

Accepted refereed manuscript of: Xu G, Fan H, Oliver DM, Dai Y, Li H, Shi Y, Long H, Xiong K & Zhao Z (2022) Decoding river pollution trends and their landscape determinants in an ecologically fragile karst basin using a machine learning model. *Environmental Research*, 214, Art. No.: 113843. <https://doi.org/10.1016/j.envres.2022.113843>
© 2022, Elsevier. Licensed under the Creative Commons Attribution–NonCommercial–NoDerivatives 4.0 International <http://creativecommons.org/licenses/by-nc-nd/4.0/>

1 Decoding river pollution trends and their 2 landscape determinants in an ecologically 3 fragile karst basin using a machine 4 learning model

5 Highlight

6 Spatial and temporal patterns of river water quality in Wu jiang River basin (WRB) were
7 analyzed from 2014 to 2019

8 Machine learning model (XGBoost) was developed to predict robust spatially-distributed
9 continuous water quality patterns

10 SHAP was used as a powerful model interpreter to decode the black box of a ML model
11 indicating the drivers of water quality deterioration

12 Geological and climatic vulnerabilities drive management decisions for control of pollution in
13 these critical areas

14 Abstract

15 Karst watersheds accommodate high landscape complexity and are influenced by both
16 human-induced and natural activity, which affects the formation and process of runoff, sediment
17 connectivity and contaminant transport and alters natural hydrological and nutrient cycling.
18 However, physical monitoring stations are costly and labor-intensive, which has confined the
19 assessment of water quality impairments on spatial scale. The geographical characteristics of
20 catchments are potential influencing factors of water quality, often overlooked in previous studies
21 of highly heterogeneous karst landscape. To solve this problem, we developed a machining learning
22 method and applied Extreme Gradient Boosting (XGBoost) to predict the spatial distribution of
23 water quality in the world's most ecologically fragile karst watershed. We used the Shapley Addition
24 interpretation (SHAP) to explain the potential determinants. Before this process, we first used the
25 water quality damage index (WQI-DET) to evaluate the water quality impairment status and
26 determined that COD_{Mn} , TN and TP were causing river water quality impairments in the WRB.
27 Second, we selected 46 watershed features based on the three key processes (sources-mobilization-
28 transport) which affect the temporal and spatial variation of river pollutants to predict water quality

29 in unmonitored reaches and decipher the potential determinants of river impairments. The predicting
30 range of COD_{Mn} spanned from 1.39 mg/L to 17.40 mg/L. The predictions of TP and TN ranged from
31 0.02 to 1.31 mg/L and 0.25 to 5.72 mg/L, respectively. In general, the XGBoost model performs
32 well in predicting the concentration of water quality in the WRB. SHAP explained that pollutant
33 levels may be driven by three factors: anthropogenic sources (agricultural pollution inputs), fragile
34 soils (low organic carbon content and high soil permeability to water flow), and pollutant transport
35 mechanisms (TWI, carbonate rocks). Our study provides key data to support decision-making for
36 water quality restoration projects in the WRB and information to help bridge the science:policy gap.
37 Keywords: Ecologically fragile karst basin; Water quality assessment; XGBoost regression; Shapley
38 additive explanations; Determinant analysis.

39

40 Introduction

41 Anthropogenic interferences have dramatically hindered natural hydrological and nutrient
42 cycles, in turn threatening river water quality in many countries and regions across the world
43 (Mandaric et al., 2018; Ockenden et al., 2017). Controlling pollutant emissions has become the focus
44 of many global environmental policies (Cardinale, 2011; Mandaric et al., 2018; Vorosmarty and
45 Sahagian, 2000). Water quality can be impacted by anthropogenic factors (such as the land use and
46 land cover changes) (Baker, 2003; Liu et al., 2018; Yan et al., 2021a). Urbanization has led to an
47 increase in impervious surfaces, which alters hydrological flow paths and deliver pollutants to the
48 river network more efficiently resulting in additional pressure and degradation of river water quality
49 (Marinoni et al. 2013). Intensification of agricultural activities may result in increased nutrient loads
50 due to fertilization and changes in surface soil properties.

51 Geographical factors (e.g. climate change, atmospheric deposition, geology and topography,
52 soil types, catchment hydrology, land use/cover and land management) are summarized as three key
53 process factors (i.e., sources, mobilisation and delivery) that define how determinants spatially and
54 temporally affect water quality in a watershed (Alvarez-Cabria et al., 2016a; Fan and Shibata, 2015;
55 Heckmann and Schwanghart, 2013; Lintern et al., 2018; Liu et al., 2021; Noori et al., 2012; Varanka
56 et al., 2015). Identifying the influences of watershed geographical characteristics on river water
57 quality is helpful to understand the evolution of river ecosystem in this region because these key
58 factors vary widely across different geographical regions (Liu et al., 2021; Mainali and Chang, 2018).

59 Karst is globally distributed landscape and supports approximately 20% of the world's
60 population (Ford and Williams, 2007; Hartmann et al., 2014). The Wu jiang River basin (WRB) is
61 located in the world's largest continuous landscapes of karst, which is deemed as an ecological
62 barrier for the Yangtze River Basin and also defined as one of the most ecologically fragile regions
63 in the world (Xu et al., 2021). Land use patterns (e.g. sloped planting and overgrazing) interact with
64 the heterogeneous karst landscape composition and configuration in complex ways (Varanka et al.,
65 2015; Xu et al., 2019). As a result, karst landscapes are more fragile and this can influence the
66 formation and processes of runoff, sediment connectivity and the delivery of pollutants from land
67 to water (Ai et al., 2015; Heckmann and Schwanghart, 2013; Yan et al., 2021b). Further,
68 hydrological processes in karst landscapes deviate from typical responses in non-karst environments
69 and this can lead to river water quality impairments differentiating from those of the plains (Deng,
70 2020; Liu et al., 2020). Due to the influence of subtropical humid monsoon, rainfall in this region
71 is seasonally unevenly distributed and heavy rainfall events are a frequent occurrence, which
72 exacerbates the mobilization and transport of pollutants (Powers et al., 2016; Singh et al., 2005a;
73 Sinha and Michalak, 2016). In recent years, the nutrient balance of the WRB has become a
74 controversial issue due to the construction of cascade dams (Li and Ji, 2016; Winemiller et al., 2016).
75 The geographical characteristics and human disturbance of WRB lead to serious water pollution and
76 complex environmental response in karst areas.

77 Field assessments of water quality can support catchment managers and stakeholders in
78 identifying spatio-temporal sensitive areas of managed landscapes and help to evaluate the benefits
79 and risks of water management strategies in priority areas (Altenburger et al., 2015; Huang et al.,
80 2021; Yi et al., 2017). However, most water quality assessments are limited to particular river
81 reaches due to the costs associated with data collection; therefore, many low-order streams are not
82 evaluated, which can limit understanding of water quality challenges in a watershed (Altenburger
83 et al., 2015; Ding et al., 2016; Mello et al., 2018). Thus, physical process-based model simulation
84 can complement field monitoring investigations. Models, e.g., HSPF and HYPE, SWAT, AGNPS
85 or semi-distributed process based model SPARROW or INCA can simulate complex nonlinear
86 interactions between nutrient transport dynamics and biogeochemical processes (Arhonditsis et al.,
87 2007; Hashemi et al., 2016; Mayorga et al., 2010; Singh et al., 2005b). Such physical process-based
88 models often preclude the identification of dominant processes operating within a watershed due to

89 uncertainties associated with parameter calibration across a large watershed (Badham et al., 2019;
90 Jakeman et al., 2006; Knoben et al., 2020). The complexity of environmental processes often results
91 in physical process-based models being costly and labor-intensive inputs of dataset collection.
92 Moreover, the karst zone under the thin soil layer in karst regions has high permeability and
93 accommodates a complex subsurface hydrological system which makes the parameterization of
94 such models difficult and hinders the transferability of mechanical process approaches to karst areas
95 (Fiorillo et al., 2015; Hartmann et al., 2015; Li et al., 2021; Malago et al., 2016). On the contrary,
96 data-driven machine learning (ML) models are recognized as an effective alternative method and
97 offer advantages for modeling complex nonlinear systems over deterministic and statistical models
98 when handling multi-source data for prediction of river water quality due to improved model
99 interpretability, prediction accuracy, and reduced computational cost (Najah Ahmed et al., 2019;
100 Sun and Scanlon, 2019; Wang et al., 2021b; Zou et al., 2019).

101 As an optimized distributed gradient lifting algorithm, Extreme Gradient Boosting (XGBoost)
102 delivers high accuracy and fast processing time (Lundberg et al., 2020). Indeed, XGBoost
103 outperformed several other machine learning techniques (e.g., Gradient Boosting and Deep Neural
104 Network, Bayesian Regularized Neural Network and Random Forest algorithm) to predict
105 probabilities, and is especially used when dealing with spatial data (Just et al., 2020; Mokoatle et
106 al., 2019). Tree-based ML models are often regarded as unexplainable black box models (Moreira
107 et al., 2020; Parsa et al., 2020). However, data-driven machine learning models suffer from several
108 drawbacks. First, ML models often require a large amount of training data to obtain robust
109 performance (Kratzert et al., 2019). Second they are still not as easily interpretable as traditionally-
110 used physics-based conceptual hydrologic models (Höge et al., 2022). Shapley Additive
111 Explanations (SHAP) is considered as a state-of-the-art model interpretation to decode the black-
112 box of ML models. It can connect optimal credit allocation with local explanations using the classic
113 Shapley values from game theory and their related extensions (Adadi and Berrada, 2018; Lundberg
114 and Lee, 2017; Molnar, 2020). This helps to understand the magnitude and direction of the influence
115 of input variables on the output variable.

116 The overarching aim of our study, therefore, was to investigate the spatial distribution of river
117 water quality impairments in the WRB and decipher how watershed features, both anthropogenic
118 and natural, impair water quality. Firstly, we compiled a complete time series trend (2014–2019)

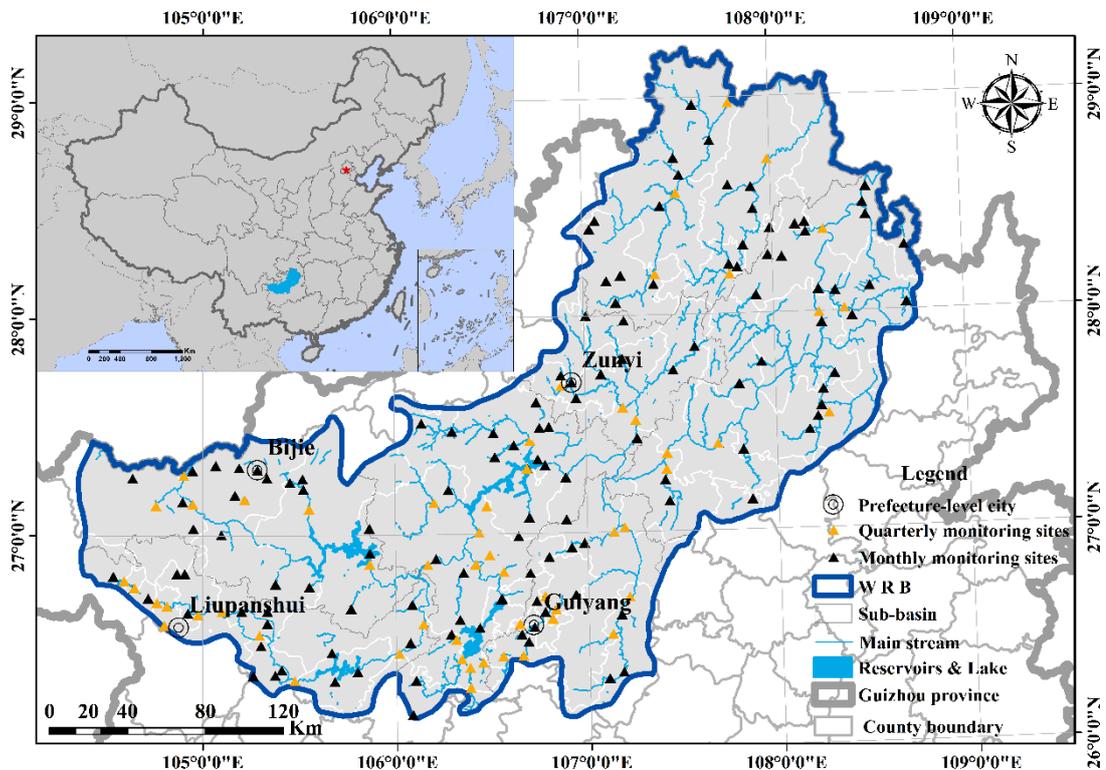
119 dataset of river water quality (14,845 records from 207 water quality sampled sites) to identify
120 spatio-temporal water quality impairments and screen the key variables that contribute to water
121 quality impairment in the WRB. We hypothesized that watershed landscape attributes are important
122 in interpreting water quality at different regional scales in the WRB. Then a powerful ML method
123 was developed to predict water pollution concentrations in unmonitored reaches and we used the
124 SHAP value to determine the significant landscape covariates of water quality in the WRB.

