembly. *Journal of Hazardous Materials* 400, 123299.
- 519 12. Grieu, S., Traoré, A., Polit, M., Colprim, J., 2005. Prediction of parameters characterizing
520 the state of a pollution removal biologic process. *Eng. Appl. Artif. Intell.* 18, 559-573.
- 521 13. Haimi, H., Mulas, M., Corona, F., & Vahala, R., 2013. Data-derived soft-sensors for
522 biological wastewater treatment plants: An overview. *Environmental Modelling &
523 Software* 47, 88-107.
- 524 14. Hill, A., Tait, S., Baillie, C., Viridis, B., & McCabe, B., 2020. Microbial electrochemical
525 sensors for volatile fatty acid measurement in high strength wastewaters: A review.
526 *Biosensors and Bioelectronics*, 112409.
- 527 15. Hosseinzadeh, A., Zhou, J. L., Altaee, A., Li, D., 2022. Machine learning modeling and
528 analysis of biohydrogen production from wastewater by dark fermentation process.
529 *Bioresource Technology* 343, 126111.
- 530 16. Jaramillo, F., Orchard, M., Muñoz, C., Antileo, C., Sáez, D., Espinoza, P., 2018. On-line
531 estimation of the aerobic phase length for partial nitrification processes in SBR based on
532 features extraction and SVM classification. *Chemical Engineering Journal* 331, 114-123.
- 533 17. Jiang, H., Zou, B., Xu, C., Xu, J., Tang, Y. Y., 2020. SVM-Boosting based on Markov
534 resampling: Theory and algorithm. *Neural Netw.* 131, 276-290.
- 535 18. Li, D., Liu, Y., Huang, D., 2020. Development of semi-supervised multiple-output soft-
536 sensors with Co-training and tri-training MPLS and MRVM. *Chemom. Intell. Lab. Syst.*
537 199, 103970.
- 538 19. Liu, Y., Xie, M., 2020. Rebooting data-driven soft-sensors in process industries: A review
539 of kernel methods. *J. Process Control* 89, 58-73.
- 540 20. Liu, Y., 2017. Adaptive just-in-time and relevant vector machine based soft-sensors with
541 adaptive differential evolution algorithms for parameter optimization. *Chem. Eng. Sci.*
542 172, 571-584.
- 543 21. Ma, J., Ding, Y., Cheng, J. C., Jiang, F., Xu, Z., 2020. Soft detection of 5-day BOD with
544 sparse matrix in city harbor water using deep learning techniques. *Water Res.* 170,
545 115350.
- 546 22. Matheri, A. N., Ntuli, F., Ngila, J. C., Seodigeng, T., Zvinowanda, C., 2021. Performance
547 prediction of trace metals and cod in wastewater treatment using artificial neural network.
548 *Computers & Chemical Engineering* 149, 107308.
- 549 23. Mehmood, T., Sæbø, S., Liland, K. H., 2020. Comparison of variable selection methods
550 in partial least squares regression. *J. Chemom.* 34, e3226.
- 551 24. Mjalli, F. S., Al-Asheh, S., Alfadala, H. E., 2007. Use of artificial neural network black-
552 box modelling for the prediction of wastewater treatment plants performance. *J. Envi.
553 Manag.* 83, 329-338.

