

# Journal Pre-proof

Hybrid decision tree-based machine learning models for short-term water quality prediction

Hongfang Lu, Xin Ma



PII: S0045-6535(20)30362-3

DOI: <https://doi.org/10.1016/j.chemosphere.2020.126169>

Reference: CHEM 126169

To appear in: *ECSN*

Received Date: 21 October 2019

Revised Date: 4 February 2020

Accepted Date: 9 February 2020

Please cite this article as: Lu, H., Ma, X., Hybrid decision tree-based machine learning models for short-term water quality prediction, *Chemosphere* (2020), doi: <https://doi.org/10.1016/j.chemosphere.2020.126169>.

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.

© 2020 Published by Elsevier Ltd.

**CRediT author statement**

Hongfang Lu: Conceptualization, Methodology, Data curation, Writing- Original draft preparation

Xin Ma: Investigation, Writing- Reviewing and Editing

# 1 Hybrid decision tree-based machine learning models for

## 2 short-term water quality prediction

3 Hongfang Lu<sup>a,b,\*</sup>, Xin Ma<sup>c</sup>

4 <sup>a</sup> State Key Laboratory of Oil and Gas Reservoir Geology and Exploitation, Southwest Petroleum University, Chengdu 610500,  
5 China

6 <sup>b</sup> School of Science, Southwest University of Science and Technology, Mianyang 621010, China

7 <sup>c</sup> Trenchless Technology Center, Louisiana Tech University, Ruston LA71270, United States

8 \*Corresponding author: [hlu006@latech.edu](mailto:hlu006@latech.edu) (H. Lu)

9 Address: 599 Dan Reneau Drive, Ruston, LA 71270, USA

10 **Abstract:** Water resources are the foundation of people's life and economic development, and are  
11 closely related to health and the environment. Accurate prediction of water quality is the key to  
12 improving water management and pollution control. In this paper, two novel hybrid decision  
13 tree-based machine learning models are proposed to obtain more accurate short-term water quality  
14 prediction results. The basic models of the two hybrid models are extreme gradient boosting  
15 (XGBoost) and random forest (RF), which respectively introduce an advanced data denoising  
16 technique - complete ensemble empirical mode decomposition with adaptive noise (CEEMDAN).  
17 Taking the water resources of Gales Creek site in Tualatin River (one of the most polluted rivers in  
18 the world) Basin as an example, a total of 1875 data (hourly data) from May 1, 2019 to July 20, 2019  
19 are collected. Two hybrid models are used to predict six water quality indicators, including water  
20 temperature, dissolved oxygen, pH value, specific conductance, turbidity, and fluorescent dissolved  
21 organic matter. Six error metrics are introduced as the basis of performance evaluation, and the results  
22 of the two models are compared with the other four conventional models. The results reveal that: (1)  
23 CEEMDAN-RF performs best in the prediction of temperature, dissolved oxygen and specific  
24 conductance, the mean absolute percentage errors (MAPEs) are 0.69%, 1.05%, and 0.90%,  
25 respectively. CEEMDAN-XGBoost performs best in the prediction of pH value, turbidity, and  
26 fluorescent dissolved organic matter, the MAPEs are 0.27%, 14.94%, and 1.59%, respectively. (2)  
27 The average MAPEs of CEEMDAN-RF and CEEMDAN-XGBoost models are the smallest,  
28 which are 3.90% and 3.71% respectively, indicating that their overall prediction performance is the  
29 best. In addition, the stability of the prediction model is also discussed in this paper. The analysis  
30 shows that the prediction stability of CEEMDAN-RF and CEEMDAN-XGBoost is higher than other  
31 benchmark models.

32 **Keywords:** decision tree-based model; short-term; water quality prediction; extreme gradient  
33 boosting; random forest; data denoising

34 **Nomenclature**

35	$\overline{IMF}_k$	$k$ -th modal component
36	$E_k(\cdot)$	$k$ -th modal component obtained by EMD decomposition
37	$O_t$	observation value at time $t$
38	$P_t$	prediction value at time $t$
39	$d^i(t)$	$i$ -th signal sequence
40	$f_i$	final prediction result
41	$f_n$	normalized prediction result
42	$g_i$ and $h_i$	the first and second derivatives of the loss function in the gradient direction
43	$wn^i(t)$	white noise
44	$x_i$	$i$ -th eigenvector
45	$\hat{y}$	predicted value
46	$y_i$	the true value of the sample
47	$z_i$	raw data
48	$z_{min}$ and $z_{max}$	minimum and maximum of the raw data, respectively
49	$z_n$	normalized data
50	$\Omega$	penalty term for model complexity
51	$F$	regression tree set
52	$I$	number of tests
53	$IMFs$	intrinsic mode functions
54	$K$	number of trees
55	$R$	final residue
56	$T$	number of leaves in the tree
57	$f$	a regression tree in tree space $F$
58	$l$	loss function
59	$r$	residue
60	$t$	current iteration
61	$w$	internal split tree weight
62	$\gamma$ and $\lambda$	configurable parameters
63	$\varepsilon$	noise standard deviation
64		

65 **Abbreviations**

66	ANN	artificial neural network
67	BPNN	back propagation neural network
68	CEEMDAN	complete ensemble empirical mode decomposition with adaptive noise
69	CPU	central processing unit
70	EEMD	ensemble empirical mode decomposition
71	EMD	empirical mode decomposition
72	FDOM	fluorescent dissolved organic matter
73	FNU	Formazin Nephelometric Units
74	GMNN	Gamma memory neural network
75	JENN	Jordan-Elman neural network
76	LSSVM	least squares support vector machine

77	LSTM	long short-term memory
78	MAE	mean absolute error
79	MAPE	mean absolute percentage error
80	Max.	maximum
81	Min.	minimum
82	ppb	parts per billion
83	PSO	particle swarm optimization
84	QSE	quinine sulfate equivalents
85	RBFNN	radial basis function neural network
86	RF	random forest
87	RMSE	root mean square error
88	RMSPE	root mean squared percentage error
89	SD	standard deviation
90	SDE	standard deviation of errors
91	SVD	singular value decomposition
92	SVM	support vector machine
93	U1	Theil U statistic 1
94	U2	Theil U statistic 2
95	USGS	U.S. Geological Survey
96	XGboost	extreme gradient boosting
97		

## 98 1. Introduction

99 Water is the resource that human beings depend on for survival. The prevention and control of  
100 water pollution have been a hot issue of concern in recent decades. Accidents of water pollution  
101 have been reported uninterruptedly all over the world in recent years (Hounslow, 2018). It can be  
102 seen that many water pollution accidents occur in the absence of information management. Many  
103 water pollution accidents can only be remedied afterward because of the lack of early prediction.  
104 Therefore, water quality prediction has become a hot topic in water environment science.

105 Water quality prediction is an essential work in water environment management. Accurate  
106 forecasting value will undoubtedly improve the management level of water resources. At present,  
107 many water resource management departments have set monitoring points to observe water quality  
108 changes, but they cannot play the role of water quality prediction. However, these monitoring data  
109 can provide a predictive basis for some data-driven models. Through real-time control of future  
110 water quality changes, water pollution degree can be accurately judged. Besides, accurate water

111 quality predictions can also provide a basis for policymakers and provide data to the environmental  
 112 management department to act as an “early warning”.

113 We reviewed the literature on water quality predictions in recent years. Zhao et al. (2007) used  
 114 BPNN to predict the water quality of Yuqiao Reservoir. Singh et al. (2009) used the ANN model to  
 115 predict the water quality of the river and analyzed it with the case of the Gomti river in India. Liu et  
 116 al. (2016) adopted a multi-task multi-view learning method to predict urban water quality. The  
 117 effectiveness of the method was verified by experiments. Palani et al. (2008) applied the ANN  
 118 model to the prediction of coastal water quality in Singapore. Water quality parameters considered  
 119 include salinity, temperature, dissolved oxygen, and chlorophyll-a. West and Dellana (2011) used  
 120 JENN and GMNN to predict the basin water quality. Chan et al. (2013) established a 3D  
 121 hydrodynamic model to predict the water quality of Hong Kong beach in real time. The experiment  
 122 proved that the accuracy of prediction is higher than 80%. Meyers et al. (2017) used several  
 123 machine learning-based models (ANN, RF, and SVM) to predict water turbidity. Peng et al. (2019)  
 124 proposed a framework for real-time prediction of daily water quality and applied it to Lake Chaohu  
 125 in China, which can better predict dissolved oxygen, total phosphorus, and other parameters. Huang  
 126 et al. (2019) established a prediction system for urban estuary water quality and used the gradient  
 127 boosting machine model to fill and predict the flow. Alba et al. (2019) used an ANN-based model  
 128 to predict the bathing water quality of the estuary, which combines laboratory analysis, machine  
 129 learning, and numerical simulation to achieve real-time water quality management. According to  
 130 the literature review, some researchers have established water quality prediction systems and used  
 131 experimental methods to predict water quality. However, physical methods are time-consuming and  
 132 labor-intensive. With the prevalence of machine learning and deep learning (Lu et al., 2020a,b;  
 133 Kong and Ma, 2018; Ma et al., 2020), more and more scholars use intelligent models to predict  
 134 water quality (Liang et al., 2020; Gao et al., 2019; Dabrowski et al., 2020; Zhang et al., 2017;  
 135 Hussein et al., 2019; Panidhapu et al., 2020). In addition, intelligent model has been applied not  
 136 only in water quality prediction, but also in the field of environment and energy engineering (Xie et  
 137 al., 2020; Wu et al., 2019; Zhang et al., 2019; Yang et al., 2020; Zhao et al., 2019). Although the  
 138 accuracy of water quality prediction is improving, because water quality is unstable and non-linear  
 139 in time series, more accurate prediction methods are worthy of further study. Therefore, this paper  
 140 proposes two hybrid models, which combine CEEMDAN method with the original models, thus  
 141 enhancing their prediction accuracy.