125 2. Materials and methods

126 2.1 Study area

127 The WRB (25°39'13"~25°41'00"N, 105°36'30"-105°46'30"E) is located in southwest China,
128 Guizhou province, which comprises a total area of 80300 km², see Fig.1. Though located in an area
129 of subtropical humid monsoon climate, with average annual precipitation of 1300 mm, there is a
130 serious shortage of clean drinking water for people and livestock (Qin et al., 2015). The WRB
131 provides agricultural irrigation, urban development, river navigation and other functions for more
132 than 35 million people in 54 counties in Guizhou Province (Xu et al., 2021b). Due to the local
133 government promoting strict farmland protection policy, the cultivated land in the WRB have
134 remained stable during the period of the Outline of the 12th Five-Year Plan and the 13th Five-year
135 Plan for National Economic and Social Development of the People's Republic of China. The total
136 use of fertilizers has increased significantly in accordance with the increase of grain demand (Li et
137 al., 2020a; Oliver et al., 2020), together with runoff and infiltration of pollutants leading to a serious
138 crisis of the river water quality in the WRB (Li et al., 2020a; Xu et al., 2021b). Due to the slow soil
139 formation of carbonate rocks, large landform slopes, low vegetation cover, water and soil
140 conservation is at risk from natural disasters and poor approaches to agricultural production have
141 exacerbated soil erosion and rock desertification (Xu et al., 2021b). Because the original surface
142 vegetation was mostly destroyed, the vegetation of the region is mainly a secondary forest,
143 consisting of subtropical evergreen and deciduous broad-leafed mixed trees, mainly composed of
144 species of *genera Cyclobalanopsis, Pinus, Betula, and Cupressus*(Sheng et al., 2018).The bedrock
145 of the WRB is mainly composed of carbonate rocks, dolomite, and limestone micaceous(Han and
146 Liu, 2004), and the main soil types are yellow loam, paddy soil, and calcareous soil. The
147 hydrogeological conditions are complex due to the unique geology of the region, which has
148 contributed to mature underground rivers. The region is covered by shallow soil (<1 m), mainly

149 composed of lime developed from dolomite (>50%), some of which is mineralized and some of
 150 which is presented in loose form (Nie et al., 2017). The soil types in the WRB are more permeable
 151 soils type A and B and may result in higher water tables and accelerated nutrient flow to the soil
 152 (Rodriguez-Galiano et al., 2014). Their texture is silty loam, sandy soils have higher porosity and
 153 therefore lower water retention, resulting in lower absorption of pollutants such as pesticides, metal
 154 ions and solutes. In addition to the aerated structure and inadequate bonding of humus to sand grains,
 155 these properties preferentially allow the infiltration of water and associated contaminants (Andry et
 156 al., 2009; Zalidis et al., 2002).



157
 158 Fig.1. The monitoring river networks and overview of WRB

159 2.2 Data resources and pre-processing

160 2.2.1 Water quality data resources

161 We applied a complete time series water dataset from 207 surface water quality monitoring
 162 sections sampled by the Bureau of Water Resources Department and Environmental Protection
 163 Department in Guizhou province respectively. A monthly sampling frequency is used for national-
 164 controlled and provincial-controlled water quality sections are once per month, while that of water
 165 functional areas is once per quarter. The time scales were across from 2014 to 2019. All the
 166 indicators were collected by mixed samples and tested in laboratory. The Sample pretreatment and

167 pollutant concentration determination methods are mainly based on "Environmental Quality
168 Standard for Surface Water" (GB3838-2002), and specific detection methods are presented in the
169 supplementary material ST3. Therefore, the complete time series pollution indexes in the WRB
170 were selected including Ammonia nitrogen ($\text{NH}_4^+\text{-N}$), Total phosphorus (TP), Five-day biochemical
171 oxygen demand (BOD_5), Dissolved oxygen (DO), Anionic surfactant (AS), Temperature (T),
172 Hydrogen ion concentration index (pH), Electrical conductivity (EC), total nitrogen (TN), sulfide
173 (SO_4^{2-}), Potassium permanganate index (COD_{Mn}). These 11 indexes were used as water environment
174 dataset in our study. We set the measured value below the detection limit as the detection limit
175 (Farnham et al., 2002).

176 2.2.2 Meteorological and landscape data sources and data pre-processing

177 We developed an integrated database including 46 watershed landscape characteristics
178 contributing to the spatial variability of river pollution according to three key driving processes
179 (sources, mobilization and delivery) put forward by multiple studies (Granger et al., 2010;
180 Hrachowitz et al., 2016; Lintern et al., 2018). More details are shown in supplementary material
181 Table.S1. Then we used the hydrological tool "Burn-in" method in ArcGIS 10.2 software with a
182 watershed pixel threshold of 15,000 to delineate 792 reaches and sub-basins of the WRB, and their
183 distribution was adjusted to be more correct according to satellite images of the basin. The
184 topographic wetness index (TWI) was calculated by using the 8-flow method proposed by Quinn
185 (Gruber and Peckham, 2009; Quinn et al., 1991). The grid soil permeability data was computed by
186 the Python-ROSSETA model according to soil texture (Zhang et al., 2018). The average annual
187 streamflow of the reaches is modeled based on the water balance Budyko model (Zhang et al., 2004)
188 through annual predication and potential evapotranspiration, and all the parameters are participated
189 in calculation referred to (Dai et al., 2021). The landscape index is calculated from Fragstats 4.3.
190 The distance between landslide geological disaster points and water body is analyzed by nearest
191 neighbor analysis ArcGIS 10.2. Industrial Point sources emission data came from fifteen thousands
192 sewage draining outlets of Guizhou Provincial Environmental Protection Bureau. Nighttime light
193 data was provided by (Li et al., 2020b). A detailed description of the data sources and processing
194 steps are provided in supplementary materials, see Fig.S1-27 and Table.ST2.

195 2.3. Modeling and database processing

196 First, we assessed water quality conditions and identified key variables deteriorating water

197 quality calculating by the Water Quality Index (WQI-DET) proposed by (Huang et al., 2019).
198 Second, we used Zonal statistics in ArcGIS 10.2 to extract the watershed characteristic data to pair
199 with geographical location of the water quality sampled sites. The integrated subbasin units (SUs)
200 database matrix containing water quality data (as model inputs) and corresponding watershed
201 characteristic data (as model output) prior to developing the prediction model are shown in Fig S1,
202 Table S1 and S2. We first detrended the water quality for use in modeling (Schwarz et al., 2006). A
203 large uncertainty would exist in time-averaged water quality on account of the temporal variability
204 within water quality datasets. Prior to performing ML models, we used the car package in R 4.03 to
205 perform the Box-Cox transformation for the site-level average mean concentrations of each water
206 quality index (Fox et al., 2012; Guo et al., 2019). The Box-Cox parameter λ was estimated
207 individually and presented in supplementary material Table S3 and all the transformed water quality
208 variables were normally distributed based on the Shapiro-Wilk's test (Box and Cox, 1964; Liu et al.,
209 2021; Wang et al., 2021a). We only selected data from 151 water quality sites for ML modeling on
210 account of the completeness of the dataset for total nitrogen. In order to improve the precision and
211 computational efficiency of the ML model, through feature selection, we first removed redundant
212 and irrelevant features according to Spearman correlation analysis, see section 3.2. Spearman
213 correlation coefficient between each pair of features were performed and Mantel test was used to
214 test the relationship between environmental factors and water quality variables (Legendre et al.,
215 2015), see Fig.5. Due to the special geological conditions, debris flow, landslide and other geological
216 disasters often occur in some parts of the WRB. We added the distance between the landslide
217 damage points to the center of water body in the ML model. Although this metric is not filtered into
218 the four models, we still included this index into the inputs of the ML model to verify if it has an
219 impact to the water quality impairment. All the prediction models in this study were performed on
220 the Jupiter notebook platform using the open sources libraries in Python3.7 (Scikit-learn, Hyperopt,
221 XGBoost). The visualization and calculation of SHAP value applied the SHAPforxgboost by Liu
222 (2019) in R 4.03 <https://liuyanguu.github.io/post/2019/07/18/Visualization-of-shap-for-xgboost/>.
223 ArcGIS®10.2 was used for process and analysis of all watershed attributes and visualization of the
224 results.

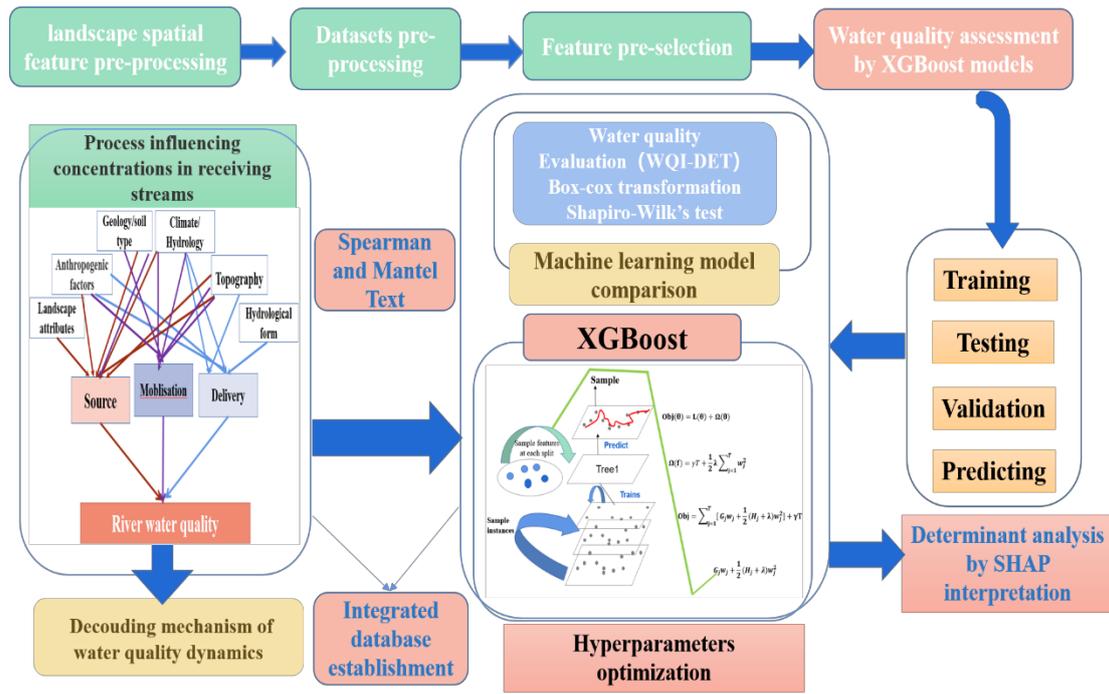


Fig.2. The schematic framework of overall methods used in this study

2.3.1 Water quality impairments evaluation

Water quality can be expressed in terms of scores calculated through integrating complex data into a mathematical expression (Nazeer et al., 2014). Water quality index (WQI) (dimensionless value) based on multiple water quality indicators has been widely used to characterize the degradation degree of surface and groundwater water quality (Lumb et al., 2011; Sutadian et al., 2016; Wu et al., 2018). We referred to the algorithm of modification water quality index (WQI-DET) to determine key variables leading to deterioration of water quality put forward by (Huang et al., 2019). The most sensitive indicators of river water quality impairments were evaluated according to the relative frequency of a variable leading to negative WQI-DET values in the WRB during 2014-2019. Indices of WQI-DET indicate extremely poor water quality (score of -∞) through to good water quality (score of 100). We can calculate the value of a WQI-DET of a single water sample by equation (1), from which monthly WQI-DET was calculated by averaging all the values within a given month. Eleven (11) water quality variables were used to calculate WQI-DET, i.e., $n = 11$, and their concentrations were evaluated against the corresponding surface water quality classes.

$$WQI_{DET}^j = \min(WQI_{DET_1}^j, \dots, WQI_{DET_i}^j, \dots, WQI_{DET_n}^j) \quad (1)$$

$$WQI_{DET-i}^j = 100 - \max\left(0, \frac{C_{ij} - C_i^I}{C_i^V - C_i^I} \times 100\right) \quad (2)$$

(WQI_{DET}^j) is the WQI-DET value for the variable i of the water sample j ; C_{ij} is the concentration of the environmental variable i of the water sample j ; C_i^V and C_i^I are the concentration of the variable i at class I and V according to (GB3838-2002), respectively.

2.3.2 Machine learning prediction method

Boosting regression Tree (Boosting) is a machine learning technique commonly used for regression and classification problems. It generates prediction models in the form of collections of weak prediction models (usually decision trees) and modelling complex phenomena (Friedman, 2001; Strobl et al., 2009). The XGBoost package is an optimized distributed gradient enhancement library that reduces the gradient of the loss function (Chen et al., 2015). The component trees using recursive binary partitioning of predictive variables are chosen to minimize the variance of residuals and segment of all predictive variables, which is considered robust to outliers (Chen and Guestrin, 2016). To be self-contained, we just provide a brief description of the XGBoosting model here, and the detailed equation can be referred to the literature elsewhere (Chen and Guestrin, 2016). Eq (1) describes the training loss and regularization which consists of the two parts of XGBoost's objective function:

$$\text{Obj}(\theta) = L(\theta) + \Omega(\theta) \quad (3)$$

where $L(\theta)$ is the training loss function employing to evaluate the model simulated performance for training data and $\Omega(\theta)$ is the regularization term aiming to control the overfitting of model (Gao et al., 2018). In addition, the complexity of each tree is often computed as the following Eq. (2):

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \quad (4)$$

In Eq (4) w_j is represented by the vector of scores on leaves while T represented by leaves respectively. The Eq.(3) is defined as the objective function of the structure score of XGBoost.

$$\text{Obj} = \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + \Gamma t \quad (5)$$

The form $G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2$ (6) is quadratic and the best w_j to a given structure $q(x)$. In each distinct round of cross-validation we tuned the hyperparameters of the XGBoost model. Grid Search CV was applied to automate the tuning of hyperparameters to determine the optimal value

271 of the given model to satisfy the model generalizability (Moriassi et al., 2007). 70% of the randomly
 272 selected data sets was used as the training set and 30% as the test set. The prediction model is only
 273 established by using the data from the training sets, and the training sets are randomly divided into
 274 five parts by non-repeated sampling. Four of them were used to train the model each time, and the
 275 remaining ones was used to verify the accuracy of the four trained models. The step was repeated
 276 five times until each subset had a chance to be used as the validation set, and the remaining subset
 277 was used as the training set. The average of the five test results was calculated as an estimate of
 278 model accuracy and as a model performance indicator of the model under the implementation of the
 279 five-fold cross-validation. Finally, we validated with the remaining 30% of the test sets. We
 280 evaluated average model predictions performance based on coefficient of determination (R^2), root
 281 mean square error (RMSE), and Nash-Sutcliffe coefficient (NSE) (Nash and Sutcliffe, 1970).