- 554 25. Moufid, M., Hofmann, M., El Bari, N., Tiebe, C., Bartholmai, M., Bouchikhi, B., 2021.
555 Wastewater monitoring by means of e-nose, VE-tongue, TD-GC-MS, and SPME-GC-
556 MS. *Talanta* 221, 121450.
- 557 26. Raud, M., Tenno, T., Jōgi, E., Kikas, T., 2012. Comparative study of semi-specific
558 *Aeromonas hydrophila* and universal *Pseudomonas fluorescens* biosensors for BOD
559 measurements in meat industry wastewaters. *Enzyme and Microbial Technology*, 50(4-5),
560 221-226.
- 561 27. Sharafati, A., Asadollah, S. B. H. S., Hosseinzadeh, M., 2020. The potential of new
562 ensemble machine learning models for effluent quality parameters prediction and related
563 uncertainty. *Process Saf. Environ.* 140, 68-78.
- 564 28. Qi, X., Guo, H., Wang, W., 2021. A reliable KNN filling approach for incomplete
565 interval-valued data. *Eng. Appl. Artif. Intell.* 100, 104175.
- 566 29. Qin, X., Gao, F., Chen, G., 2012. Wastewater quality monitoring system using sensor
567 fusion and machine learning techniques. *Water Res.* 46, 1133-1144.
- 568 30. Qiu, J., Lü, F., Zhang, H., Shao, L., He, P., 2021. Data mining strategies of molecular
569 information for inspecting wastewater treatment by using UHRMS. *Trends in*
570 *Environmental Analytical Chemistry*, e00134.
- 571 31. Wang, G., Jia, Q. S., Zhou, M., Bi, J., Qiao, J., 2021. Soft-sensing of wastewater
572 treatment process via deep belief network with event-triggered learning. *Neurocomputing*
573 436, 103-113.
- 574 32. Wang, T., Xu, Z., Huang, Y., Dai, Z., Wang, X., Lee, M., ... Li, B., 2020. Real-time in
575 situ auto-correction of K⁺ interference for continuous and long-term NH₄⁺ monitoring in
576 wastewater using solid-state ion selective membrane (S-ISM) sensor assembly.
577 *Environmental Research* 189, 109891.
- 578 33. Wang, X., Kvaal, K., Ratnaweera, H., 2019. Explicit and interpretable nonlinear soft
579 sensor models for influent surveillance at a full-scale wastewater treatment plant.
580 *J. Process Control* 77, 1-6.
- 581 34. Wu, X., Kumar, V., Quinlan, J.R., Ghosh, J., Yang, Q., Motoda, H. ... Zhou, Z.H., 2008.
582 Top 10 algorithms in data mining. *Knowl. Inf. Syst.* 14, 1-37.
- 583 35. Xiao, H., Bai, B., Li, X., Liu, J., Liu, Y., Huang, D., 2019. Interval multiple-output soft
584 sensors development with capacity control for wastewater treatment applications: A
585 comparative study. *Chemom. Intell. Lab. Syst.* 184, 82-93.
- 586 36. Ye, Z., Yang, J., Zhong, N., Tu, X., Jia, J., Wang, J., 2020. Tackling environmental
587 challenges in pollution controls using artificial intelligence: A review. *Sci. Total Environ.*
588 699, 134279.
- 589 37. Zaghoul, M. S., Hamza, R. A., Iorhemen, O. T., Tay, J. H., 2020. Comparison of
590 adaptive neuro-fuzzy inference systems (ANFIS) and support vector regression (SVR) for
591 data-driven modelling of aerobic granular sludge reactors. *J. Environ. Chem. Eng.* 8,
592 103742.
- 593 38. Zhu, S., Han, H., Guo, M., Qiao, J., 2017. A data-derived soft-sensor method for
594 monitoring effluent total phosphorus. *Chin. J. Chem. Eng.* 25, 1791-1797.

Tables**Table 1** RMSE (mg/L) of the model trained on the (a) UCI Machine Learning Repository dataset and (b) Hong Kong dataset (using different methods of handling missing values)

	(a) Case 1: UCI Machine Learning Repository				(b) Case 2: Hong Kong Dataset			
	Influent BOD		Effluent BOD		Influent BOD		Effluent BOD	
	Ave	Std. Dev.	Ave	Std. Dev.	Ave	Std. Dev.	Ave	Std. Dev.
Samples without missing values	52.41	9.06	10.59	7.96	67.79	17.52	0.47	0.20
Missing values filled in with kNN	52.07	8.96	10.59	8.01	70.64	16.64	0.77	0.86
Missing values processed by XGBoost	51.93	9.31	10.55	7.98	68.60	19.97	1.17	1.48

Table 2 RMSE (mg/L) of 10-fold cross-validation for models developed using the (a) UCI Machine Learning Repository dataset and (b) Hong Kong dataset.