142 The rest of the paper is organized as follows: Section 2 introduces six water quality indicators,  
 143 Section 3 describes the collected data, data cleaning technology--CEEMDAN and two decision  
 144 tree-based machine learning models--XGBoost and RF. Section 4 describes the process of

145 prediction and error metrics. Section 5 reveals the prediction results and discussions. Section 6  
146 summarizes the main conclusions and future works.

147

148 **2. Water quality indicators**

149 There are many water quality indicators, which can be divided into chemical indicators,  
150 physical indicators, biological indicators, radioactive indicators, and so on (Valdivia-Garcia et al.,  
151 2019). Water for different purposes usually has different evaluation indicators. For example, the  
152 critical indicators of domestic water use include temperature, pH value, biochemical oxygen  
153 demand, while the indicators of food industry water use include suspended solid, pH value and  
154 number of Escherichia coli. This paper introduces several common water quality indicators (see  
155 Table 1), including temperature, dissolved oxygen, pH value, specific conductance, turbidity, and  
156 FDOM.

157

Table 1. Common water quality indicators and explanations.

Indicator	Explanation	Reference(s)
Water temperature	Temperature is a critical physical indicator of water. The sudden rise in water temperature during daily monitoring indicates that the water body may be contaminated by new sources of pollution. Thermal pollution may also cause biological growth to increase and cause biological pollution in the water.	Tao et al., 2020; Graf et al. 2019
Dissolved oxygen	Free oxygen dissolved in water is called dissolved oxygen, it is a parameter to measure the quality of water. In general, the concentration of dissolved oxygen in water is called equilibrium concentration when it approaches saturation. The saturation value of dissolved oxygen is 9.17 mg/L at 20°C. When water is polluted by oxygen-consuming pollutants, dissolved oxygen decreases. Oxygen-consuming pollutants include carbohydrates, proteins, oils, amino acids, fatty acids, esters and other organic compounds. These pollutants mainly come from domestic sewage and some industrial wastewater.	Larsen et al., 2019
pH value	The pH of natural water is generally between 6.5 and 8.5. A suitable pH range for drinking water is from 7 to 8.5, with a limit range of 6.5 to 9.2. Generally, fish live normally in water with a pH of 6.5 to 8.5. The crop is suitable for growth in water with a pH of 6 to 7.5. Long-term irrigation of water with a pH lower than 5.5 will cause the nitrifying bacteria in the soil to be inhibited, the nitrification will be weakened, and the nitrogen fertilizer will not be fully released.	Mosley et al., 2010
Specific conductance	Specific conductance represents the ability to conduct current in aqueous solution, and it is also an index for routine monitoring of water quality with multi-parameters. The specific conductance is proportional to the ion content in the solution, so the total soluble matter content can be estimated indirectly.	Makarewicz et al., 2012
Turbidity	Turbidity indicates the obstruction extent of light sources by suspended matter in water. The reason for the increase in turbidity in rivers and lakes is because the river water contains many suspended substances. When the turbidity is large, it will affect the photosynthesis of aquatic organisms and reduce the self-purification ability of water.	Kerr et al., 2018
FDOM	FDOM is a fast and simple method for tracking dissolved organic matter in natural water, it is often used as a biological indicator of lake water.	Liu et al., 2019

158 **3. Material and methods**

159 In this paper, two new hybrid models, CEEMDAN-XGBoost and CEEMMDAN-RF, are used to  
160 predict water quality indicators. This section describes the collected data and introduces the relevant  
161 theories of CEEMDAN, XGBoost, and RF.

162 *3.1. Collected data description*

163 The water quality data for this paper is from the Tualatin River in Oregon, USA, it is called one  
164 of the ten most dangerous rivers in the world. The Tualatin River drains 712 square miles in the  
165 northwest corner of Oregon, it is a sub-basin of the Willamette River Basin. It is about 134  
166 kilometers long, and the slope of most sections is very flat. The main tributaries of the Tualatin  
167 River include Scoggins, Gales, Dairy, Rock and Fanno Creek. The summer water flow is discharged  
168 by water from the Scoggins Reservoir and the Barney Reservoir, which diverts water into the upper  
169 Tualatin River. Wastewater from wastewater treatment plants accounts for a large portion of  
170 summer river flows.

171 Before the 1970s, wastewater treatment plants discharged high concentrations of ammonia,  
172 nitrogen, and phosphorus into the mainstream of the Tualatin River. High ammonia concentration  
173 usually causes obvious nitrification in rivers, resulting in low dissolved oxygen concentration. In  
174 addition, in summer, the abundance of phytoplankton in the Tualatin River results in the river's  
175 water quality violating the requirements of minimum dissolved oxygen and maximum pH value.  
176 Later, in 1970, the Unified Sewerage Agency of Washington County was established. They used  
177 various water treatment methods to control the pollution and health problems of the Tualatin River.  
178 By 2002, many water bodies in the Tualatin River Basin had been confirmed to be damaged.

179 This paper chooses the time series data of the Gales Creek site of Tualatin River as the research  
180 object. The raw data of the monitoring point is from USGS (<https://www.usgs.gov>). The collected  
181 water quality data are temperature, dissolved oxygen, pH value, specific conductance, turbidity, and  
182 FDOM from 0:00 on May 1, 2009 to 23:00 on July 20, 2019 (the data interval is one hour). Each  
183 water quality indicator has 1875 data, and their statistical descriptions are shown in Table 2. It can  
184 be seen that the six datasets have multiple probability density function types such as Triangular,  
185 Johnson SB and Lognormal (3P), indicating that the data used in this paper is diverse and extensive,  
186 and lays a foundation for more convincing conclusions.

187

**Table 2.** Statistical descriptions of the data in this paper.

Water quality indicator	Unit	Data period (month/day/year)	Amount of data	Test set data amount	Statistical distribution			Statistical characteristics			
					Distribution	Parameters		Max.	Min.	Mean	SD
Temperature	°C	05/01/2019 00:00-07/20/2019 23:00	1875	187	Triangular	$m=17.73$ , $a=9.853$ , $b=23.722$		23.55	9.98	17.08	2.79
Dissolved oxygen	mg/L				Johnson SB	$\gamma=0.24388$ , $\delta=1.121$ , $\lambda=4.2411$ , $\xi=7.2869$		11.11	7.12	9.21	0.81
pH	Dimensionless				Error function	$k=3.0988$ , $\sigma=0.08784$ , $\mu=7.3089$		7.52	7.08	7.31	0.09
Specific conductance	uS/cm				Johnson SB	$\gamma=-0.25905$ , $\delta=0.51603$ , $\lambda=35.138$ , $\xi=94.969$		133.9	93.8	115.24	10.70
Turbidity	FNU				Lognormal (3P)	$\sigma=0.96807$ , $\mu=0.22892$ , $\gamma=0.45261$		16.0	0.5	2.25	1.72
FDOM	ppb QSE				Dagum	$k=0.82881$ , $\alpha=15.915$ , $\beta=12.418$ , $\gamma=0$		32.62	8.46	12.28	1.75

188

189    *3.2. Decision tree-based machine learning models*

190       The basic models of the two hybrid models used in this paper are XGBoost and RF. They are  
 191       all belong to decision tree-based machine learning models. The decision tree-based model has many  
 192       advantages:

193       a) Ability to handle both data and regular attributes; b) Insensitive to missing values; c) High  
 194       efficiency, the decision tree only needs to be built once. In fact, there are other models in the field  
 195       of machine learning, such as ANN and SVM. Compared to them, decision tree-based models may  
 196       have faster calculation speed and are more conducive to short-term prediction. Moreover, water  
 197       quality monitoring data sometimes have missing values due to equipment failure, the decision  
 198       tree-based model has an advantage in forecasting.

199

200    **3.2.1. eXtreme Gradient Boosting (XGBoost)**

201       XGBoost was proposed by Chen and Guestrin (2016) and is based on the C++ language (Nobre  
 202       and Neves, 2019). The model has achieved great success since its appearance, and it is always seen  
 203       in the top models in various data mining competitions. XGBoost is able to integrate multiple weak  
 204       learning machines into one strong learning machine by iterating and generating multiple trees, and  
 205       it has the following features: a) It can automatically utilize the multithreading of the CPU for  
 206       parallelism, while improving the algorithm to improve accuracy, and this is the most prominent  
 207       feature of XGBoost; b) It is a lifting learning algorithm based on the decision tree model and can  
 208       process sparse data automatically; c) large amounts of data can be processed at high speed  
 209       according to block technology.