282 2.3.3. SHAP analysis

283 Shapley Additive Explanations (SHAP) is a unified approach to create interpretable machine
 284 learning models. It helps to explain the output of any ML model and to visualize and describe the
 285 complex causal relationship between driving forces and the prediction target (Li et al., 2018). SHAP,
 286 an additive explanation model, is inspired by the theoretically optimal Shapley value of
 287 cooperative game theory, with all the characteristics treated as "contributors" (Lundberg and Lee,
 288 2017; Strumbelj and Kononenko, 2014). Shapely values are determined according to several axioms
 289 to help allocate the contribution fairly for a group N (with N features) . (Lundberg et al., 2020;
 290 Lundberg et al., 2018). A linear function of binary features g is defined based on the following
 291 additive feature attribution method in equation (6):

$$292 \quad \phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F|-|S|-1)!}{|F|!} [f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S)] \quad (6)$$

$$293 \quad g(z') = \phi_0 + \sum_{i=1}^M \phi_i Z_i^1 \quad (7)$$

294 where z' , equals to 1 when a feature is observed, otherwise it equals to 0, and M is the number
 295 of input features. In this study, we apply TreeExplainer proposed by Lundberg and Lee (2017) to
 296 accurately calculate TreeSHAP values of the tree integration models. Hyperparameters tuning of
 297 the XGBoosting model are performed separately in each round of cross-validation, and the overall
 298 RMSE was calculated based on the out-of-sample prediction after cross-validation. The SHAP
 299 values of a given prediction variable and observation value exist differences in the outputs, i.e. the

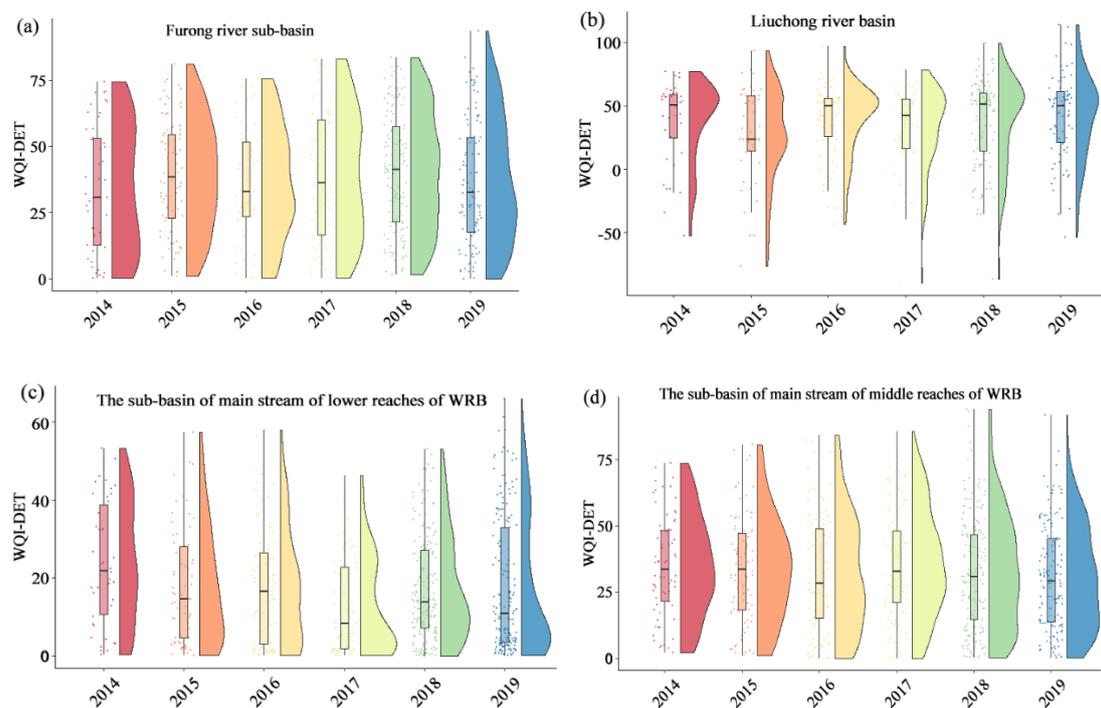
300 predicted water quality indicator, up to whether the model is suitable for using or not using the
301 prediction variable after performing each observation. The mean absolute SHAP values of all
302 observed values summarize the importance of global features, and can be interpreted by more local
303 models through scatter plots of individual predictive variables and their SHAP values (Just et al.,
304 2020).

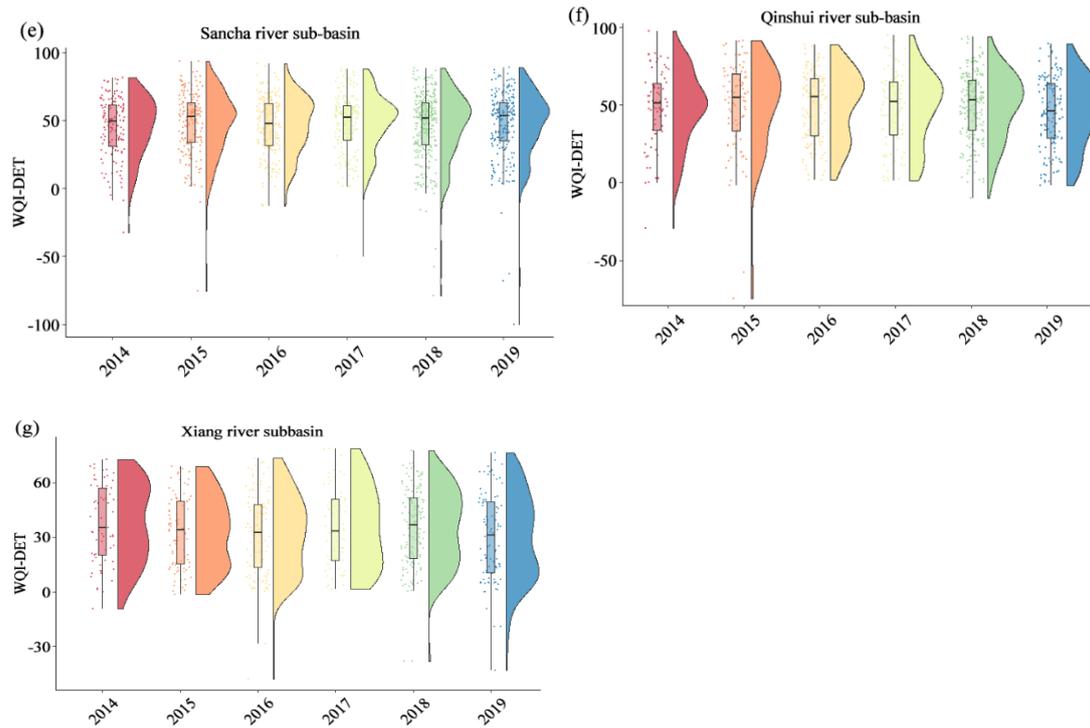
305 **3. Results**

306 3.1 Assessment of water quality impairments in WRB

307 According to the method of section 2.3.1, river water quality in the WRB shows temporal and
308 spatial variation. Water quality is roughly consistent with the distribution of population density and
309 the geographical line of terrain (decreasing from southwest to northeast along the elevation),
310 presenting a complicated characteristic of fractal phenomenon (Xu et al., 2021). To be specific, from
311 2014 to 2019, 37.2% of monitoring sections in Furong River basin in the northeast of the WRB
312 showed good water quality. While only 25.7 % and 17.1 % of monitoring sections in the middle and
313 southwest part of the WRB showed good water quality. Furong River Basin has the best water
314 quality and the mean value of median of WQI-DET is 62.15. While the Sancha River Basin (SCRB),
315 Liuchong River Basin (LCRB) and Qingshui River Basin (QRB) had the worst water quality, with
316 the mean values of median of WQI-DET from 2014 to 2019 being 42.75, 44.15 and 49.81,
317 respectively. The mean value of the median of WQI-DET values of the Xiangjiang River basin
318 (XJRB) and the middle reaches of main stream of the WRB (MRMS of WRB) and lower reaches of
319 main stream of the WRB (LRMM of WRB) were 44.71, 51.65 and 46.65, respectively. The worst
320 water quality monitored section in the WRB was found mainly in the SCRB and LCRB, among
321 which 74.1% (23) of 31 sampled sites and 76.44% (26) of 34 sampled points are worse than Class
322 V on the grounds of water quality standards (GB3838-2002). The water quality of 18 sampling sites
323 (71.76%) in the LCRB was very poor, which was significantly higher than that in central regions
324 (XJRB and LRWRB) and northeast regions (FRB). The sampled sites with extremely poor water
325 quality accounted for 46.1%, 36.4% and 18.19%, respectively. During the year of 2014-2019, the
326 median trend line of WQI-DET (Figure 4) of the FRB, middle reaches and lower reaches of the
327 WRB illustrated a positive slope (k) ($P < 0.01$). It showed a small improvement of water quality in
328 general, increasing 0.52, 0.13 and 0.34 of WQI-DET per year. However, the water quality of the
329 SRB and XRB showed a negative slope, and the overall water quality decreased slightly. The K

330 value of WQI-DET decreased by -1.71 and -1.14 per year respectively. It is worth noting that the
 331 river reaches in the rural areas around the cities, are being seriously polluted. The WQI-DET of
 332 water quality decreased from upstream to downstream reaches of the WRB. The water quality in the
 333 middle reaches of the WRB was slightly better than that in the lower reaches of WRB, which may
 334 be due to the large discharge and the cumulative effect of pollutants from upstream to downstream.
 335 The absolute number of WQI-DET showed a slight downward trend from 2014 to 2019, but it began
 336 to increase after 2017, mainly due to the increase of water quality sampled sites and samples. In all
 337 sampled sites of the whole WRB, 61.3% (127/207) of water quality conditions were seriously
 338 impaired. We used the relative frequency of negative WQI-DET value caused by each variable to
 339 determine the most sensitive indices of river water quality impairment in the WRB during 2014-
 340 2019, and these were COD_{Mn} , TN, and TP. In particular, the contribution of COD_{Mn} (reflecting
 341 organic pollutants) and TP to water quality impairment increased, especially in the river reaches
 342 around densely populated urban and rural residential areas.

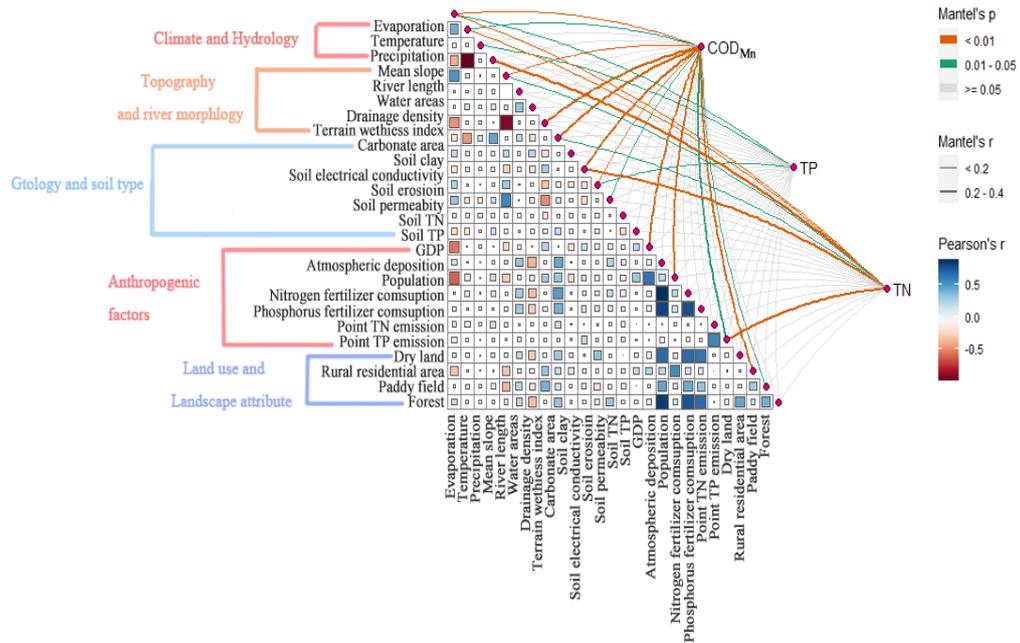




343 Fig.4. Inter-annual variability of the WQI-DET values for seven sub-basins of WRB during the
 344 2014–2019 period. The cloud and rain map represents the WQI-DET distribution in the seven sub-
 345 basins of WRB. X-axis 1-7 represents Furong River Basin (coral pink), Liuchong River Basin
 346 (orange), Qingshui River Basin (straw orange) and the low reaches of main stream of WRB (apple
 347 green), and Xiang River Basin(blue) the middle reaches of main stream of the WRB (olivine). N is
 348 the number of sampling data. Note: To better visualize Fig.4, WQI-DET<-100 was omitted

349 3.2 The results of water quality prediction based on XGBoost

350 The Pearson correlation coefficients between each pair of watershed features were
 351 calculated initially, and only those with a Spearman correlation larger than 0.5 were kept (Fig.5).
 352 The Mantel test of mutual information (Mantel text) is a nonlinear correlation metric for pairs of
 353 geographic characteristics or environmental factors(Legendre et al., 2015). About 50% of features
 354 with low correlation (for COD_{Mn} , TN and TP and $p < 0.05$) were further discarded. Meanwhile, a
 355 Partial Mantel test can eliminate the interference of autocorrelation between environmental factors.
 356 The larger the correlation coefficient of Mantel test, the smaller the P value is. It indicates the greater
 357 the impact of geographical landscape factors on a water quality index(Zeller et al., 2016).