	(a) Case 1: UCI Machine Learning Repository		(b) Case 2: Hong Kong Dataset	
	Influent BOD			
	Average	Std. Dev.	Average	Std. Dev.
XGBoost*	51.93	9.31	67.79	17.52
ANN with kNN	50.51	11.04	67.58	19.60
SVM with kNN	50.51	11.04	67.49	23.58
	Effluent BOD			
	Average	Std. Dev.	Average	Std. Dev.
XGBoost *	10.55	7.98	0.47	0.20
ANN	10.80	9.95	0.48	0.30
SVM	11.98	12.87	0.51	0.41

* Note: For Case 1, the XGBoost were analyzed with Sparsity Awareness Algorithm

Figures

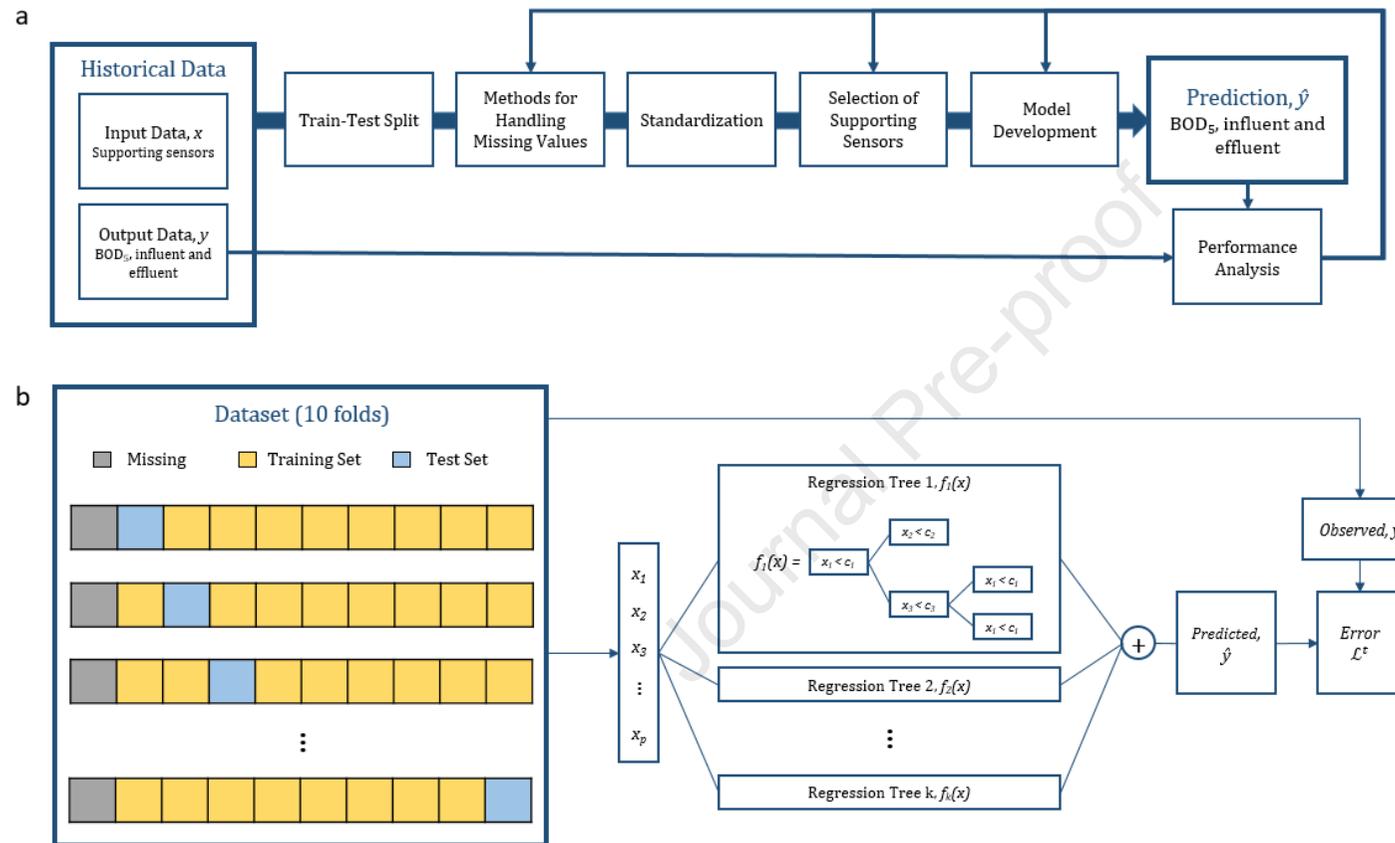
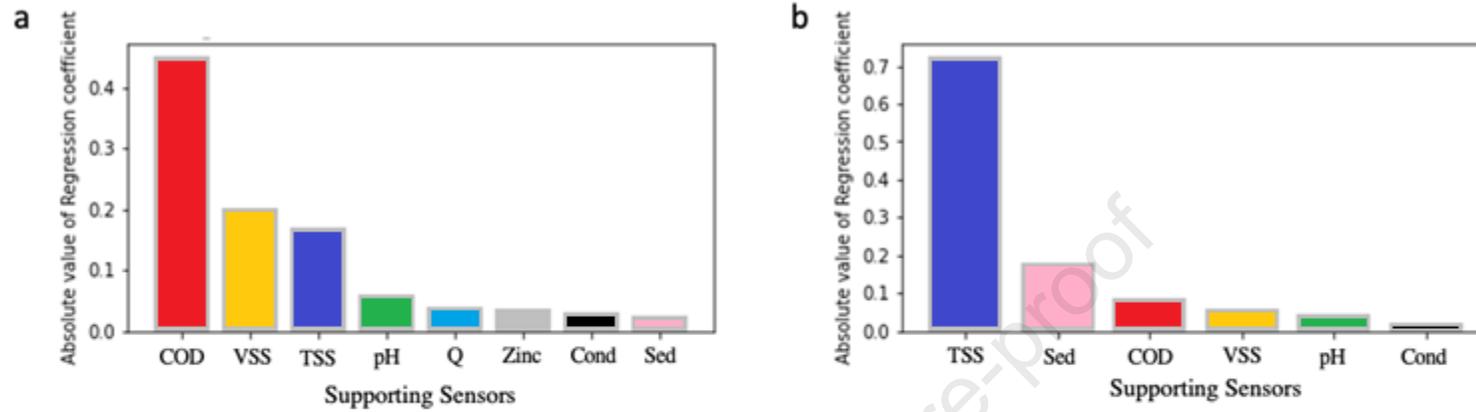