210       In the XGBoost model, tree model adopts additive model

$$\hat{y} = \sum_{k=1}^K f_k(x_i), f_k \in F \quad (1)$$

212       The objective function is

$$L(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k) \quad (2)$$

213       where  $\Omega(f) = \gamma T + 0.5\lambda\|w\|^2$ ,  $w = (w_1, w_2, \dots, w_k)$ .

215       Because learning all tree parameters at once is challenging, XGBoost uses an additive strategy  
 216       that learns the parameters of one tree at a time:

$$\begin{aligned} \hat{y}_i^{(0)} &= 0 \\ \hat{y}_i^{(1)} &= \hat{y}_i^{(0)} + f_1(x_i) \\ \hat{y}_i^{(2)} &= \hat{y}_i^{(1)} + f_2(x_i) \\ &\vdots \end{aligned}$$

217  $\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + f_t(x_i)$  (3)

218 The XGBoost algorithm uses the stepwise forward additive model as the gradient boosting  
 219 algorithm. The difference is that the gradient boosting algorithm is a negative gradient that learns a  
 220 weak learner to approximate the loss function. The XGBoost algorithm first finds the second-order  
 221 Taylor approximation of the loss function at that point, and then minimizes the approximation loss  
 222 function to train the weak learner. Therefore, the objective function can be expressed as

223  $L^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t)$  (4)

224 Using the second-order Taylor expansion, the following function can be obtained

225  $L^{(t)} = \sum_{i=1}^n [g_i f_t(x_i) + 0.5 h_i f_t^2(x_i)] + \Omega(f_t)$  (5)

227 **3.2.2. Random forest**

228 RF is an integrated learning method for classification and regression (Liaw and Wiener, 2002).  
 229 It is one of the representative of ensemble learning, and it is an additive model based on bagging  
 230 algorithm. Different from bagging, when constructing each tree, RF uses a random sample predictor  
 231 before each node segmentation, which can reduce bias. It has the following characteristics: a) The  
 232 introduction of two randomness makes RF not easy to fall into overfitting, and has excellent noise  
 233 immunity; b) It can process data of high dimension (many features) without feature selection; c) It has  
 234 fast training speed and is easy to be parallelized, so it is relatively simple to implement (Wu et al.,  
 235 2019). More details about RF can be found in the literature (Liaw and Wiener, 2002).

236

237 *3.3. Data denoising method: CEEMDAN*

238 Short-term water quality may be affected by factors such as temperature, industrial wastewater  
 239 discharge and so on, so that the data may have large fluctuations in the time series and exhibit a  
 240 high degree of nonlinear characteristics, which undoubtedly increases the difficulty of prediction.  
 241 Therefore, many scholars use EMD, SVD, EEMD, wavelet decomposition, and other methods to  
 242 extract feature values. Although these methods can improve prediction accuracy to some extent,  
 243 they all have some limitations (see Table 2). For example, EMD is prone to mode mixing, while  
 244 wavelet decomposition is often not ideal in some practical problems.

245 CEEMDAN is based on EEMD by adding a limited number of adaptive white noise (Torres et  
 246 al., 2011; Zhou et al., 2019). Its implementation process is as follows:

247 (1) Add a white noise sequence to the raw signal to generate a noisy signal set

248  $d^i(t) = d(t) + \varepsilon_0 w n^i(t), i = 1, 2, \dots I$  (6)

249 (2) EMD decomposition operation is performed on the signal set to obtain

250  $\overline{IMF_1}(t) = I^{-1} \sum_{i=1}^I IMF_1^i(t)$  (7)

251 (3) Calculate the margin signal of the first stage ( $k = 1$ )

252  $r_1(t) = d(t) - \overline{IMF_1}(t)$  (8)

253 (4) Calculate the second modal component

254  $\overline{IMF_2}(t) = I^{-1} \sum_{i=1}^I E_1 \{r_1(t) + \varepsilon_1 E_1 [wn^i(t)]\}$  (9)

255 (5) For the following stages, calculate the  $k$ -th margin signal in the same way

256  $r_k(t) = r_{k-1}(t) - \overline{IMF_k}(t)$  (10)

257 (6) Calculate the  $(k + 1)$ -th modal component

258  $\overline{IMF_{k+1}}(t) = I^{-1} \sum_{i=1}^I E_1 \{r_k(t) + \varepsilon_k E_k [wn^i(t)]\}$  (11)

259 (7) Repeat Eq. (10) until the residual component no longer satisfies the decomposition condition.

260 Finally, the original signal  $d(t)$  can be expressed as

261  $d(t) = \sum_{i=1}^K \overline{IMF_i}(t) + R(t)$  (12)

262

## 263 4. Prediction process and error metrics

### 264 4.1. Prediction process

265 (1) Data decomposition

266 CEEMDAN is used for data decomposition and denoising, so that raw data with large  
267 fluctuations is decomposed into multiple datasets with less fluctuations. In other words, the data in  
268 the same dataset has more obvious similar features, as can be seen from Appendix 1.

269 (2) Data normalization

270 In order to eliminate the dimensional influence of the data indicators, the data after the  
271 decomposition is normalized and limited to the range of [0,1], the equation is

272  $z_n = \frac{z_i - z_{min}}{z_{max} - z_{min}}$  (13)

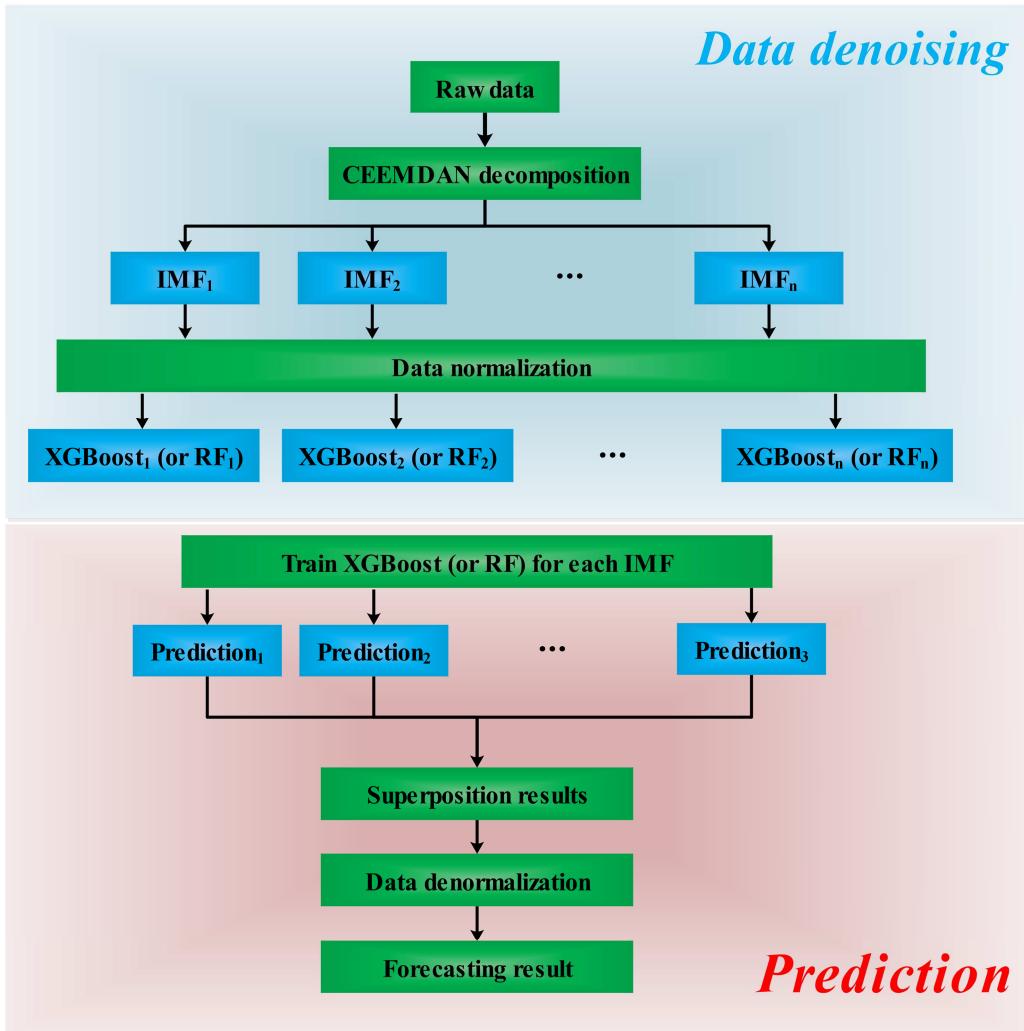
273 (3) Divide data into the training set and test set

274 In this paper, all decomposed datasets are divided into the training set and test set, and their  
275 ratios are 9:1. The sliding window length is 7, that is, the data of the first 6 hours are used to predict  
276 the next one.