358

359 Fig.5. Pearson correlation between environmental factors is shown in the lower right corner,
 360 and Mantel correlation between watershed geographical factors and water environmental factors is
 361 shown in the upper right corner

362 We used XGBoost to model the relationship between 23 watershed landscape variables and
 363 three water quality indices (Fig.5). Four other popular machine learning techniques were
 364 implemented prior to this work, with adoption of XGBoost, the best predictor, to improve the
 365 performance of machine learning-based water quality predictions. Results are reported in
 366 supplementary material Table S4. According to the Nash-Sutcliffe efficiency coefficient (ranging
 367 from 0.15 to 0.77), the models for COD_{Mn} , and TN were significantly improved after feature
 368 selection, whereas the performance of the TP model was slightly better before feature selection (see
 369 supplementary material Table.S5 for more details). Among the three models, the training datasets
 370 of TP and COD_{Mn} were well fitted in the cross-validation, indicating that these two models have the
 371 highest accuracy. R^2 of CV were 0.68 and 0.57, RMSE were up to 0.79 and 0.94. These two models
 372 were applied to the TP and COD_{Mn} test datasets, with R^2 of CV of 0.71 and 0.65, and RMSE of 0.79
 373 and 0.94, respectively.

374 However, the results of training the model for TN data fell short of expectations, the R^2 of CV
 375 training datasets of TN are 0.37, RMSE were 1.48. The R^2 of the test datasets of TN were 0.39, and
 376 RMSE was 1.37. Through hyperparameter optimization, the TP model was slightly improved while

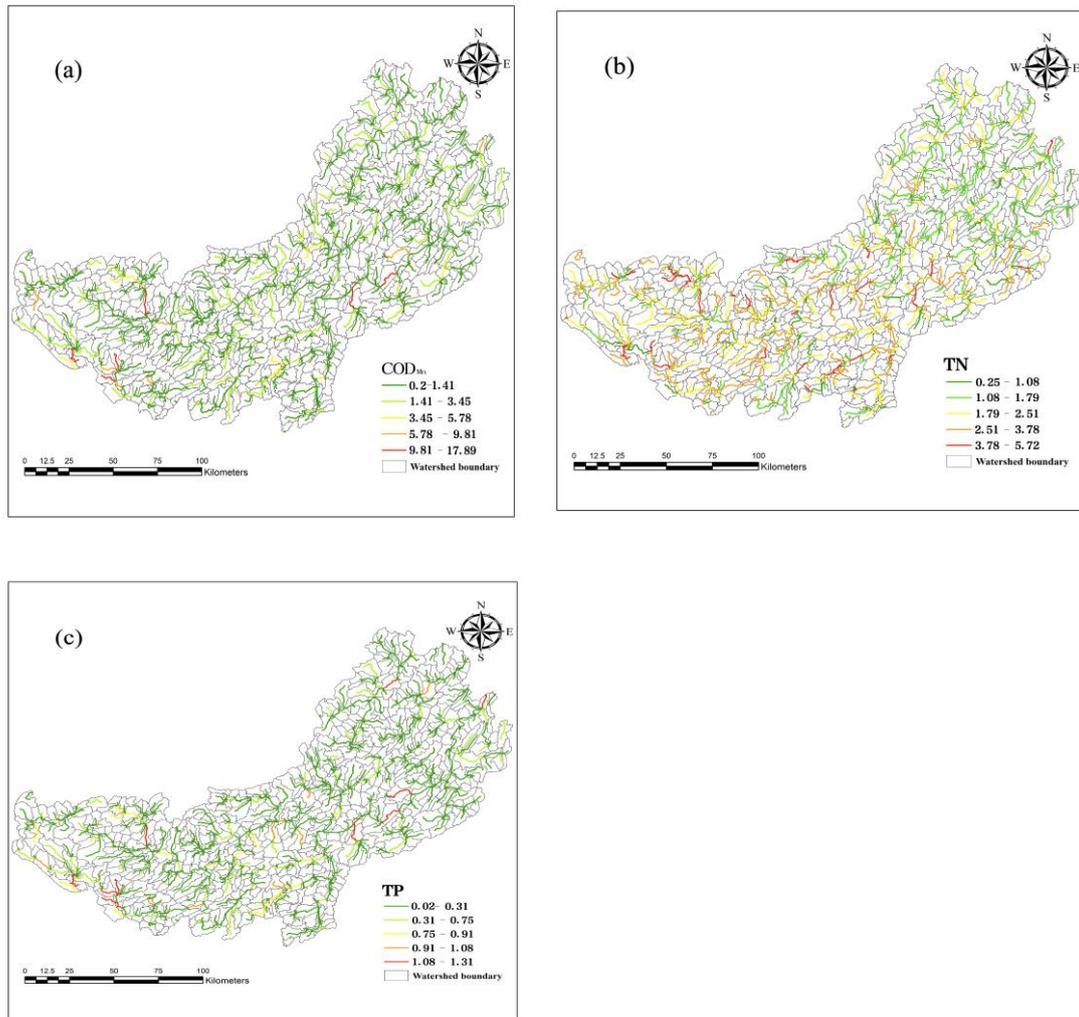
377 the COD_{Mn} and TN were greatly improved. The R² of COD_{Mn} training dataset and test dataset were
 378 increased to 0.67 and 0.73, RMSE was reduced to 0.79 mg/L and 0.65mg/L. R² of TP model training
 379 datasets and training datasets were 0.79 and 0.81, RMSE of which were 0.54 mg/L and 0.66 mg/L.
 380 Nash-Sutcliffe efficiency coefficients of those three models were improved after feature selection
 381 and parameter optimization ranging from 0.54 to 0.77, more details are provided in Table 1 and
 382 supplementary material (Table. S5 and Fig.S3).

383 Table 1. Performance of XGBoost models before and after hyperparameter optimization

Water quality index		R ²					
RMSE		NSE		Train		Test	
		Train	Test	Train	Test	Train	Test
	COD _{Mn}	0.57	0.65	1.96	0.94	0.37	0.45
Default	TN	0.37	0.39	1.48	1.37	0.41	0.23
Parameters	TP	0.68	0.71	0.93	0.79	0.46	0.51
	COD _{Mn}	0.67	0.73	0.79	0.65	0.64	0.77
	TN	0.57	0.61	0.81	0.78	0.54	0.61
Optimized	TP	0.79	0.81	0.54	0.66	0.77	0.74

384 Note: WQI = water quality index, R² = coefficient of determination, NSE =Nash-Sutcliffe efficiency coefficient,
 385 RMSE= root mean square error

386 The concentrations of COD_{Mn}, TN and TP of 792 SUBs were predicted by the XGBoost model.
 387 The predictions show that COD_{Mn} concentration ranges from 0.2 to 17.31 mg/L, with an average
 388 concentration of 15.84 mg/L, See Fig. 6 (a) to (c). The reaches with higher COD_{Mn} concentration
 389 were distributed in densely populated urban reaches of QSRB, XJRB and the MRMS. The TN and
 390 TP concentrations in the WRB ranged from 0.25 to 5.72 mg/L and 0.02 to 1.31 mg/L, with a mean
 391 concentration of 3.83 mg/L and 0.56 mg/L respectively. The central and southeast portions of WRB
 392 are the most contaminated, with significant amounts of TN and TP, which is consistent with the
 393 spatial distribution of agricultural non-point source losses documented in this watershed (Dai et al.,
 394 2021; Xu et al., 2021a).

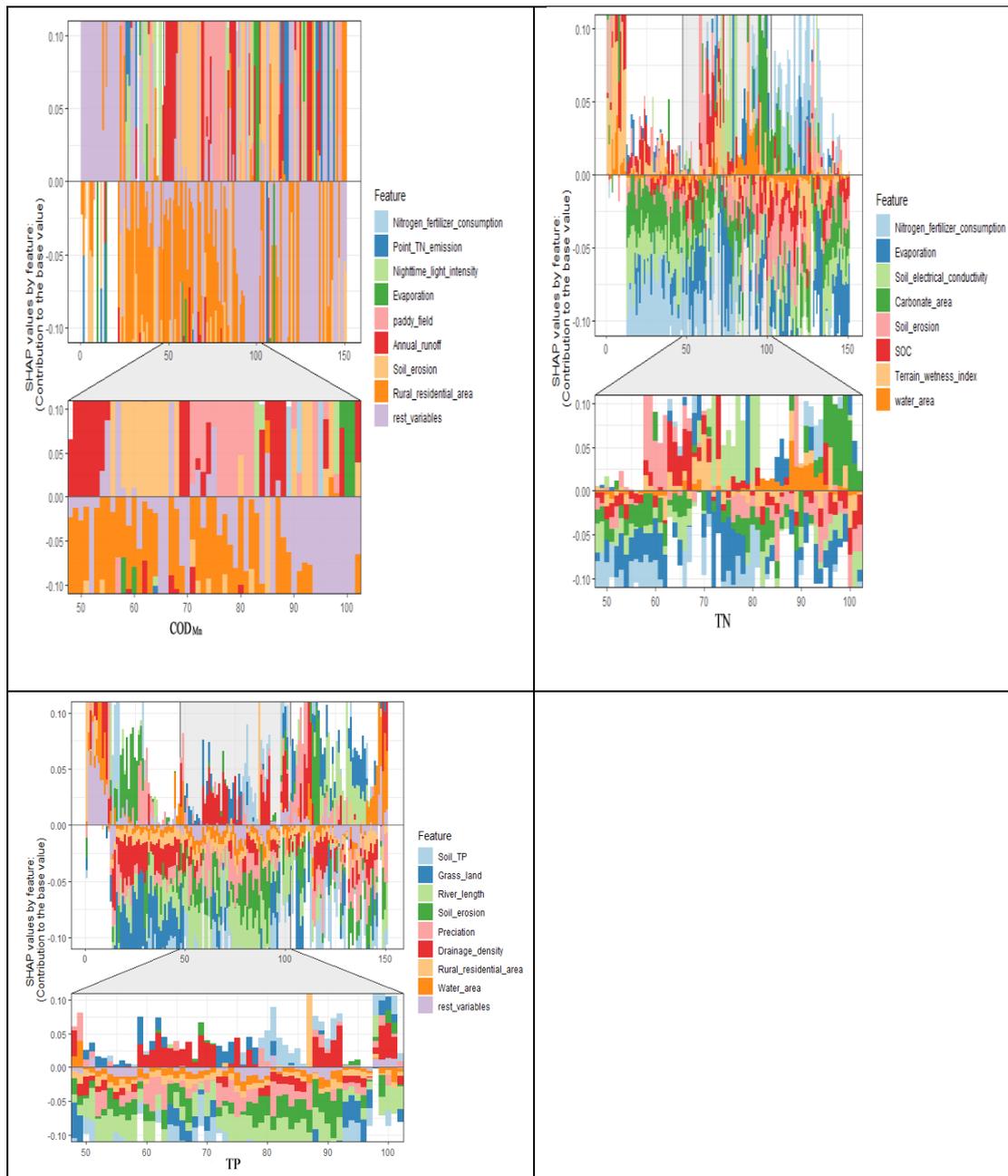


395 Fig.6. The XGBOOST models projected the following four water quality parameters: (a)COD_{Mn},
 396 (b)TN, and (c)TP

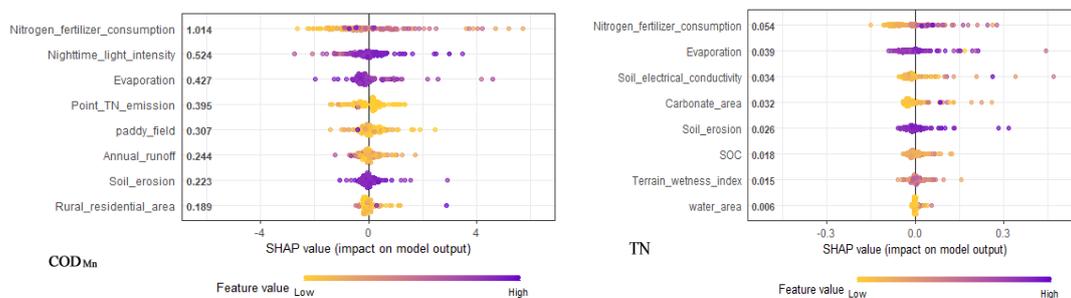
397 3.3 Analysis of determinants of water quality

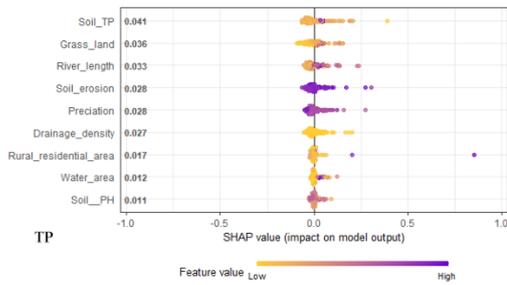
398 We used a 5-fold cross-validation split to evaluate the average absolute SHAP value as a
 399 measure of global feature importance. By performing each round of cross-validation, a recursive
 400 stepwise procedure was employed to order and remove features by increasing importance. Variable
 401 selection was run in three water quality databases using feature importance from SHAP values. Eight
 402 key features leading to water quality deterioration were selected to draw SHAP force plots according
 403 to six clusters of sub-groups for COD_{Mn}, TN and TP models, see our supplementary material
 404 (Fig.S4). The SHAP values of COD_{Mn} ranged from 0.189 to 1.014, The SHAP values of TN and TP
 405 ranged from 0.006 to 0.054 and 0.011 to 0.041. We then pooled the features from the previous
 406 ranking and sorted the importance of the features from lowest to highest according to the average