Figure 1 Soft sensor development frameworks: (a) methods applied, and (b) XGBoost model structure.

Case 1: UCI Machine Learning Repository Dataset



Case 2: SWHSTW Dataset

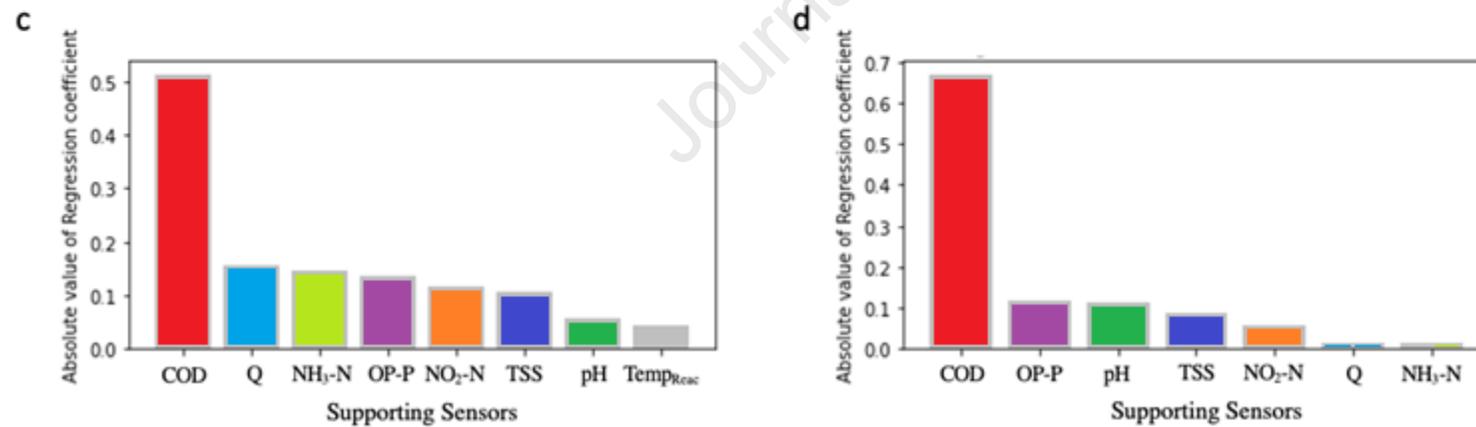
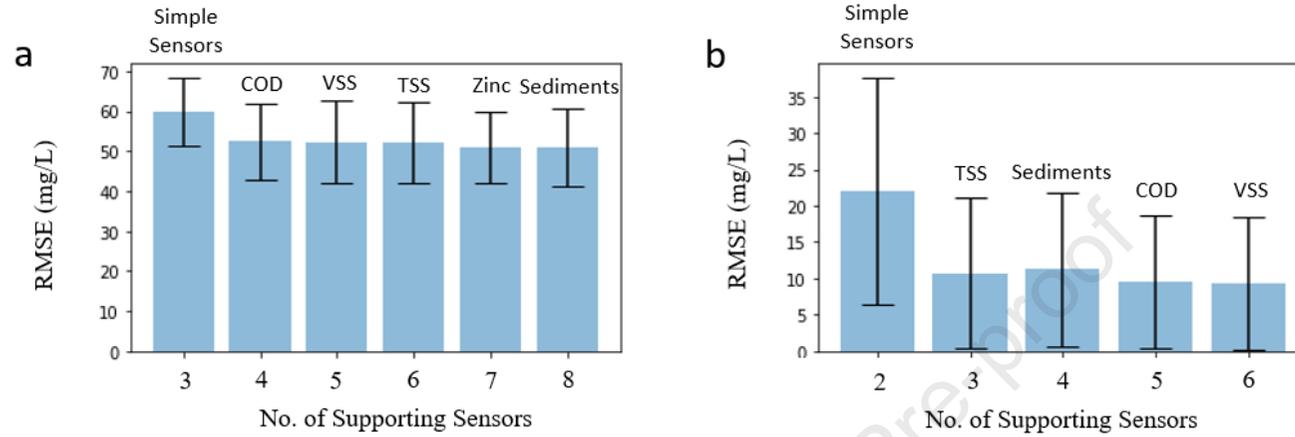