277 (4) Prediction

278 XGBoost and RF are used to do the prediction, the prediction results are summarized, then  
279 denormalize the summarized data using Eq. (14) to get the ultimate result, as shown in Fig.1.

280  $f_i = f_n(z_{max} - z_{min}) + z_{min}$  (14)



281

282 **Fig.1.** Data denoising process and prediction.

283

284 *4.2. Error metrics*

285 In this paper, six error metrics are used to evaluate the prediction performance, and their  
 286 expressions are shown in Eqs. (15)-(20). Among them, MAE, RMSE, MAPE, RMSPE are four  
 287 common error evaluation indicators (Ma et al., 2019), and the smaller their values, the smaller the  
 288 error. U1 and U2 respectively represent prediction accuracy and prediction quality. The smaller the  
 289 value, the higher the prediction accuracy and the better the prediction quality. In 1982, Lewis rated  
 290 the prediction performance based on MAPE. The MAPE less than 10% can be considered as  
 291 “excellent”, the MAPE between 10% and 20% can be evaluated as “good”, and the prediction  
 292 performance is “reasonable” when the MAPE is in the range of 20%-50%. If the MAPE is greater  
 293 than 50%, the prediction result is “inaccurate”.

$$294 \quad \text{MAE} = \frac{1}{n} \sum_{t=1}^n |O_t - P_t| \quad (15)$$

295

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{t=1}^n (O_t - P_t)^2} \quad (16)$$

296

$$\text{MAPE} = \frac{100\%}{n} \sum_{t=1}^n \left| \frac{O_t - P_t}{O_t} \right| \quad (17)$$

297

$$\text{RMSPE} = \sqrt{\frac{1}{n} \sum_{t=1}^n \left( \frac{O_t - P_t}{O_t} \right)^2} \quad (18)$$

298

$$U1 = \frac{\sqrt{\frac{1}{n} \sum_{t=1}^n (O_t - P_t)^2}}{\sqrt{\frac{1}{n} \sum_{t=1}^n O_t^2} + \sqrt{\frac{1}{n} \sum_{t=1}^n P_t^2}} \quad (19)$$

299

$$U2 = \sqrt{\frac{\sum_{t=1}^n (O_t - P_t)^2}{\sum_{t=1}^n O_t^2}} \quad (20)$$

300

301 **5. Results and discussions**

302 This section presents the prediction results and error analysis results of six water quality  
 303 indicators, and discusses the stability of the prediction model.

304 *5.1. Results*

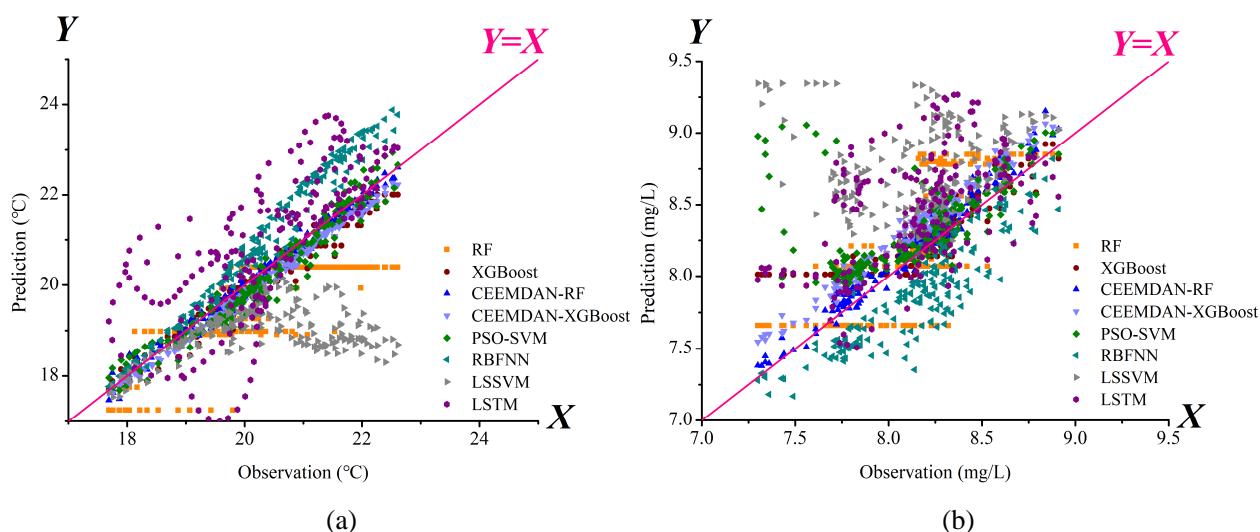
305 Fig.2 shows the prediction results of six water quality indicators using various models, namely  
 306 RF, XGboost, CEEMDAN-RF, CEEMDAN-XGBoost, PSO-SVM, RBFNN, LSSVM, and LSTM.  
 307 The X-axis represents the observation value, and the Y-axis represents the prediction value. In  
 308 general, the ideal prediction results will be distributed over Y=X or evenly distributed on both sides  
 309 of the line, thus indicating that the error is basically obeying the Gaussian distribution. Therefore,  
 310 the closer the point to the Y=X line, the smaller the error. In the prediction of the six indicators, it  
 311 can be clearly seen that the points of RF, LSTM, and RBFNN are far away from Y=X. It can be  
 312 seen from (d) and (f) of Fig.2 that the prediction results of some models are distributed on the same  
 313 side of Y=X, which indicates that the deviation is large, and there may be problems of over-fitting  
 314 or under-fitting. Table 3 lists the prediction errors of each model for each indicator. It can be  
 315 concluded that CEEMDAN-RF or CEEMDAN-XGboost models have the best performance.  
 316 CEEMDAN-RF performs best in the prediction of temperature, dissolved oxygen, and specific

317 conductance, the MAPEs are 0.69%, 1.05% and 0.90%, respectively. CEEMDAN-XGBoost  
 318 performs best in the prediction of pH value, turbidity and FDOM, the MAPEs are 0.27%, 14.94%  
 319 and 1.59%, respectively.

320 In the prediction of turbidity, the MAPEs of the eight models are generally larger. The MAPEs  
 321 of RF, XGboost, CEEMDAN-RF, CEEMDAN-XGBoost, PSO-SVM, RBFNN, LSSVM and LSTM  
 322 models are 18.91%, 19.13%, 18.34%, 14.94%, 44.08%, 24.55%, 21.29% and 24.22%, respectively.  
 323 Moreover, we compared the average values of MAPE and RMSPE of the six water quality  
 324 indicators (without considering other error metrics, because the magnitude of different water quality  
 325 indicators is different, the average values for MAE, RMSE, U1, U2 are not available for reference).  
 326 The results show that the average MAPEs of RF, XGboost, CEEMDAN-RF, CEEMDAN-XGBoost,  
 327 PSO-SVM, RBFNN, LSSVM and LSTM models are 6.41%, 4.60%, 3.90%, 3.71%, 8.96%, 6.93%,  
 328 6.96% and 7.35%, respectively. The RMSPEs of RF, XGboost, CEDAN-XGBoost, PSO-SVM,  
 329 RBFNN, SVM and SVM models are 6.41%, 4.60%, 3.90%, 3.71%, 8.96%, 6.93%, 6.96% and  
 330 7.35%, respectively. Therefore, in general, CEEMDAN-XGBoost has the best prediction  
 331 performance, followed by CEEMDAN-RF.