407 absolute SHAP value of all the features in the model. Repeating this step, least important features
408 in each step were discarded. After plotting the overall RMSE predicted by cross-validation based
409 on the chosen features, we finally selected the model with the lowest RMSE for each of the three
410 water qualities. The SHAP Force plot (Fig.7) essentially superimposed these SHAP values for each
411 observation and shows how the final output is obtained as a sum of the attributes of each predictive
412 variable. The X-axis is set to -1 to 1 to facilitate the comparison of the three models. The Y-axis
413 shows the order of the average absolute value of all observations (Fig.7). The eigenvalue is of
414 absolute SHAP value is higher, the influence of the eigenvalue on the model output is greater. We
415 used a bee swarm plot to illustrate and rank the watershed factors driving water quality in the average
416 absolute SHAP value. COD_{Mn} was driven by anthropogenic factors, and the average absolute SHAP
417 values were: (1.014) of nitrogen fertilizer consumption and (0.524) of night light intensity, and
418 (0.395) of point source nitrogen emission. The land use types such as the paddy land, dry land, grass
419 land and rural residential areas are sorted by COD_{Mn}, TP, and the mean absolute SHAP value ranged
420 from 0.036 to 0.307. Descriptive scatter plots representing watershed characteristics and their SHAP
421 scores are provided in supplementary material (Fig.S5), approximating their contribution (local
422 feature importance) to the prediction of the Y-axis (three water quality characteristics). As important
423 meteorological and hydrological factors, rainfall, evaporation and runoff drive the variation of TP,
424 COD_{Mn} and TN, the average absolute of SHAP values were within the range of 0.028 to 0.427.
425 River morphology factors such as river length, drainage density and water area play a pivotal role
426 in influencing river water quality. Lithologic features (e.g., carbonate rocks) and soil property (soil
427 erosion, Terrain wetness index (TWI), soil permeability and soil electrical conductivity, soil organic
428 carbon) are also important determinants that can influence the deterioration of river water quality,
429 the range of the mean absolute SHAP values were from 0.028 to 0.323 (Fig.8). The descriptive
430 scatter plots representing watershed characteristics and their SHAP scores are provided in
431 supplementary material (Fig.S4), approximating their contributions (local feature importance) to the
432 prediction of the Y-axis (three water quality characteristics).



433 Fig.7. SHAP force plots show how the final output is the sum of the attributes of each predictive
 434 variable





435 Fig.8. The Bees warm plots show SHAP values for watershed characteristics of observation
 436 using each water quality indicator. The Y-axis represents the rank of the average absolute SHAP
 437 values of the observed values (COD_{Mn}, TN, TP) of all watershed features

438 4 Discussion

439 4.1 The performance of the model to predict pollutant concentration in river water quality of WRB

440 The XGBOOST model developed in our study constructed a nonlinear mapping relationship
 441 between the multi-source data and the concentrations of COD_{Mn}, TN and TP. This provided accurate
 442 prediction of the concentrations in unmonitored reaches of the river network. The predicted
 443 concentration of pollutants in the study area is consistent with the measured results, the average R²
 444 were higher than 0.78. It is worth noting that the input variables used to construct the model in this
 445 study were available to obtain. However, the model, constructed based on these input variables, can
 446 successfully predict river pollution. This provides a solution to time constraints imposed by sample
 447 collection, transportation and detection of traditional river water pollution concentrations via
 448 analysis methods, but also solves the problem that traditional monitoring methods cannot conduct
 449 rapid on-site analysis (Shuhong et al., 2019). More importantly, the predictions (Fig.S4-S7) show
 450 that the model constructed in this study can predict the concentration of pollutants in rivers better
 451 than the other 4 ML methods. Of course, the watershed characteristics used in the model will greatly
 452 affect the model performance of different water quality parameters, especially when the water
 453 quality parameters have different sources and migration/transformation processes (Alvarez-Cabria
 454 et al., 2016b). The watershed characteristics adopted by our model can explain the variability of TP
 455 and COD_{Mn} concentrations, but are inferior when predicting TN concentrations. Due the leakage of
 456 carbonate, groundwater systems can act as a net sink for dissolved Nitrogen by increasing the
 457 residence time and reducing the loads of TN through biochemical processes (Zhang et al., 2020B),
 458 thereby reducing the loads from surface water. Meanwhile, TN leaking from carbonate aquifers can

459 be stored and converted to other nitrogen forms (e.g., by nitrification and/or N_2 to NO_3^-N by
460 anaerobic ammonia treatment) (Dai et al., 2021; Zhang et al., 2020). However, weak statistical
461 significance does not necessarily mean that those variables represented by watershed landscape
462 characteristics are inherently unimportant in determining the sources, dynamics and transport of
463 pollutants. Secondly, the accuracy of the model is not only affected by the predictors, but also
464 impacted by the environmental behavior of the predicted targets. The determination of parameters
465 is based not only on the statistical significance of the coefficients, but also on the overall model
466 fitting and the physical importance of the parameters (Wang et al., 2021b). At the same time, the
467 ML model will have better simulation performance when the sample size increases. Or it might be
468 possible that as we do not consider the effects of seasonal variations on water quality resulting in
469 the characteristics considered, our study could not fully explain the water quality impairments.
470 Moreover, the diversification and spatial differences of various water quality parameters are
471 determinants of different ecological, socio-economic and policy influences in the basin that can
472 contribute to uncertainties in the model accuracy.

473 4.2 How does natural characteristics and anthropologic factors generate covariances on water
474 quality in the WRB?

475 In our study, to obtain a more unbiased model, influential watershed characteristics and
476 variables are considered as much as possible which include variables that are not readily available
477 e.g. industrial point source pollution, fertilizer application data, sewage treatment plants, etc. River
478 water quality is covaried with many geographical factors such as topography, land cover,
479 biogeochemical reactivity, climate etc. In section 3.3, we apply SHAP values to decipher how these
480 factors shape the water quality of rivers. Temperature has been identified as a key factor that directly
481 influences riverine thermal regimes and biogeochemical processes, such as nitrification,
482 denitrification, ammonification, and sediment diagenesis rates (Lintern et al., 2018; Sardans et al.,
483 2008; van Vliet et al., 2013). High temperature conditions will increase the growth and degradation
484 of algae and the capacity of sediments to adsorb phosphorus, which leads to higher COD_{Mn}
485 concentration and higher concentrations of total phosphorus in water (Xia et al., 2015). Rainfall can
486 profoundly modulate the flow-concentration relationship (Green et al., 2007), especially during a
487 few short but intense precipitation events, where particulate matter and bioavailable phosphorus

488 loads may differ by an order of magnitude between wet and dry conditions (Long et al., 2014).
489 Environmental land use conflict caused by land use deviating from land capacity (natural use) is the
490 root cause of accelerated deterioration of water quality, and it may lead to continuous changes in
491 precipitation-runoff-infiltration processes, which in turn lead to extensive soil erosion and nutrient
492 loss (Blevins et al., 1998; Pacheco et al., 2018; Suescun et al., 2017; Thomas et al., 2016; Valle
493 Junior et al., 2014). In karst areas, as agriculture encroaches on natural lands, rapid land cover
494 change often leads to long-term damage to soil and water conservation and other important
495 ecosystem services (Li et al., 2021) Low natural vegetation covers owing to improper land use
496 practices (cropping on sloping land), livestock grazing, and environmental hazards (rocky
497 desertification) may reduce contaminant attenuation in karst area (Jiang et al., 2014).

498 In highly permeable carbonate karst aquifers, where there are widespread formations of fissures,
499 fractures, and conduits, fast (e.g. conduit) and slow (e.g. fracture and matrix) flow transfer pathways
500 will operate (Clifford and Williams, 2007). This leads to the rapid infiltration of rainwater that
501 carries pollutants (e.g., from livestock, domestic and industrial discharge effluents) and
502 contaminates groundwater (Wang et al., 2020; Yue et al., 2019). Geology and soil type determine
503 the sources of sediment and natural nutrients in the catchment (Bostanmaneshrad et al., 2018;
504 Grayson et al., 1997; Juracek and Ziegler, 2009). The erodibility of soil and rock and the adsorption
505 capacity of soil affect the flow of water components in a watershed. The mobilization of sediments
506 is closely correlated to the susceptibility of the geological deposit and the soil within the catchment
507 to erosion and weathering (Meybeck et al., 1990; Perry and Vanderklein, 2009). Lithology
508 determines the alkalinity (pH) and conductivity of water and the concentration of different ions
509 associated with many biogeochemical processes (Doherty et al., 2014). Soil adsorption capacity also
510 affects nutrient mobilization in catchments. The transport of dissolved phosphorus, nitrogen and
511 salts from catchments to recipient waters via underground flow pathways is strongly influenced by
512 the hydrological characteristics of the soil. When the soil saturated conductivity of aquifers in
513 watershed areas is low, the residence time of dissolved components in groundwater flow in the
514 catchment area is increased (Lintern et al., 2018). This provides more opportunity for components
515 to be lost from flow paths through nutrient absorption or biogeochemical processes such as
516 denitrification (Hasani Sangani et al., 2015). The soil permeability (soil hydraulic properties) can
517 affect soil quality and moisture, thereby altering the input of nutrients or organic matter to

518 groundwater and river systems (Rodriguez-Blanco et al., 2015). TWI reflects topographic control
519 of groundwater surface and soil moisture, while high TWI values indicate shallow groundwater
520 table and high soil moisture (Rodhe and Seibert, 1999). Soil organic nitrogen is the main source of
521 nitrates in rivers, soil moisture can promote the production of $\text{NO}_3\text{-N}$ from the nitrification of soil
522 organic nitrogen. This can explain that TWI entered TN model (Li et al., 2019). Soil pH plays an
523 important role in determining the morphology of phosphate in soils because phosphate can bind to
524 different iron when pH changes. Phosphates tend to form insoluble compounds in the presence of
525 high concentrations of exchanged calcium. A reduction in soluble phosphorus usually occurs at
526 higher pH (Sierra et al., 2017). As a result, high soil pH reduces water transport of phosphorus from
527 land to rivers.

528 Anthropogenic activities has altered the river morphological conditions (e.g., changes in
529 hydrological connectivity due to dam construction) and can severely impair river water quality
530 (Maavara et al., 2020; Rodriguez-Blanco et al., 2015). In our study, (river length, drainage density,
531 water areas) are all important covariables affecting water quality in TN, TP models. In general,
532 higher drainage density may increase the likelihood that terrestrial pollutants carried by surface
533 runoff will enter the water body (Alexander et al., 2002; Prasad et al., 2005). The river length
534 determines that the transportation time of pollutants in the stream follows first-order reaction
535 kinetics related to hydraulic residence time (Smith et al., 1997). Reservoirs play an important
536 ecological function by hydrologically connecting upland and downstream river networks and
537 influencing the biological cycle of nutrients. They have strong nutrient removal/interception
538 capabilities which can be sinks of incoming nutrients or, if water quality is poor, they become
539 sources of pollutants in downstream river reaches. Due to the need to improve engineering water
540 shortage and flood control, the local government, most rivers in WRB are impounded (Dai et
541 al.,2020). According to the Bureau of Hydrology and Resources of Guizhou Province, besides a
542 dam cascade, there are 19,652 small reservoirs in the entire WRB (Dai, 2019). Dam cascade and
543 small reservoirs have altered the hydrological regime, river morphology and lateral connectivity and
544 increased longitudinal fragmentation of the basin (Viaroli et al., 2018), which has further amplified
545 the instability of the biogeochemical processes and extended the range of resulting environment
546 damage. These structural modifications have also increased regional hydraulic retention times and
547 slowed the flow rate of rivers, in turn hindering river metabolism, amplifying nutrient transport and

548 delivery, but also triggering eutrophication in rivers themselves (Dodds, 2006; Nizzoli et al., 2018).
549 GDP, industrial point emission, rural residential area and night light index were considered as
550 key factors result in COD_{Mn} deterioration of water quality, which may also compensate for the index
551 of population and urban development were filtered by XGBOOST. In the past decade, the WRB has
552 made great efforts to promote the construction of municipal sewage treatment and sewage discharge
553 standards has been strictly enforced under the background of China promoting huge investments to
554 total environmental restoration (Xu et al.,2021). However, there is still a considerable gap in the
555 design principle and operation performance due to treatment facilities and the sewer system lagging
556 behind. Effluent discharge standards and sludge disposal are severely inconsistent with local
557 conditions and environmental requirements (Lu et al., 2019; Qu et al., 2019). It might be the reason
558 that COD_{Mn} concentrations in many urban reaches of WRB were higher than the acceptable limits.
559 Moreover, rural residential area was an important determinant of water quality (COD_{Mn} and TP).
560 Due to pursuing economic development of rural areas and agricultural intensification, the demand
561 and consumption of water has been increasing and in turn runoff from fields and farms has increased
562 in accord with the increases of discharge of domestic sewage, animal waste, leachate from manure
563 storage facilities or green feed (Skinner et al., 1997). Moreover, the buildings in rural areas are
564 spatially scattered, and the high construction costs are very unfavorable for the construction of
565 public water supply and sewage treatment systems (Kupiec et al., 2021). In addition, karst rural
566 areas not only lack the knowledge of proper manure management, but also lack proper manure
567 storage facilities or poor technical standards (Gao et al., 2014; Norse and Ju, 2015; Oliver et al.,
568 2020).