Figure 2 PLS regression coefficients for (a and c) influent BOD and (b and d) effluent BOD, with common parameters indicated by color.

Case 1: UCI Machine Learning Repository Dataset



Case 2: SWHSTW Dataset

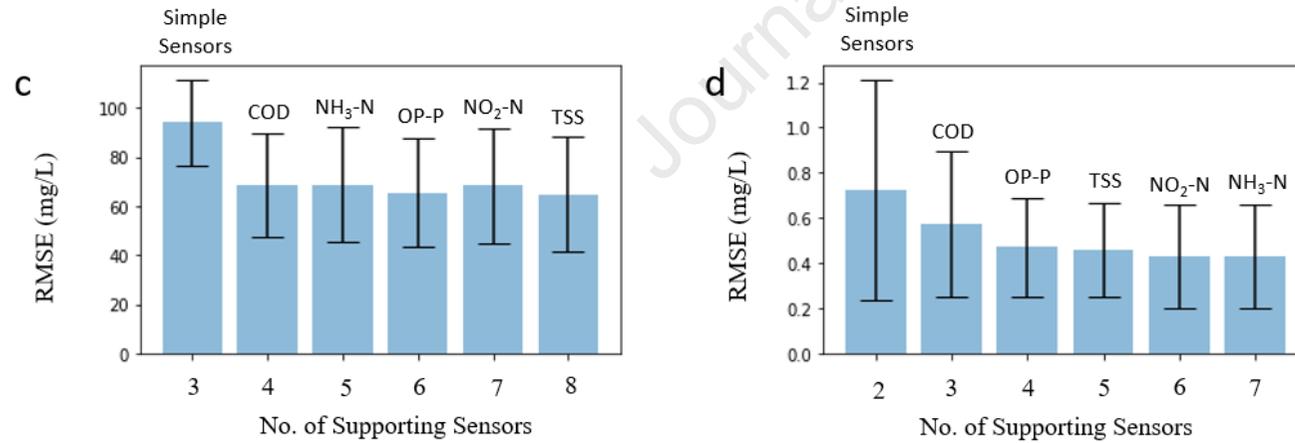


Figure 3 Change in RMSE (mg/L) supporting sensors are incrementally added to the soft sensor for (a, c) influent BOD and (b, d) effluent BOD.