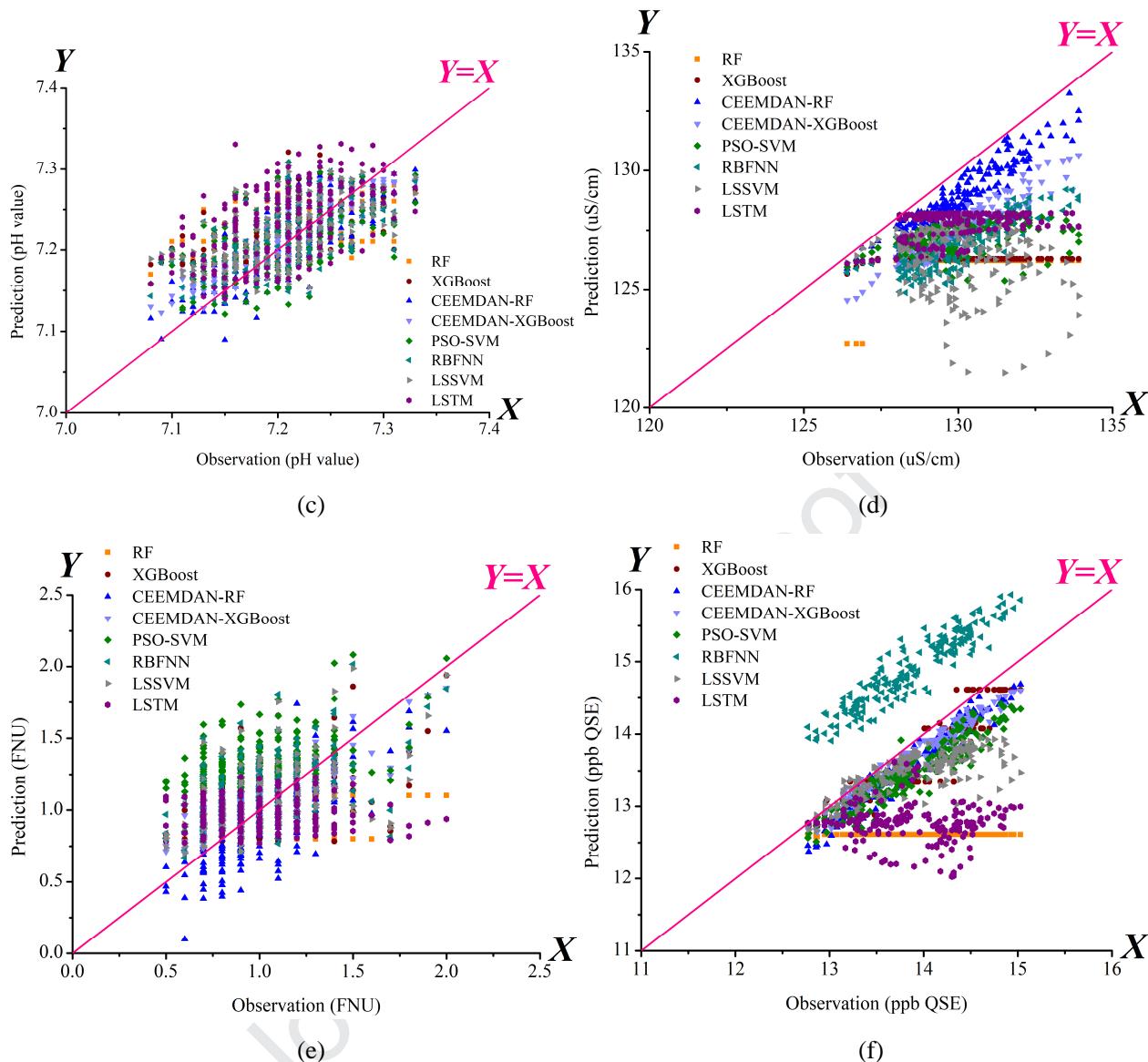
332 For the prediction results of RF, CEEMDAN-RF, XGBoost, and CEEMDAN-XGBoost, the  
 333 prediction results of CEEMDAN-RF are better than RF, and the improvement effect of prediction  
 334 performance is noticeable. The average MAPE of CEEMDAN-RF is 38.10% lower than RF.  
 335 However, some error indicators of CEEMDAN-XGBoost are not as good as XGBoost. For example,  
 336 in the prediction of temperature and dissolved oxygen, the MAPEs of XGBoost are lower than that  
 337 of CEEMDAN-XGBoost. However, on the whole (comprehensive consideration of the results of  
 338 the six water quality indicators), the prediction performance of CEEMDAN-XGBoost is better than  
 339 that of XGBoost because its average MAPE is 19.35% lower than XGBoost.

340



341

342



343

344

(c)

345

(d)

346

347

**Fig.2.** Water quality prediction results (test set). (a) Temperature; (b) Dissolved oxygen; (c) pH value; (d) Specific

348 conductance; (e) Turbidity; (f) FDOM.

349

**Table 3.** Water quality prediction errors.

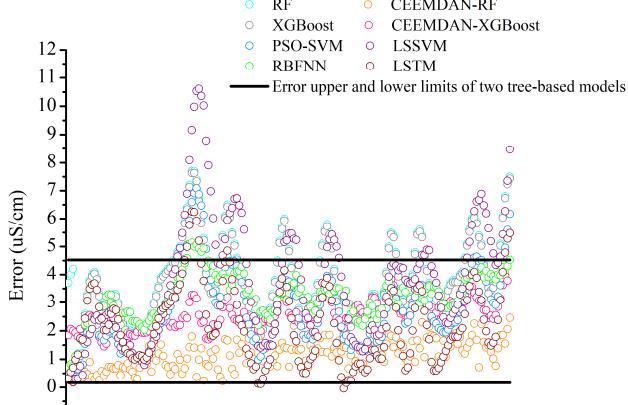
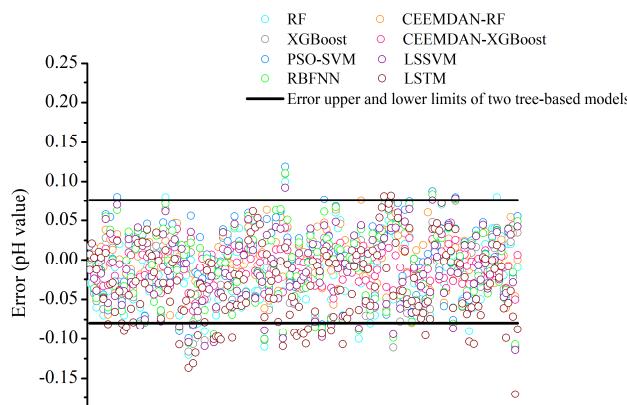
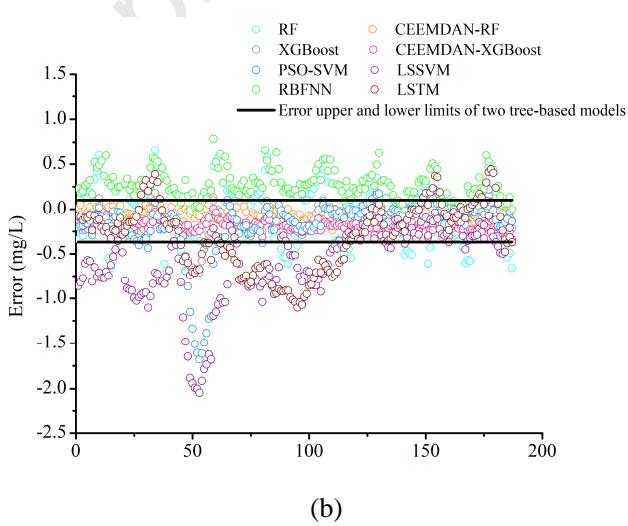
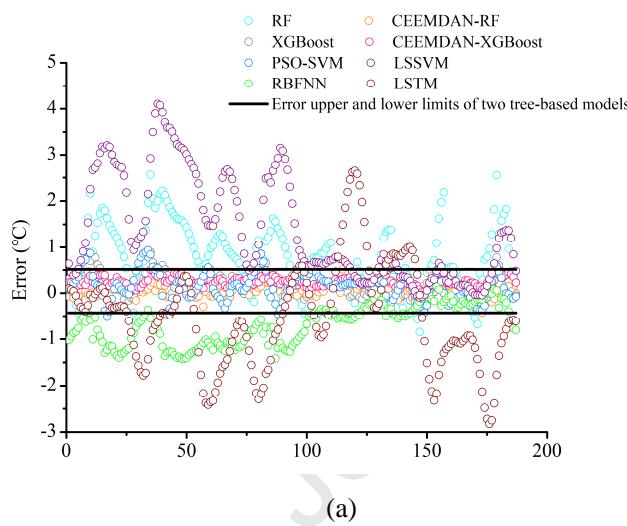
Indicator	Models	Error metrics						Indicator	Models	Error metrics						
		MAE	RMSE	MAPE (%)	RMSPE (%)	U1	U2			MAE	RMSE	MAPE (%)	RMSPE (%)	U1	U2	
Temperature (Unit: $^{\circ}\text{C}$ )	RF	0.85	1.04	4.08	4.97	0.026	0.051	(uS/cm)	Specific conductance	RF	3.83	4.08	2.94	3.12	0.016	0.031
	XGBoost	0.22	0.29	1.07	1.38	0.007	0.014		XGBoost	3.72	3.99	2.85	3.04	0.016	0.031	
	CEEMDAN-RF	<b>0.14</b>	<b>0.17</b>	<b>0.69</b>	<b>0.86</b>	<b>0.004</b>	<b>0.008</b>		CEEMDAN-RF	<b>1.17</b>	<b>1.27</b>	<b>0.90</b>	<b>0.98</b>	<b>0.005</b>	<b>0.010</b>	
	CEEMDAN-XGBoost	0.24	0.26	1.18	1.26	0.006	0.013		CEEMDAN-XGBoost	2.59	2.65	1.99	2.03	0.010	0.020	
	PSO-SVM	0.27	0.34	1.31	1.66	0.008	0.017		PSO-SVM	2.74	3.06	2.10	2.33	0.012	0.024	
	RBFNN	0.69	0.81	3.30	3.82	0.020	0.040		RBFNN	3.22	3.32	2.47	2.55	0.013	0.026	
	LSSVM	1.36	1.77	6.43	8.21	0.045	0.087		LSSVM	3.72	4.32	2.85	3.29	0.017	0.033	
	LSTM	0.98	1.23	4.93	6.23	0.030	0.060		LSTM	2.21	2.60	1.69	1.97	0.010	0.020	
Dissolved oxygen(Unit: mg/L)	RF	0.25	0.30	3.03	3.72	0.019	0.037	Turbidity (FNU)	RF	0.20	0.28	18.91	24.05	0.143	0.2669	
	XGBoost	0.15	0.21	1.96	2.72	0.013	0.026		XGBoost	0.18	0.24	19.13	25.10	0.112	0.2270	
	CEEMDAN-RF	<b>0.09</b>	<b>0.10</b>	<b>1.05</b>	<b>1.26</b>	<b>0.006</b>	<b>0.013</b>		CEEMDAN-RF	0.18	0.23	18.34	22.76	0.114	0.2186	
	CEEMDAN-XGBoost	0.19	0.20	2.30	2.42	0.012	0.024		CEEMDAN-XGBoost	<b>0.13</b>	<b>0.16</b>	<b>14.94</b>	<b>19.18</b>	<b>0.075</b>	<b>0.1534</b>	
	PSO-SVM	0.22	0.38	2.76	5.08	0.023	0.047		PSO-SVM	0.38	0.42	44.08	54.09	0.176	0.4063	
	RBFNN	0.27	0.31	3.34	3.83	0.020	0.039		RBFNN	0.22	0.28	24.55	32.54	0.127	0.2683	
	LSSVM	0.61	0.74	7.62	9.55	0.044	0.091		LSSVM	0.20	0.26	21.29	28.04	0.122	0.2478	
	LSTM	0.38	0.48	4.67	5.96	0.029	0.058		LSTM	0.24	0.31	24.22	31.01	0.152	0.2936	
pH value	RF	0.04	0.05	0.54	0.65	0.003	0.006	FDOM (ppb QSE)	RF	1.26	1.37	8.97	9.65	0.052	0.099	
	XGBoost	0.04	0.05	0.52	0.65	0.003	0.006		XGBoost	0.29	0.33	2.07	2.38	0.012	0.024	
	CEEMDAN-RF	<b>0.02</b>	0.03	0.33	0.41	0.002	0.004		CEEMDAN-RF	0.29	0.32	2.11	2.27	0.011	0.023	
	CEEMDAN-XGBoost	<b>0.02</b>	<b>0.02</b>	<b>0.27</b>	<b>0.33</b>	<b>0.001</b>	<b>0.003</b>		CEEMDAN-XGBoost	<b>0.22</b>	<b>0.25</b>	<b>1.59</b>	<b>1.77</b>	<b>0.009</b>	<b>0.018</b>	
	PSO-SVM	0.04	0.04	0.51	0.62	0.003	0.006		PSO-SVM	0.42	0.47	3.02	3.30	0.017	0.034	
	RBFNN	0.04	0.04	0.50	0.61	0.003	0.006		RBFNN	1.02	1.04	7.39	7.54	0.036	0.075	
	LSSVM	0.04	0.04	0.51	0.62	0.003	0.006		LSSVM	0.44	0.56	3.08	3.90	0.021	0.041	
	LSTM	0.05	0.06	0.65	0.80	0.004	0.008		LSTM	1.12	1.28	7.95	8.98	0.048	0.092	