569 4.3 Management implications and future challenges

570 However, deterioration of water quality can be caused by many factors, such as complex
571 geographical environment and intensive human intervention (including mining, intensive
572 agriculture activity), inadequate sewage treatment measures (Xu et al.,2021) and poor groundwater
573 environment (Li et al., 2020a; Zeng et al., 2020). In addition, damming and the construction of
574 multiple small reservoirs have drastically reduced surface runoff, limiting the river's ability to dilute
575 effluent from sewage treatment plants, see section 4.2. As mentioned above, our findings support
576 the development of strategies by identifying key characteristics of pollutant sources and
577 incorporating them into regional planning (e.g. changing land use, improving industrial structure

578 and distribution). In addition, since the Chinese government has been promoting ecological
579 rehabilitation projects to restore rocky desertification and improve local poverty, the river water
580 quality has been neglected in the WRB. We also suggest that soil, water processes and environmental
581 effects should be incorporated into a unified scientific management framework to best communicate
582 the trade-offs between policy options and promotion of pollution control and ecological restoration.
583 This will help realize ecological value and promote green development management of the WRB
584 (Xu et al., 2021a). Our approach can not only effectively promote pollutant sources control, but also
585 decelerate the pollutant migration and transformation process. It is imperative to adjust local
586 economic structure and develop low-pollution water-saving industry. Water-saving irrigation
587 schemes also appear to be a necessary measure to reduce pollutant infiltration into the soil.

588 And some micro-policy proposals were advocated; promoting BMPs is a good choice in this
589 case, it will allow policy makers to mitigate non-point source pollutants and further restore river
590 ecosystems in the agricultural areas of WRB. BMPs include improving the efficiency of fertilization,
591 improving manure management and buffering the pollutant delivery processes between land and
592 water (e.g. restrict livestock farming near rivers, plant more vegetation near river banks).
593 Promoting soil remediation is also important for restoration of the water environment of WRB and
594 requires management of vulnerable geological areas with well-drained soils, high recharge and low
595 soil organic carbon characteristics. It is necessary to implement integrated management of surface
596 water and groundwater to alleviate the contradiction between intensive water use and geographical
597 environment constraints.

598 Although results have been achieved using ML methods to detect and evaluate water quality,
599 we still need to consider the potential disconnects between macro-scale simulations and local social,
600 economic, and environmental realities, as well as catchment-scale constraints for on-site water
601 quality management. But we also need to conduct field assessments to assess the extent to which
602 reductions in pollutants concentrations are actually achieved, based on best management practices
603 for site-specific nutrient sources combining the landscape characteristics (Jarvie et al.,
604 2018; Sharpley et al., 2016). In the short term, it may be unrealistic to expect pollutant
605 concentrations to be reduced to the compliance and restricted target concentrations, especially in
606 highly impaired karst basins with multiple complex pollutant sources and long-term legacy nutrient
607 contributions (Jarvie et al., 2018; Sharpley et al., 2013; Xu et al., 2021a). However, our simulation

608 to assess the nutrient limitations combined with assessment of compliance and limitation gaps,
609 provides a basis for developing targeted approaches to nutrient water quality compliance in future
610 work.

611 5 Conclusion

612 Understanding of the multiple forces determining river water quality and the complexity and
613 interaction of these forces is necessary to develop successful water quality management strategies.
614 Those knowledges can be used to develop predictive models that will help to predict river water
615 quality. In this study, we evaluated those important factors affecting the spatio-temporal variation
616 of water quality (COD_{Mn},TN,TP) in an ecologically fragile watershed with high landscape
617 heterogeneity by adopting a data-driven machine learning approach. Machine learning can take
618 advantage of all the crossover effects between variables to improve the accuracy of model
619 predictions, which is an advantage over traditional statistical models. the Nash efficiency coefficient
620 are ranging from 0.54 to 0.8, which indicates that our prediction is reliable and robust. Through the
621 analysis of powerful model interpreter (SHAP), though anthropogenic factors such as land use are
622 closely related to river pollutant concentrations, the effects of key hydroclimatic, soil types and
623 vegetation conditions vary across different components and regions. XGBOOST can be used to
624 identify potential water quality hot spots in unmonitored locations; this suggests that catchments
625 with steep gradients, fragile soils or areas with widespread carbonate rocks should be sampled more
626 frequently. Our study underlines the needs to highlight soil and water processes and integrate
627 environmental effects into a unified scientific management framework when implementing
628 ecological engineering restoration in karst areas. Therefore, as more land management surveys are
629 been promoting and ongoing water quality monitoring data are available, an extended temporal or
630 spatio-temporal modeling framework may be used to assess the success of recovery measures in the
631 future. In the meanwhile, we should consider combining the assessment of simulated nutrient
632 limits with the assessment of compliance and limitation gaps to provide a basis for developing a
633 targeted approach to river water quality compliance that focuses on closing the gap between current
634 and target concentrations.

635 Reference

636 Adadi, A., Berrada, M., 2018. Peeking Inside the Black-Box: A Survey on Explainable Artificial

637 Intelligence(XAI).IeeeAccess.652138-52160.[https://doi.org/10.1109/Access.2018.](https://doi.org/10.1109/Access.2018.2870052)
638 2870052

639 Ai, L., et al., 2015. Spatial and seasonal patterns in stream water contamination across
640 mountainous watersheds: Linkage with landscape characteristics. Journal of Hydrology.
641 523 398-408.<https://doi.org/10.1016/j.jhydrol.2015.01.082>

642 Alexander, R. B., et al., 2002. Estimating the sources and transport of nutrients in the Waikato
643 River Basin, New Zealand. Water Resources Research. 38 4-1-4-
644 23.<https://doi.org/Artn1268>
645 [10.1029/2001wr000878](https://doi.org/10.1029/2001wr000878)

646 Altenburger, R., et al., 2015. Future water quality monitoring—adapting tools to deal with
647 mixtures of pollutants in water resource management. Science of the total environment.
648 512 540-551.<https://doi.org/10.1016/j.scitotenv.2014.12.057>

649 Alvarez-Cabria, M., et al., 2016a. Modelling the spatial and seasonal variability of water quality
650 for entire river networks: Relationships with natural and anthropogenic factors. Science
651 of the Total Environment. 545 152-162.<https://doi.org/10.1016/j.scitotenv.2015.12.109>

652 Alvarez-Cabria, M., et al., 2016b. Modelling the spatial and seasonal variability of water quality
653 for entire river networks: Relationships with natural and anthropogenic factors. Sci Total
654 Environ. 545-546 152-62.<https://doi.org/10.1016/j.scitotenv.2015.12.109>

655 Andry, H., et al., 2009. Water retention, hydraulic conductivity of hydrophilic polymers in sandy
656 soil as affected by temperature and water quality. Journal of Hydrology. 373 177-
657 183.<https://doi.org/10.1016/j.jhydrol.2009.04.020>

658 Arhonditsis, G. B., et al., 2007. Eutrophication risk assessment using Bayesian calibration of

659 process-based models: Application to a mesotrophic lake. *Ecological Modelling*. 208
660 215-229.<https://doi.org/10.1016/j.ecolmodel.2007.05.020>

661 Badham, J., et al., 2019. Effective modeling for Integrated Water Resource Management: A
662 guide to contextual practices by phases and steps and future opportunities.
663 *Environmental Modelling & Software*. 116 40-
664 56.<https://doi.org/10.1016/j.envsoft.2019.02.013>

665 Baker, A., 2003. Land use and water quality. *Hydrological Processes*. 17 2499-
666 2501.<https://doi.org/10.1002/hyp.5140>

667 Blevins, R., et al., 1998. Conservation tillage for erosion control and soil quality. Chapter in"
668 *Advances in Soil and Water Conservation*". Ann Arbor Press, Chelsea, MI.
669 <https://doi.org/10.1201/9781315136912-4>

670 Bostanmaneshrad, F., et al., 2018. Relationship between water quality and macro-scale
671 parameters (land use, erosion, geology, and population density) in the Siminehrood
672 River Basin. *Sci Total Environ*. 639 1588-
673 1600.<https://doi.org/10.1016/j.scitotenv.2018.05.244>

674 Box, G. E. P., Cox, D. R., 1964. An Analysis of Transformations. *Journal of the Royal Statistical*
675 *Society Series B-Statistical Methodology*. 26 211-252.[https://doi.org/10.1111/j.2517-](https://doi.org/10.1111/j.2517-6161.1964.tb00553.x)
676 [6161.1964.tb00553.x](https://doi.org/10.1111/j.2517-6161.1964.tb00553.x)

677 Cardinale, B. J., 2011. Biodiversity improves water quality through niche partitioning. *Nature*.
678 472 86-9.<https://doi.org/10.1038/nature09904>

679 Chen, T., Guestrin, C., Xgboost: A scalable tree boosting system. *Proceedings of the 22nd acm*
680 *sigkdd international conference on knowledge discovery and data mining*, 2016, pp.

681 785-794. <https://doi.org/10.1145/2939672.2939785>

682 Chen, T., et al., 2015. Xgboost: extreme gradient boosting. R package version 0.4-2. 1 1-4

683 Clifford, F., Williams, P., 2007. Karst Hydrogeology & Geomorphology. [https://doi.org/ 10.1002/](https://doi.org/10.1002/)

684 [9781118684986](https://doi.org/10.1002/9781118684986)

685 Dai, Y. B., et al., 2021. Modelling the sources and transport of ammonium nitrogen with the

686 SPARROW model: A case study in a karst basin. Journal of Hydrology. 592

687 125763. <https://doi.org/ARTN12576310.1016/j.jhydrol.2020.125763>

688 Deng, X. J., 2020. Influence of water body area on water quality in the southern Jiangsu Plain,

689 eastern China. Journal of Cleaner Production. 254

690 120136. <https://doi.org/ARTN120136>

691 [10.1016/j.jclepro.2020.120136](https://doi.org/10.1016/j.jclepro.2020.120136)

692 Ding, J., et al., 2016. Influences of the land use pattern on water quality in low-order streams

693 of the Dongjiang River basin, China: A multi-scale analysis. Science of the Total

694 Environment. 551 205-216. <https://doi.org/10.1016/j.scitotenv.2016.01.162>

695 Dodds, W. K., 2006. Eutrophication and trophic state in rivers and streams. Limnology and

696 Oceanography. 51 671-680. https://doi.org/10.4319/lo.2006.51.1_part_2.0671

697 Doherty, J. M., et al., 2014. Hydrologic Regimes Revealed Bundles and Tradeoffs Among Six

698 Wetland Services. Ecosystems. 17 1026-1039. [https://doi.org/10.1007/s10021-014-](https://doi.org/10.1007/s10021-014-9775-3)

699 [9775-3](https://doi.org/10.1007/s10021-014-9775-3)

700 Fan, M., Shibata, H., 2015. Simulation of watershed hydrology and stream water quality under

701 land use and climate change scenarios in Teshio River watershed, northern Japan.

702 Ecological Indicators. 50 79-89. <https://doi.org/10.1016/j.ecolind.2014.11.003>

703 Farnham, I. M., et al., 2002. Treatment of nondetects in multivariate analysis of groundwater
704 geochemistry data. *Chemometrics and Intelligent Laboratory Systems*. 60 265-
705 281.[https://doi.org/10.1016/S0169-7439\(01\)00201-5](https://doi.org/10.1016/S0169-7439(01)00201-5)

706 Fiorillo, F., et al., 2015. A model to simulate recharge processes of karst massifs. *Hydrological*
707 *Processes*. 29 2301-2314.<https://doi.org/10.1002/hyp.10353>

708 Ford, D., Williams, P., *Karst Hydrogeology and Geomorphology*. 2007.
709 <https://doi.org/10.1002/9781118684986>

710 Fox, J., et al., 2012. Package 'car'. Vienna: R Foundation for Statistical Computing.

711 Friedman, J. H., 2001. Greedy function approximation: A gradient boosting machine. *The*
712 *Annals of Statistics*. 29 1189-1232.<https://doi.org/10.1214/aos/1013203451>

713 Gao, M. F., et al., 2014. Modeling nitrogen loading from a watershed consisting of cropland and
714 livestock farms in China using Manure-DNDC. *Agriculture Ecosystems & Environment*.
715 185 88-98.<https://doi.org/10.1016/j.agee.2013.10.023>

716 Granger, S. J., et al., Chapter 3 - Towards a Holistic Classification of Diffuse Agricultural Water
717 Pollution from Intensively Managed Grasslands on Heavy Soils. *Advances in*
718 *Agronomy*. Academic Press, 2010, pp. 83-115. [https://doi.org/10.1016/S0065-](https://doi.org/10.1016/S0065-2113(10)05003-0)
719 [2113\(10\)05003-0](https://doi.org/10.1016/S0065-2113(10)05003-0)

720 Grayson, R. B., et al., 1997. Preferred states in spatial soil moisture patterns: Local and
721 nonlocal controls. *Water Resources Research*. 33 2897-
722 2908.<https://doi.org/10.1029/97wr02174>

723 Gruber, S., Peckham, S., 2009. Land-surface parameters and objects in hydrology.
724 *Developments in Soil Science*. 33 171-194 <https://doi.org/10.1016/S0166->

725 2481(08)00007-X

726 Guo, D., et al., 2019. Key Factors Affecting Temporal Variability in Stream Water Quality. Water
727 Resources Research. 55 112-129.<https://doi.org/10.1029/2018wr023370>

728 Han, G. L., Liu, C. Q., 2004. Water geochemistry controlled by carbonate dissolution: a study
729 of the river waters draining karst-dominated terrain, Guizhou Province, China.
730 Chemical Geology. 204 1-21.<https://doi.org/10.1016/j.chemgeo.2003.09.009>