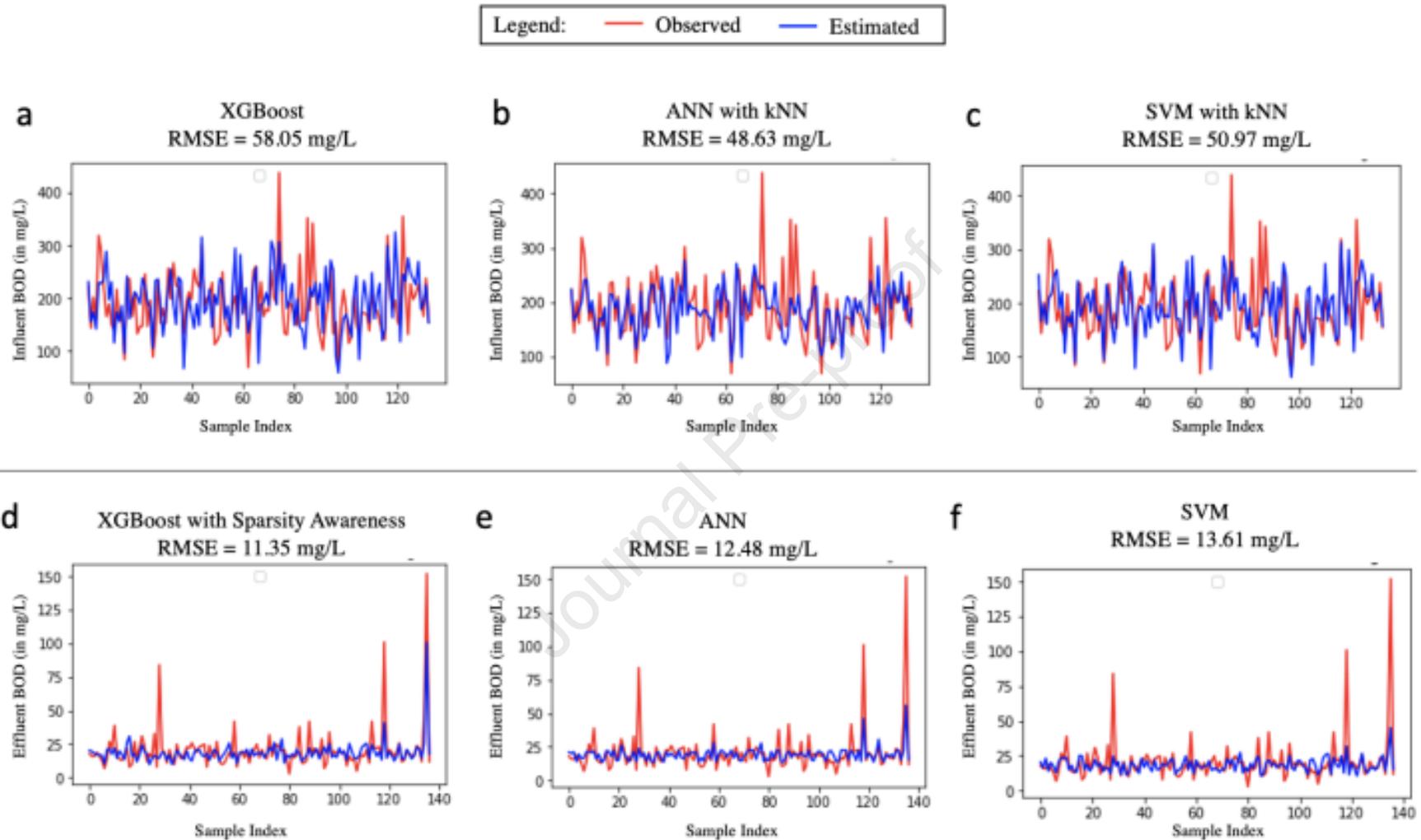


Figure 4 Visual comparison and RMSE (mg/L) of (a) BOD_{inf} estimated by XGBoost; (b) BOD_{inf} estimated by ANN with kNN; (c) BOD_{eff} estimated by SVM with kNN; (d) BOD_{eff} estimated by XGBoost; (e) BOD_{eff} estimated by ANN with kNN; and (e) BOD_{eff} estimated by SVM with kNN, modeled using the UCI Machine Learning Repository dataset.