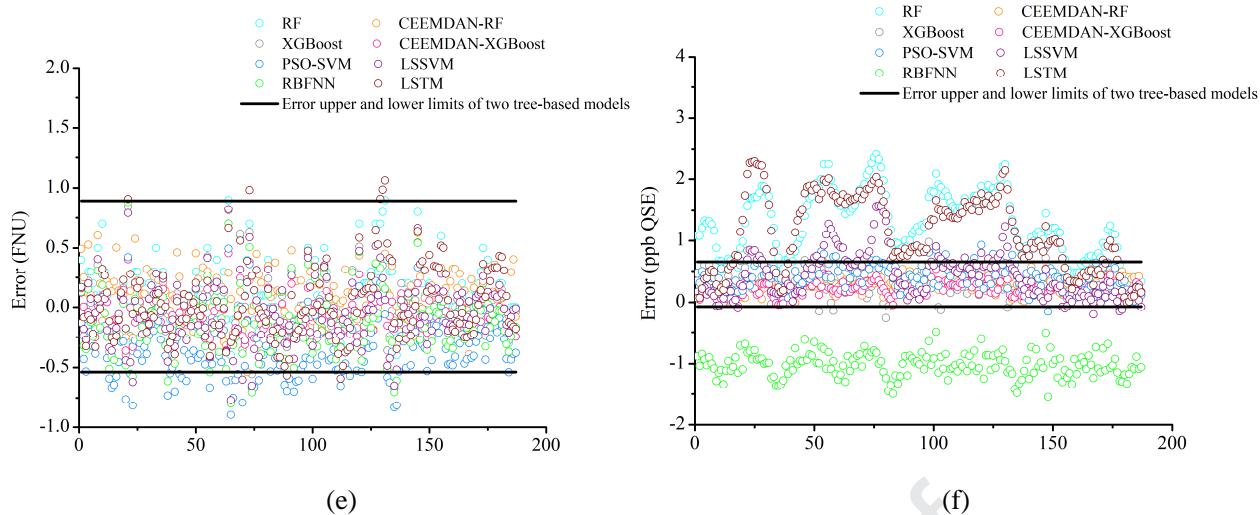
350

Note: Bold represents the data with best performance in the current dataset.

## 351 5.2. Discussions

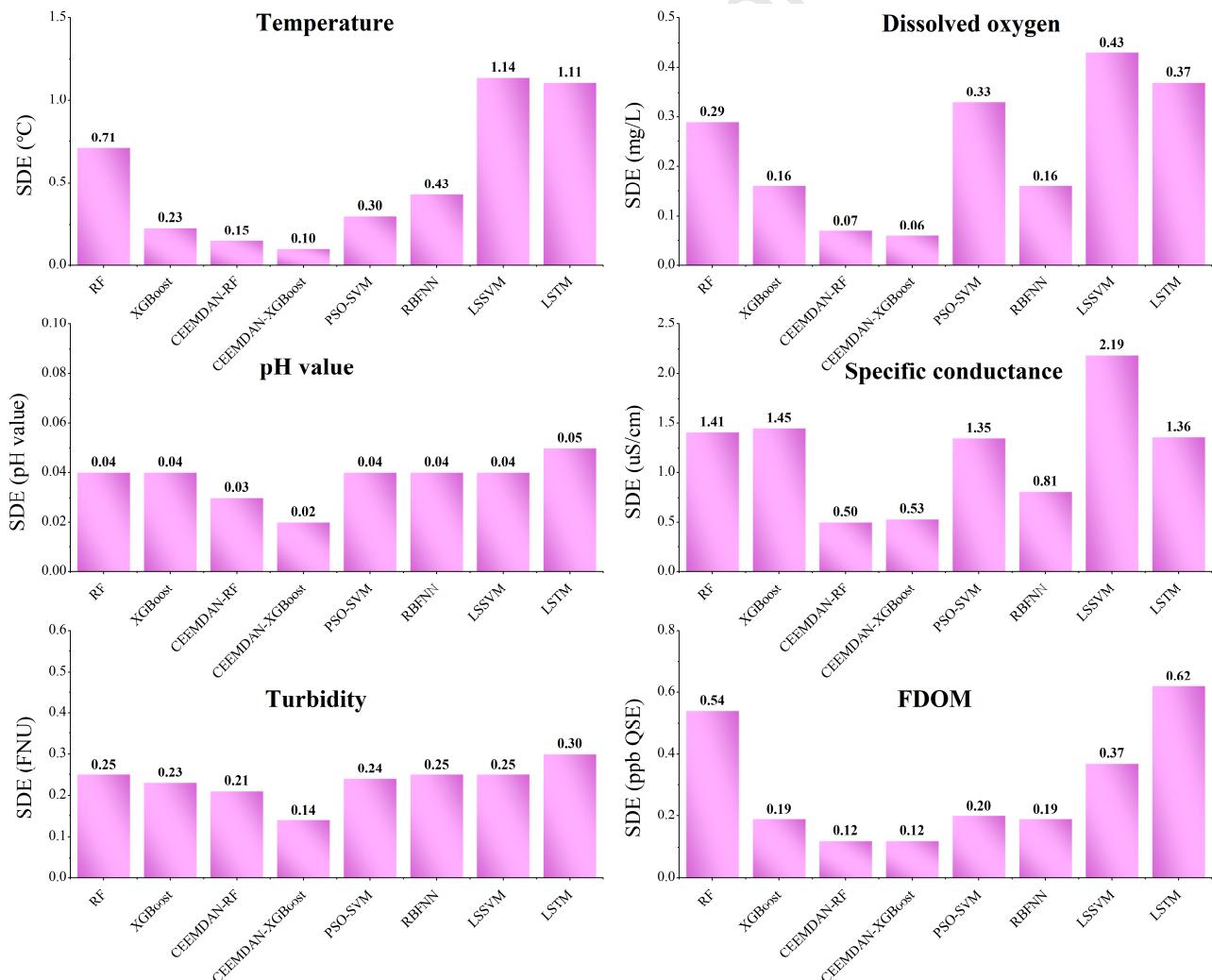
352 Error analysis results can only obtain the overall performance of the prediction. However,  
 353 short-term water quality prediction sometimes requires accurate results at various time points. Some  
 354 models have different characteristics in different datasets. For example, some models have  
 355 difficulty ensuring the accuracy of time series with mutation points. Therefore, the stability of  
 356 prediction is particularly critical. In this paper, the SDE is used as the criterion to evaluate the  
 357 prediction stability of each model. Fig.3 shows the prediction error of each water quality indicator.  
 358 The black lines represent the error upper and lower limits of CEEMDAN-XGBoost model and  
 359 CEMDAN-RF model. It can be seen that many errors of the other six models lie outside the two  
 360 error ranges. Moreover, Fig.4 shows the SDEs of the prediction results of various water quality  
 361 indicators. It can be seen that the SDEs of CEEMDAN-RF and CEEMDAN-XGBoost are smaller  
 362 than the other models, indicating that the stability of the two novel models is better.