731 Hartmann, A., et al., 2015. A large-scale simulation model to assess karstic groundwater
732 recharge over Europe and the Mediterranean. Geoscientific Model Development. 8
733 1729-1746.<https://doi.org/10.5194/gmd-8-1729-2015>

734 Hartmann, A., et al., 2014. Karst water resources in a changing world: Review of hydrological
735 modeling approaches. Reviews of Geophysics. 52 218-
736 242.<https://doi.org/10.1002/2013rg000443>

737 Hasani Sangani, M., et al., 2015. Modeling relationships between catchment attributes and river
738 water quality in southern catchments of the Caspian Sea. Environ Sci Pollut Res Int.
739 22 4985-5002.<https://doi.org/10.1007/s11356-014-3727-5>

740 Hashemi, F., et al., 2016. Review of scenario analyses to reduce agricultural nitrogen and
741 phosphorus loading to the aquatic environment. Sci Total Environ. 573 608-
742 626.<https://doi.org/10.1016/j.scitotenv.2016.08.141>

743 Heckmann, T., Schwanghart, W., 2013. Geomorphic coupling and sediment connectivity in an
744 alpine catchment — Exploring sediment cascades using graph theory. Geomorphology.
745 182 89-103.<https://doi.org/10.1016/j.geomorph.2012.10.033>

746 Höge, M., et al., 2022. Improving hydrologic models for predictions and process understanding

747 using Neural ODEs. Hydrol. Earth Syst. Sci. Discuss. 2022 1-
748 29.<https://doi.org/10.5194/hess-2022-56>

749 Hrachowitz, M., et al., 2016. Transit times-the link between hydrology and water quality at the
750 catchment scale. Wiley Interdisciplinary Reviews: Water. 3 629-
751 657.<https://doi.org/10.1002/wat2.1155>

752 Huang, J., et al., 2019. How successful are the restoration efforts of China's lakes and
753 reservoirs? Environ Int. 123 96-103.<https://doi.org/10.1016/j.envint.2018.11.048>

754 Huang, J., et al., 2021. Characterizing the river water quality in China: Recent progress and on-
755 going challenges. Water Res. 201
756 117309.<https://doi.org/10.1016/j.watres.2021.117309>

757 Jakeman, A. J., et al., 2006. Ten iterative steps in development and evaluation of environmental
758 models. Environmental Modelling & Software. 21 602-
759 614.<https://doi.org/10.1016/j.envsoft.2006.01.004>

760 Jarvie, H. P., et al., 2018. Phosphorus and nitrogen limitation and impairment of headwater
761 streams relative to rivers in Great Britain: A national perspective on eutrophication. Sci
762 Total Environ. 621 849-862.<https://doi.org/10.1016/j.scitotenv.2017.11.128>

763 Jiang, Z. C., et al., 2014. Rocky desertification in Southwest China: Impacts, causes, and
764 restoration. Earth-Science Reviews. 132 1-
765 12.<https://doi.org/10.1016/j.earscirev.2014.01.005>

766 Juracek, K. E., Ziegler, A. C., 2009. Estimation of sediment sources using selected chemical
767 tracers in the Perry lake basin, Kansas, USA. International Journal of Sediment
768 Research. 24 108-125.[https://doi.org/10.1016/S1001-6279\(09\)60020-2](https://doi.org/10.1016/S1001-6279(09)60020-2)

769 Just, A. C., et al., 2020. Gradient boosting machine learning to improve satellite-derived column
770 water vapor measurement error. Atmos Meas Tech. 13 4669-
771 4681.<https://doi.org/10.5194/amt-13-4669-2020>

772 Knoben, W. J. M., et al., 2020. A Brief Analysis of Conceptual Model Structure Uncertainty
773 Using 36 Models and 559 Catchments. Water Resources Research. 56
774 e2019WR025975.<https://doi.org/10.1029/2019wr025975>

775 Kratzert, F., et al., 2019. Toward improved predictions in ungauged basins: Exploiting the power
776 of machine learning. Water Resources Research. 55 11344-11354
777 <https://doi.org/10.1029/2019WR026065>

778 Kupiec, J. M., et al., 2021. Assessment of the impact of land use in an agricultural catchment
779 area on water quality of lowland rivers. PeerJ. 9
780 e10564.<https://doi.org/10.7717/peerj.10564>

781 Legendre, P., et al., 2015. Should the Mantel test be used in spatial analysis? Methods in
782 Ecology and Evolution. 6 1239-1247.<https://doi.org/10.1111/2041-210x.12425>

783 Li, C., Ji, H., 2016. Chemical weathering and the role of sulfuric and nitric acids in carbonate
784 weathering: Isotopes (^{13}C , ^{15}N , ^{34}S , and ^{18}O) and chemical constraints. Journal of
785 Geophysical Research: Biogeosciences. 121 1288-
786 1305.<https://doi.org/10.1002/2015jg003121>

787 Li, C., et al., 2019. Identification of sources and transformations of nitrate in the Xijiang River
788 using nitrate isotopes and Bayesian model. Sci Total Environ. 646 801-
789 810.<https://doi.org/10.1016/j.scitotenv.2018.07.345>

790 Li, S. L., et al., 2021. Karst ecosystem and environment: Characteristics, evolution processes,

791 and sustainable development. Agriculture Ecosystems & Environment. 306
792 107173.<https://doi.org/ARTN10717310.1016/j.agee.2020.107173>

793 Li, S. L., et al., 2020a. Effects of agricultural activities coupled with karst structures on riverine
794 biogeochemical cycles and environmental quality in the karst region. Agriculture
795 Ecosystems & Environment. 303
796 107120.[https://doi.org/ARTN10712010.1016/j.agee.2020.](https://doi.org/ARTN10712010.1016/j.agee.2020.107120)
797 [107120](https://doi.org/ARTN10712010.1016/j.agee.2020.107120)

798 Li, X., et al., 2018. Watershed System Model: The Essentials to Model Complex Human-Nature
799 System at the River Basin Scale. Journal of Geophysical Research-Atmospheres. 123
800 3019-3034.<https://doi.org/10.1002/2017jd028154>

801 Li, X., et al., 2020b. A harmonized global nighttime light dataset 1992-2018. Sci Data. 7
802 168.<https://doi.org/10.1038/s41597-020-0510-y>

803 [Lintern, A., Webb, J., Ryu, D., Liu, S., Bende-Michl, U., Waters, D., Leahy, P., Wilson, P. and](https://doi.org/10.1038/s41597-020-0510-y)
804 [Western, A.W. \(2018\), Key factors influencing differences in stream water quality](https://doi.org/10.1038/s41597-020-0510-y)
805 [across space. WIREs Water, 5: e1260. https://doi.org/10.1002/wat2.1260](https://doi.org/10.1038/s41597-020-0510-y)

806 Liu, J., et al., 2018. Assessing how spatial variations of land use pattern affect water quality
807 across a typical urbanized watershed in Beijing, China. Landscape and Urban Planning.
808 176 51-63.<https://doi.org/10.1016/j.landurbplan.2018.04.006>

809 Liu, L., et al., 2020. Insights into the long-term pollution trends and sources contributions in
810 Lake Taihu, China using multi-statistic analyses models. Chemosphere. 242
811 125272.<https://doi.org/10.1016/j.chemosphere.2019.125272>

812 Liu, S., et al., 2021. A multi-model approach to assessing the impacts of catchment

813 characteristics on spatial water quality in the Great Barrier Reef catchments. Environ
814 Pollut. 288 117337. <https://doi.org/10.1016/j.envpol.2021.117337>

815 Long, T. Y., et al., 2014. Evaluation of stormwater and snowmelt inputs, land use and
816 seasonality on nutrient dynamics in the watersheds of Hamilton Harbour, Ontario,
817 Canada. Journal of Great Lakes Research. 40 964-
818 979. <https://doi.org/10.1016/j.jglr.2014.09.017>

819 Lu, J. Y., et al., 2019. Optimizing operation of municipal wastewater treatment plants in China:
820 The remaining barriers and future implications. Environ Int. 129 273-278.
821 <https://doi.org/10.1016/j.envint.2019.05.057>

822 Lumb, A., et al., 2011. A Review of Genesis and Evolution of Water Quality Index (WQI) and
823 Some Future Directions. Water Quality Exposure and Health. 3 11-
824 24. <https://doi.org/10.1007/s12403-011-0040-0>

825 Lundberg, S. M., et al., 2020. From Local Explanations to Global Understanding with
826 Explainable AI for Trees. Nat Mach Intell. 2 56-67. <https://doi.org/10.1038/s42256-019-0138-9>

827 [0138-9](https://doi.org/10.1038/s42256-019-0138-9)

828 Lundberg, S. M., et al., 2018. Consistent individualized feature attribution for tree ensembles.
829 arXiv preprint arXiv:1802.03888. <https://doi.org/10.48550/arXiv.1705.07874>

830 Lundberg, S. M., Lee, S.-I., A unified approach to interpreting model predictions. Proceedings
831 of the 31st international conference on neural information processing systems, 2017,
832 pp. 4768-4777. <https://doi.org/10.48550/arXiv.1705.07874>

833 Maavara, T., Chen, Q., Van Meter, K. et al. River dam impacts on biogeochemical cycling. Nat
834 Rev Earth Environ 1, 103–116 (2020). <https://doi.org/10.1038/s43017-019-0019-0>

835 Mainali, J., Chang, H., 2018. Landscape and anthropogenic factors affecting spatial patterns of
836 water quality trends in a large river basin, South Korea. *Journal of Hydrology*. 564 26-
837 40.<https://doi.org/10.1016/j.jhydrol.2018.06.074>

838 Malago, A., et al., 2016. Regional scale hydrologic modeling of a karst-dominant
839 geomorphology: The case study of the Island of Crete. *Journal of Hydrology*. 540 64-
840 81.<https://doi.org/10.1016/j.jhydrol.2016.05.061>

841 Mandaric, L., et al., 2018. Impact of urban chemical pollution on water quality in small, rural and
842 effluent-dominated Mediterranean streams and rivers. *Sci Total Environ*. 613-614 763-
843 772.<https://doi.org/10.1016/j.scitotenv.2017.09.128>

844 Mayorga, E., et al., 2010. Global Nutrient Export from WaterSheds 2 (NEWS 2): Model
845 development and implementation. *Environmental Modelling & Software*. 25 837-
846 853.<https://doi.org/10.1016/j.envsoft.2010.01.007>

847 Mello, K. d., et al., 2018. Effects of land use and land cover on water quality of low-order streams
848 in Southeastern Brazil: Watershed versus riparian zone. *Catena*. 167 130-
849 138.<https://doi.org/10.1016/j.catena.2018.04.027>

850 Meybeck, M., et al., Global freshwater quality: a first assessment. *Global freshwater quality: A*
851 *first assessment, 1990*, pp. 306-306. <https://doi.org/10.1577/1548-8659-121.1.141>

852 Mokoatle, M., et al., Predicting road traffic accident severity using accident report data in South
853 Africa. *Proceedings of the 20th Annual International Conference on Digital Government*
854 *Research, 2019*, pp. 11-17. <https://doi.org/10.1145/3325112.3325211>

855 Molnar, C., 2020. 5.10 SHAP (SHapley Additive exPlanations) Interpretable machine learning.
856 Lulu. com, 2020.

857 Moreira, C., et al., 2020. An Interpretable Probabilistic Approach for Demystifying Black-box
858 Predictive Models.<https://doi.org/10.48550/arXiv.2007.10668>

859 Moriasi, D. N., et al., 2007. Model evaluation guidelines for systematic quantification of
860 accuracy in watershed simulations. Transactions of the Asabe. 50 885-
861 900.[https://doi.org/Doi 10.13031/2013.23153](https://doi.org/Doi%2010.13031/2013.23153)

862 Najah Ahmed, A., et al., 2019. Machine learning methods for better water quality prediction.
863 Journal of Hydrology. 578 124084.<https://doi.org/10.1016/j.jhydrol.2019.124084>

864 Nash, J. E., Sutcliffe, J. V., 1970. River flow forecasting through conceptual models part I—A
865 discussion of principles. Journal of hydrology. 10 282-290.
866 [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6)

867 Nazeer, S., et al., 2014. Heavy metals distribution, risk assessment and water quality
868 characterization by water quality index of the River Soan, Pakistan. Ecological
869 Indicators. 43 262-270.<https://doi.org/10.1016/j.ecolind.2014.03.010>

870 Nie, Y., et al., 2017. Comparison of Rooting Strategies to Explore Rock Fractures for Shallow
871 Soil-Adapted Tree Species with Contrasting Aboveground Growth Rates: A
872 Greenhouse Microcosm Experiment. Front Plant Sci. 8
873 1651.<https://doi.org/10.3389/fpls.2017.01651>

874 Nizzoli, D., et al., 2018. Denitrification in a meromictic lake and its relevance to nitrogen flows
875 within a moderately impacted forested catchment. Biogeochemistry. 137 143-
876 161.<https://doi.org/10.1007/s10533-017-0407-9>

877 Noori, R., et al., 2012. Chemometric Analysis of Surface Water Quality Data: Case Study of the
878 Gorganrud River Basin, Iran. Environmental Modeling & Assessment. 17 411-

879 420.<https://doi.org/10.1007/s10666-011-9302-2>

880 Norse, D., Ju, X., 2015. Environmental costs of China's food security. *Agriculture, Ecosystems*
881 & *Environment*. 209 5-14. <https://doi.org/10.1016/j.agee.2015.02.014>

882 Ockenden, M. C., et al., 2017. Major agricultural changes required to mitigate phosphorus
883 losses under climate change. *Nat Commun*. 8 161. [https://doi.org/10.1038/s41467-](https://doi.org/10.1038/s41467-017-00232-0)
884 [017-00232-0](https://doi.org/10.1038/s41467-017-00232-0)