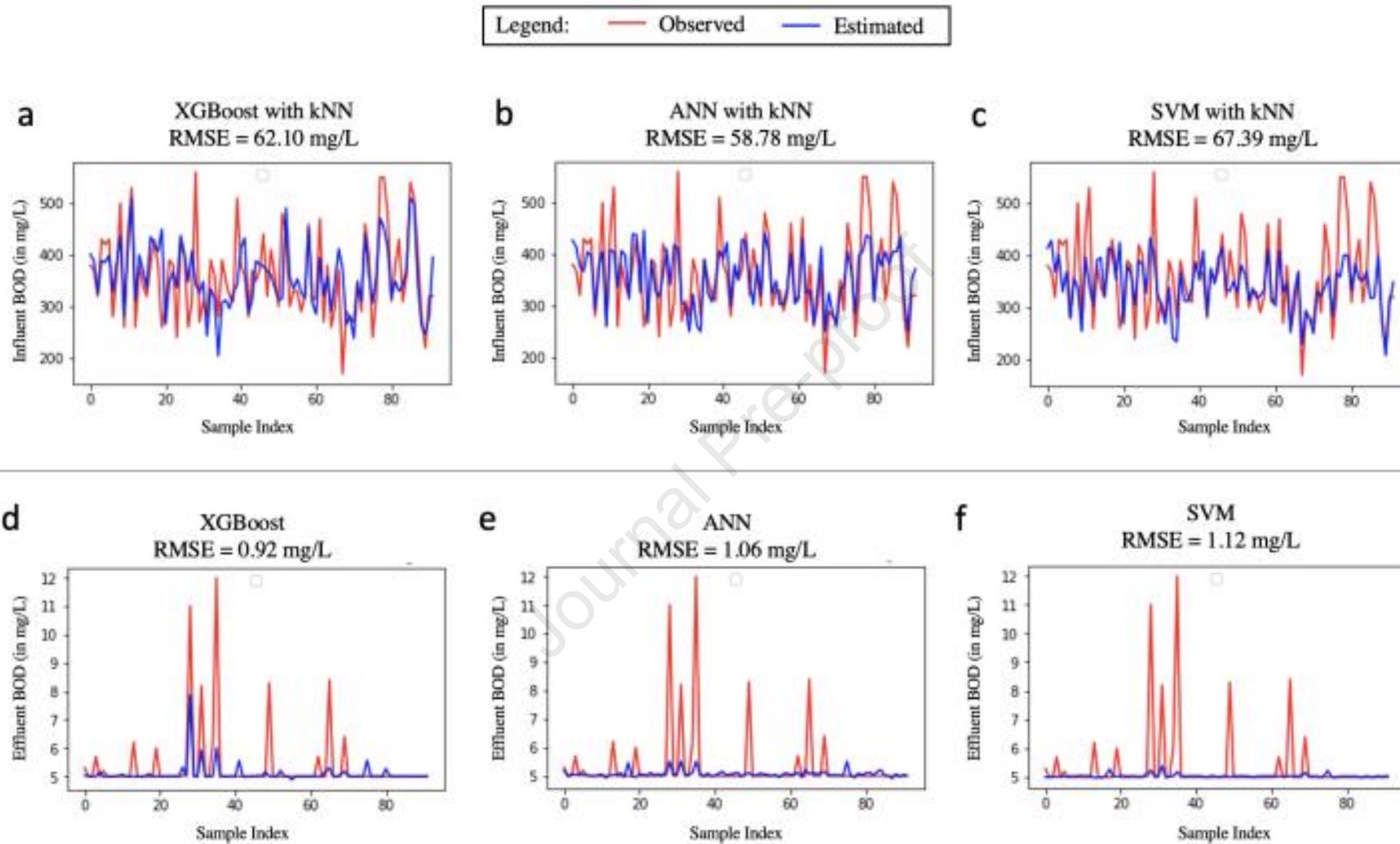


Figure 5 Visual comparison and RMSE (mg/L) of (a) BOD_{inf} estimated by XGBoost; (b) BOD_{inf} estimated by ANN with kNN; (c) BOD_{eff} estimated by SVM with kNN; (d) BOD_{eff} estimated by XGBoost; (e) BOD_{eff} estimated by ANN with kNN; and (f) BOD_{eff} estimated by SVM with kNN, modelled using the Hong Kong dataset.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this review paper.

Journal Pre-proof