367  
368  
369  
370  
371

**Fig.3.** Water quality prediction error (test set). (a) Temperature; (b) Dissolved oxygen; (c) pH value; (d) Specific conductance; (e) Turbidity; (f) FDOM.



372  
373  
374

**Fig.4.** SDEs for water quality indicators.

375 **6. Conclusions and future works**

376 This paper proposes two hybrid decision tree-based models (CEEMDAN-XGBoost and  
 377 CEEMDAN-RF) for the water quality prediction. The CEEMDAN in the two models is used to  
 378 decompose the raw data with large fluctuations, so that the prediction performance of XGBoost and  
 379 RF can be better. This paper takes the water quality of the Gales Creek site of Tualatin River in  
 380 Oregon, USA as the research object, collects the data from May 1st to July 20th, 2019, and divides  
 381 the raw data into training sets and test sets according to the ratio of 9:1. Two models are used to  
 382 predict water temperature, dissolved oxygen, pH value, specific conductance, turbidity, and FDOM,  
 383 and the prediction results are compared with other benchmark models. Besides, this paper takes  
 384 SDE as an evaluation index to analyze the stability of the model. The results indicate:

385 a) CEEMDAN-RF performs best in the prediction of temperature, dissolved oxygen and  
 386 specific conductance, the MAPEs are 0.69%, 1.05% and 0.90%, respectively.  
 387 CEEMDAN-XGBoost performs best in the prediction of pH value, turbidity and FDOM, the  
 388 MAPEs are 0.27%, 14.94% and 1.59%, respectively.

389 b) The average MAPEs of RF, XGboost, CEEMDAN-RF, CEEMDAN-XGBoost, PSO-SVM,  
 390 RBFNN, LSSVM and LSTM models are 6.41%, 4.60%, 3.90%, 3.71%, 8.96%, 6.93%, 6.96% and  
 391 7.35%, respectively. Therefore, in general, CEEMDAN-XGBoost has the best predictive  
 392 performance, followed by CEEMDAN-RF.

393 c) The SDEs of CEEMDAN-RF and CEMDAN-XGBoost is smaller than other benchmark  
 394 models in predicting each water quality indicator, which implies that the stability of the two new  
 395 models is better.

396 Although the models proposed in this paper can already achieve high prediction accuracy, the  
 397 following aspects can be considered in future research: (1) Consider other factors affecting water  
 398 quality in the model, rather than purely time series issues; (2) Due to the high demand of short-term  
 399 prediction on computing time, parallel computing (Li et al., 2019) can be considered in the future.

400 **Declaration of interest**

401 None

402

403 **Acknowledgments**

404 This article is funded by Open Fund of State Key Laboratory of Oil and Gas Reservoir Geology and  
 405 Exploitation (Southwest Petroleum University) (PLN201710), National Natural Science Foundation

406 of China (71901184), Humanities and Social Science Fund of Ministry of Education of China  
 407 (19YJCZH119), and China Scholarship Council (201708030006).

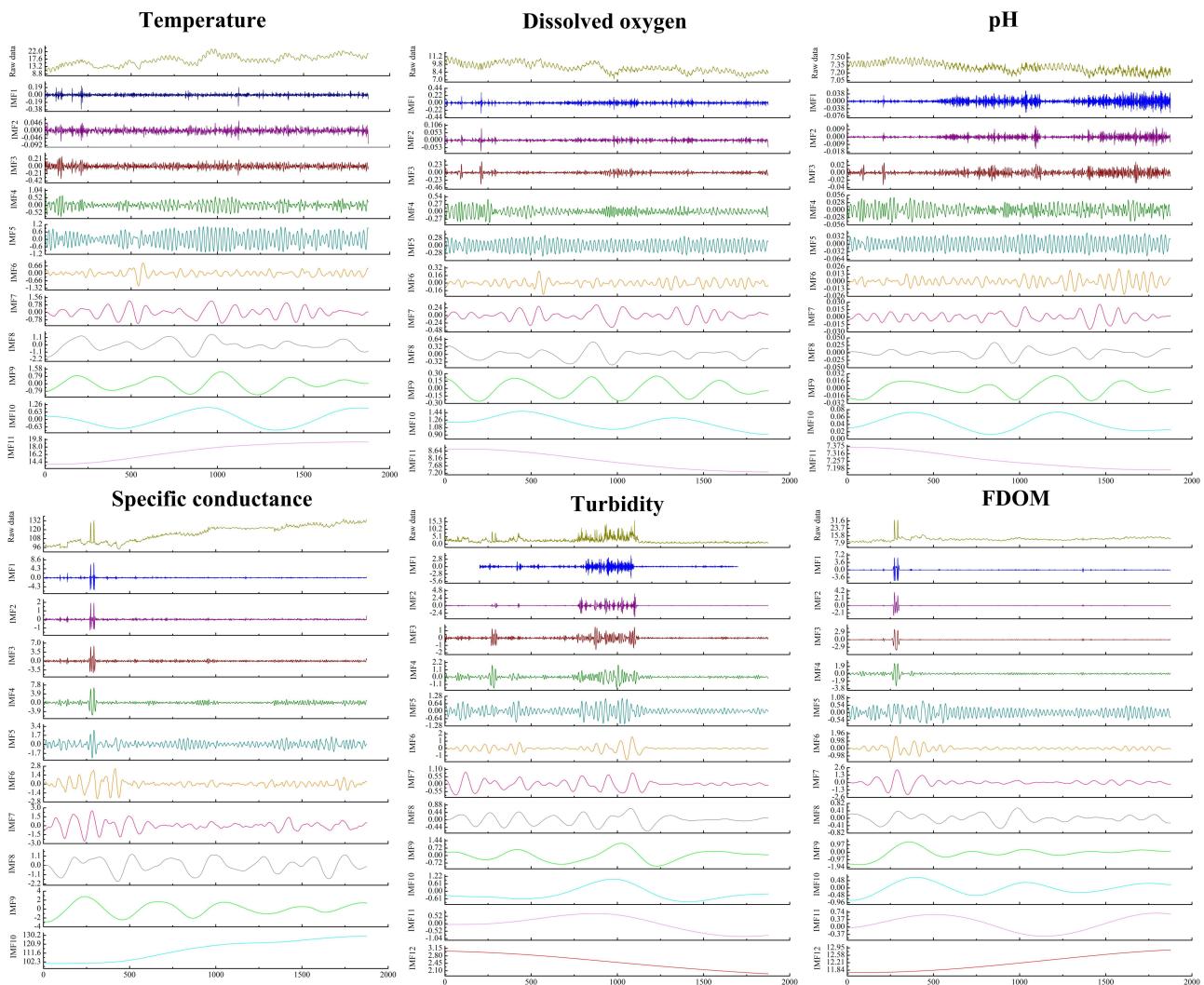
408

409 **References**

- 410 Chan, S. N., Thoe, W., & Lee, J. H. W. (2013). Real-time forecasting of Hong Kong beach water quality by 3D  
 411 deterministic model. *Water research*, 47(4), 1631-1647.
- 412 Chen, T., & Guestrin, C. (2016, August). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd*  
 413 *acm sigkdd international conference on knowledge discovery and data mining* (pp. 785-794). ACM.
- 414 Dabrowski, J. J., Rahman, A., Pagendam, D. E., & George, A. (2020). Enforcing mean reversion in state space  
 415 models for prawn pond water quality forecasting. *Computers and Electronics in Agriculture*, 168, 105120.
- 416 Gao, G., Xiao, K., & Chen, M. (2019). An intelligent IoT-based control and traceability system to forecast and  
 417 maintain water quality in freshwater fish farms. *Computers and Electronics in Agriculture*, 166, 105013.
- 418 García-Alba, J., Bárcena, J. F., Ugarteberu, C., & García, A. (2019). Artificial neural networks as emulators of  
 419 process-based models to analyse bathing water quality in estuaries. *Water research*, 150, 283-295.
- 420 Graf, R., Zhu, S., & Sivakumar, B. (2019). Forecasting river water temperature time series using a wavelet-neural  
 421 network hybrid modelling approach. *Journal of Hydrology*, 578, 124115.
- 422 Hounslow, A. (2018). *Water quality data: analysis and interpretation*. CRC press.
- 423 Huang, P., Trayler, K., Wang, B., Saeed, A., Oldham, C. E., Busch, B., & Hipsey, M. R. (2019). An integrated  
 424 modelling system for water quality forecasting in an urban eutrophic estuary: The swan-canning estuary  
 425 virtual observatory. *Journal of Marine Systems*, 103218.
- 426 Hussein, A. M., Elaziz, M. A., Wahed, M. S. A., & Sillanpää, M. (2019). A new approach to predict the missing  
 427 values of algae during water quality monitoring programs based on a hybrid moth search algorithm and the  
 428 random vector functional link network. *Journal of Hydrology*, 575, 852-863.
- 429 Kerr, J. G., Zettel, J. P., McClain, C. N., & Kruk, M. K. (2018). Monitoring heavy metal concentrations in turbid  
 430 rivers: Can fixed frequency sampling regimes accurately determine criteria exceedance frequencies,  
 431 distribution statistics and temporal trends?. *Ecological Indicators*, 93, 447-457.
- 432 Kong, L., & Ma, X. (2018). Comparison study on the nonlinear parameter optimization of nonlinear grey  
 433 Bernoulli model (NGBM (1, 1)) between intelligent optimizers. *Grey Systems: Theory and Application*, 8(2),  
 434 210-226.
- 435 Larsen, S. J., Kilminster, K. L., Mantovanelli, A., Goss, Z. J., Evans, G. C., Bryant, L. D., & McGinnis, D. F.  
 436 (2019). Artificially oxygenating the Swan River estuary increases dissolved oxygen concentrations in the  
 437 water and at the sediment interface. *Ecological Engineering*, 128, 112-121.
- 438 Li, F., Chen, J., & Wang, Z. (2019). Wireless MapReduce distributed computing. *IEEE Transactions on*  
 439 *Information Theory*, 65(10), 6101-6114.
- 440 Liang, Z., Zou, R., Chen, X., Ren, T., Su, H., & Liu, Y. (2020). Simulate the forecast capacity of a complicated  
 441 water quality model using the long short-term memory approach. *Journal of Hydrology*, 581, 124432.
- 442 Liaw, A., & Wiener, M. (2002). Classification and regression by randomForest. *R news*, 2(3), 18-22.

- 443 Liu, W. X., He, W., Wu, J. Y., Wu, W. J., & Xu, F. L. (2019). Effects of fluorescent dissolved organic matters  
 444 (FDOMs) on perfluoroalkyl acids (PFAAs) in lake and river water. *Science of the Total Environment*, 666,  
 445 598-607.
- 446 Liu, Y., Zheng, Y., Liang, Y., Liu, S., & Rosenblum, D. S. (2016). Urban water quality prediction based on  
 447 multi-task multi-view learning.
- 448 Lu, H., Ma, X., & Azimi, M. (2020). US natural gas consumption prediction using an improved kernel-based  
 449 nonlinear extension of the Arps decline model. *Energy*, 116905.
- 450 Lu, H., Ma, X., Huang, K., & Azimi, M. (2020a). Carbon trading volume and price forecasting in China using  
 451 multiple machine learning models. *Journal of Cleaner Production*, 249, 119386.
- 452 Ma, X., Mei, X., Wu, W., Wu, X., & Zeng, B. (2019). A novel fractional time delayed grey model with Grey Wolf  
 453 Optimizer and its applications in forecasting the natural gas and coal consumption in Chongqing China.  
 454 *Energy*, 178, 487-507.
- 455 Ma, X., Wu, W., Zeng, B., Wang, Y., & Wu, X. (2020). The conformable fractional grey system model. *ISA*  
 456 transactions. DOI: 10.1016/j.isatra.2019.07.009.
- 457 Makarewicz, J. C., Lewis, T. W., Boyer, G. L., & Edwards, W. J. (2012). The influence of streams on nearshore  
 458 water chemistry, Lake Ontario. *Journal of Great Lakes Research*, 38, 62-71.
- 459 Meyers, G., Kapelan, Z., & Keedwell, E. (2017). Short-term forecasting of turbidity in trunk main networks.  
 460 *Water research*, 124, 67-76.
- 461 Mosley, L. M., Peake, B. M., & Hunter, K. A. (2010). Modelling of pH and inorganic carbon speciation in  
 462 estuaries using the composition of the river and seawater end members. *Environmental Modelling &*  
 463 *Software*, 25(12), 1658-1663.
- 464 Nobre, J., & Neves, R. F. (2019). Combining principal component analysis, discrete wavelet transform and  
 465 XGBoost to trade in the financial markets. *Expert Systems with Applications*, 125, 181-194.
- 466 Palani, S., Liong, S. Y., & Tkalich, P. (2008). An ANN application for water quality forecasting. *Marine Pollution*  
 467 *Bulletin*, 56(9), 1586-1597.
- 468 Panidhapu, A., Li, Z., Aliashrafi, A., & Peleato, N. M. (2020). Integration of weather conditions for predicting  
 469 microbial water quality using Bayesian Belief Networks. *Water Research*, 170, 115349.
- 470 Peng, Z., Hu, W., Liu, G., Zhang, H., Gao, R., & Wei, W. (2019). Development and evaluation of a real-time  
 471 forecasting framework for daily water quality forecasts for Lake Chaohu to Lead time of six days. *Science of*  
 472 *The Total Environment*, 687, 218-231.
- 473 Singh, K. P., Basant, A., Malik, A., & Jain, G. (2009). Artificial neural network modeling of the river water  
 474 quality—a case study. *Ecological Modelling*, 220(6), 888-895.
- 475 Tao, Y., Wang, Y., Rhoads, B., Wang, D., Ni, L., & Wu, J. (2020). Quantifying the impacts of the Three Gorges  
 476 Reservoir on water temperature in the middle reach of the Yangtze River. *Journal of Hydrology*, 582,  
 477 124476.
- 478 Torres, M. E., Colominas, M. A., Schlotthauer, G., & Flandrin, P. (2011, May). A complete ensemble empirical  
 479 mode decomposition with adaptive noise. In 2011 IEEE international conference on acoustics, speech and  
 480 signal processing (ICASSP) (pp. 4144-4147). IEEE.

- 481 Valdivia-Garcia, M., Weir, P., Graham, D. W., & Werner, D. (2019). predicted Impact of Climate Change on  
482 trihalomethanes Formation in Drinking Water treatment. *Scientific reports*, 9(1), 9967.
- 483 West, D., & Dellana, S. (2011). An empirical analysis of neural network memory structures for basin water  
484 quality forecasting. *International Journal of Forecasting*, 27(3), 777-803.
- 485 Wu, L. F., Li, N., & Zhao, T. (2019). Using the seasonal FGM (1, 1) model to predict the air quality indicators in  
486 Xingtai and Handan. *Environmental Science and Pollution Research*, 26(14), 14683-14688.
- 487 Wu, L., Huang, G., Fan, J., Zhang, F., Wang, X., & Zeng, W. (2019). Potential of kernel-based nonlinear  
488 extension of Arps decline model and gradient boosting with categorical features support for predicting daily  
489 global solar radiation in humid regions. *Energy conversion and management*, 183, 280-295.
- 490 Xie, M., Wu, L., Li, B., & Li, Z. (2020). A novel hybrid multivariate nonlinear grey model for forecasting the  
491 traffic-related emissions. *Applied Mathematical Modelling*, 77, 1242-1254.
- 492 Yang, W., Wang, J., Niu, T., & Du, P. (2020). A novel system for multi-step electricity price forecasting for  
493 electricity market management. *Applied Soft Computing*, 88, 106029.
- 494 Zhang, P., Ma, X., & She, K. (2019). A Novel Power-Driven Grey Model with Whale Optimization Algorithm  
495 and Its Application in Forecasting the Residential Energy Consumption in China. *Complexity*, 2019.
- 496 Zhao, J., Wang, J., Guo, Z., Guo, Y., Lin, W., & Lin, Y. (2019). Multi-step wind speed forecasting based on  
497 numerical simulations and an optimized stochastic ensemble method. *Applied Energy*, 255, 113833.
- 498 Zhao, Y., Nan, J., Cui, F. Y., & Guo, L. (2007). Water quality forecast through application of BP neural network  
499 at Yuqiao reservoir. *Journal of Zhejiang University-Science A*, 8(9), 1482-1487.
- 500 Zhou, L., Xu, C., Yuan, Z., Lu, T. (2019). Dam Deformation Prediction Based on CEEMDAN-PSR-KELM  
501 Model. *Yellow River*, 41(6), 138-142.

502 **Appendix**

503

504 **Appendix 1.** Decomposition of raw data by CEEMDAN.

505

## Highlights

1. Two hybrid decision tree-based models are proposed to predict the water quality.
2. An advanced denoising method is used to preprocess raw data.
3. The case study was conducted on the most polluted river Tualatin River in Oregon, USA.
4. The prediction stability of the model is analyzed.

**Declaration of interests**

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 paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