885 Oliver, D. M., et al., 2020. How does smallholder farming practice and environmental
886 awareness vary across village communities in the karst terrain of southwest China?
887 *Agriculture, Ecosystems & Environment*. 288 106715.
888 <https://doi.org/10.1016/j.agee.2019.106715>

889 Pacheco, F. A. L., et al., 2018. An approach to validate groundwater contamination risk in rural
890 mountainous catchments: the role of lateral groundwater flows. *MethodsX*. 5 1447-
891 1455.<https://doi.org/10.1016/j.mex.2018.11.002>

892 Parsa, A. B., et al., 2020. Toward safer highways, application of XGBoost and SHAP for real-
893 time accident detection and feature analysis. *Accid Anal Prev*. 136
894 105405.<https://doi.org/10.1016/j.aap.2019.105405>

895 Perry, J., Vanderklein, E. L., 2009. *Water quality: management of a natural resource*. John
896 Wiley & Sons. [https://doi.org/10.1016/S0025-326X\(97\)00155-0](https://doi.org/10.1016/S0025-326X(97)00155-0)

897 Powers, S. M., et al., 2016. Long-term accumulation and transport of anthropogenic
898 phosphorus in three river basins. *Nature Geoscience*.
899 9353.<https://doi.org/10.1038/Ngeo2693>

900 Prasad, V. K., et al., 2005. Exploring the relationship between hydrologic parameters and

901 nutrient loads using digital elevation model and GIS—a case study from Sugarcreek
902 headwaters, Ohio, USA. Environmental monitoring and assessment. 110 141-169
903 <https://doi.org/10.1007/s10661-005-6688-9>

904 Qin, N., et al., 2015. Impacts of climate change on regional hydrological regimes of the Wujiang
905 River watershed in the Karst area, Southwest China. Geoenvironmental Disasters. 2
906 [10.https://doi.org/10.1186/s40677-015-0013-x](https://doi.org/10.1186/s40677-015-0013-x)

907 Qu, J. H., et al., 2019. Municipal wastewater treatment in China: Development history and future
908 perspectives. Frontiers of Environmental Science & Engineering. 13
909 [88.https://doi.org/ARTN 8810.1007/s11783-019-1172-x](https://doi.org/ARTN 8810.1007/s11783-019-1172-x)

910 Quinn, P., et al., 1991. The Prediction of Hillslope Flow Paths for Distributed Hydrological
911 Modeling Using Digital Terrain Models. Hydrological Processes. 5 59-79.
912 <https://doi.org/DOI 10.1002/hyp.3360050106>

913 Rodhe, A., Seibert, J., 1999. Wetland occurrence in relation to topography: a test of topographic
914 indices as moisture indicators. Agricultural and Forest Meteorology. 98-9 325-
915 [340.https://doi.org/10.1016/S0168-1923\(99\)00104-5](https://doi.org/10.1016/S0168-1923(99)00104-5)

916 Rodriguez-Blanco, M. L., et al., 2015. Relating nitrogen export patterns from a mixed land use
917 catchment in NW Spain with rainfall and streamflow. Hydrological Processes. 29 2720-
918 [2730.https://doi.org/10.1002/hyp.10388](https://doi.org/10.1002/hyp.10388)

919 Rodriguez-Galiano, V., et al., 2014. Predictive modeling of groundwater nitrate pollution using
920 Random Forest and multisource variables related to intrinsic and specific vulnerability:
921 a case study in an agricultural setting (Southern Spain). Sci Total Environ. 476-477
922 [189-206.https://doi.org/10.1016/j.scitotenv.2014.01.001](https://doi.org/10.1016/j.scitotenv.2014.01.001)

923 Sardans, J., et al., 2008. Changes in soil enzymes related to C and N cycle and in soil C and
924 N content under prolonged warming and drought in a Mediterranean shrubland. Applied
925 Soil Ecology. 39 223-235.<https://doi.org/10.1016/j.apsoil.2007.12.011>

926 Schwarz, G., et al., Section 3. The SPARROW Surface Water-Quality Model—Theory,
927 application and user documentation. Techniques and Methods, Reston, VA, 2006. e B.
928 Hoos , RB Alexander 和 RA Smith <https://doi.org/10.3133/tm6B3>

929 Sharpley, A., et al., 2013. Phosphorus legacy: overcoming the effects of past management
930 practices to mitigate future water quality impairment. J Environ Qual. 42 1308-
931 26.<https://doi.org/10.2134/jeq2013.03.0098>

932 Sheng, M. Y., et al., 2018. Response of soil physical and chemical properties to Rocky
933 desertification succession in South China Karst. Carbonates and Evaporites. 33 15-
934 28.<https://doi.org/10.1007/s13146-016-0295-4>

935 Sierra, C. A., et al., 2017. Monitoring ecological change during rapid socio-economic and
936 political transitions: Colombian ecosystems in the post-conflict era. Environmental
937 Science & Policy. 76 40-49.<https://doi.org/10.1016/j.envsci.2017.06.011>

938 Singh, J., et al., 2005a. Hydrological modeling of the Iroquois river watershed using HSPF and
939 SWAT 1. JAWRA Journal of the American Water Resources Association. 41 343-360
940 <https://doi.org/10.1111/j.1752-1688.2005.tb03740.x>

941 Singh, J., et al., 2005b. Hydrological Modeling of the Iroquois River Watershed Using Hspf and
942 Swat. Journal of the American Water Resources Association. 41 343-
943 360.<https://doi.org/10.1111/j.1752-1688.2005.tb03740.x>

944 Sinha, E., Michalak, A. M., 2016. Precipitation Dominates Interannual Variability of Riverine

945 Nitrogen Loading across the Continental United States. Environ Sci Technol. 50 12874-
946 12884. <https://doi.org/10.1021/acs.est.6b04455>

947 Skinner, J. A., et al., 1997. An overview of the environmental impact of agriculture in the UK.
948 Journal of Environmental Management. 50 111-
949 128. <https://doi.org/10.1006/jema.1996.0103>

950 Smith, R. A., et al., 1997. Regional interpretation of water-quality monitoring data. Water
951 resources research. 33 2781-2798 <https://doi.org/10.1029/97WR02171>

952 Strobl, C., et al., 2009. An introduction to recursive partitioning: rationale, application, and
953 characteristics of classification and regression trees, bagging, and random forests.
954 Psychol Methods. 14 323-48. <https://doi.org/10.1037/a0016973>

955 Strumbelj, E., Kononenko, I., 2014. Explaining prediction models and individual predictions with
956 feature contributions. Knowledge and Information Systems. 41 647-
957 665. <https://doi.org/10.1007/s10115-013-0679-x>

958 Suescun, D., et al., 2017. Vegetation cover and rainfall seasonality impact nutrient loss via
959 runoff and erosion in the Colombian Andes. Regional Environmental Change. 17 827-
960 839. <https://doi.org/10.1007/s10113-016-1071-7>

961 Sun, A. Y., Scanlon, B. R., 2019. How can Big Data and machine learning benefit environment
962 and water management: a survey of methods, applications, and future directions.
963 Environmental Research Letters. 14
964 073001. <https://doi.org/ARTN07300110.1088/1748-9326/ab1b7d>

965 Sutadian, A. D., et al., 2016. Development of river water quality indices—a review.
966 Environmental monitoring and assessment. 188 58 <https://doi.org/10.1007/s10661->

967 015-5050-0

968 Thomas, I. A., et al., 2016. Improving the identification of hydrologically sensitive areas using

969 LiDAR DEMs for the delineation and mitigation of critical source areas of diffuse

970 pollution. *Sci Total Environ.* 556 276-90.<https://doi.org/10.1016/j.scitotenv.2016.02.183>

971 Valle Junior, R. F., et al., 2014. Environmental land use conflicts: A threat to soil conservation.

972 *Land Use Policy.* 41 172-185.<https://doi.org/10.1016/j.landusepol.2014.05.012>

973 van Vliet, M. T. H., et al., 2013. Global river discharge and water temperature under climate

974 change. *Global Environmental Change.* 23 450-464.

975 <https://doi.org/10.1016/j.gloenvcha.2012.11.002>

976 Varanka, S., et al., 2015. Geomorphological factors predict water quality in boreal rivers. *Earth*

977 *Surface Processes and Landforms.* 40 1989-1999.<https://doi.org/10.1002/esp.3601>

978 Viaroli, P., et al., 2018. Space and time variations of watershed N and P budgets and their

979 relationships with reactive N and P loadings in a heavily impacted river basin (Po river,

980 Northern Italy). *Sci Total Environ.* 639 1574-

981 1587.<https://doi.org/10.1016/j.scitotenv.2018.05.233>

982 Vorosmarty, C. J., Sahagian, D., 2000. Anthropogenic disturbance of the terrestrial water cycle.

983 *Bioscience.* 50 753-765.[https://doi.org/10.1641/0006-](https://doi.org/10.1641/0006-3568(2000)050[0753:Adottw2.0.Co;2)

984 [3568\(2000\)050\[0753:Adottw2.0.Co;2](https://doi.org/10.1641/0006-3568(2000)050[0753:Adottw2.0.Co;2)

985 Wang, F., et al., 2021a. Spatial heterogeneity modeling of water quality based on random forest

986 regression and model interpretation. *Environ Res.* 202

987 111660.<https://doi.org/10.1016/j.envres.2021.111660>

988 Wang, F., et al., 2021b. Spatial heterogeneity modeling of water quality based on random forest

989 regression and model interpretation. *Environmental Research.* 202

990 111660.<https://doi.org/https://doi.org/10.1016/j.envres.2021.111660>

991 Wang, Z.-J., et al., 2020. Rainfall driven nitrate transport in agricultural karst surface river
992 system: Insight from high resolution hydrochemistry and nitrate isotopes. Agriculture,
993 Ecosystems & Environment. 291 106787.<https://doi.org/10.1016/j.agee.2019.106787>

994 Winemiller, K. O., et al., 2016. DEVELOPMENT AND ENVIRONMENT. Balancing hydropower
995 and biodiversity in the Amazon, Congo, and Mekong. Science. 351 128-
996 9.<https://doi.org/10.1126/science.aac7082>

997 Wu, Z., et al., 2018. Assessing river water quality using water quality index in Lake Taihu Basin,
998 China. Sci Total Environ. 612 914-922.<https://doi.org/10.1016/j.scitotenv.2017.08.293>

999 Xia, X. H., et al., 2015. Potential Impacts of Climate Change on the Water Quality of Different
1000 Water Bodies. Journal of Environmental Informatics. 25 85-98. [https://doi.org](https://doi.org/10.3808/jei.201400263)
1001 [/10.3808/jei.201400263](https://doi.org/10.3808/jei.201400263)

1002 Xu, G., et al., 2021a. Spatio-temporal characteristics and determinants of anthropogenic
1003 nitrogen and phosphorus inputs in an ecologically fragile karst basin: Environmental
1004 responses and management strategies. Ecological Indicators. 133 108453.
1005 [https://doi.org/ 10.1016/j.eco lind.](https://doi.org/10.1016/j.ecoind.)

1006 Xu, G. Y., et al., 2019. Influence of Landscape Structures on Water Quality at Multiple Temporal
1007 and Spatial Scales: A Case Study of Wujiang River Watershed in Guizhou. Water. 11
1008 159.<https://doi.org/ARTN15910.3390/w11010159>

1009 Xu, G. Y., et al., 2021b. Spatio-temporal characteristics and determinants of anthropogenic
1010 nitrogen and phosphorus inputs in an ecologically fragile karst basin: Environmental
1011 responses and management strategies. Ecological Indicators. 133 108453.<https://doi.org/10.1016/j.ecoind.2021.108453>

1012 [rg/ARTN_10845310.1016/j.ecolind.2021.108453](https://doi.org/10.1016/j.ecolind.2021.108453)

1013 Yan, W., et al., 2021a. The effect of landscape complexity on water quality in mountainous
1014 urbanized watersheds: a case study in Chongqing, China. *Landscape and Ecological*
1015 *Engineering*. 17 165-193.<https://doi.org/10.1007/s11355-021-00448-9>

1016 Yan, W. T., et al., 2021b. The effect of landscape complexity on water quality in mountainous
1017 urbanized watersheds: a case study in Chongqing, China. *Landscape and Ecological*
1018 *Engineering*. 17 165-193.<https://doi.org/10.1007/s11355-021-00448-9>

1019 Yi, Q. T., et al., 2017. Tracking Nitrogen Sources, Transformation, and Transport at a Basin
1020 Scale with Complex Plain River Networks. *Environmental Science & Technology*. 51
1021 5396-5403.<https://doi.org/10.1021/acs.est.6b06278>

1022 Yue, F. J., et al., 2019. Land use interacts with changes in catchment hydrology to generate
1023 chronic nitrate pollution in karst waters and strong seasonality in excess nitrate export.
1024 *Science of the Total Environment*. 696 134062.
1025 <https://doi.org/ARTN13406210.1016/j.scitotenv34062>

1026 Zalidis, G., et al., 2002. Impacts of agricultural practices on soil and water quality in the
1027 Mediterranean region and proposed assessment methodology. *Agriculture*
1028 *Ecosystems & Environment*. 88 137-146.[https://doi.org/10.1016/S0167-](https://doi.org/10.1016/S0167-8809(01)00249-3)
1029 [8809\(01\)00249-3](https://doi.org/10.1016/S0167-8809(01)00249-3)

1030 Zeller, K. A., et al., 2016. Using simulations to evaluate Mantel-based methods for assessing
1031 landscape resistance to gene flow. *Ecol Evol*. 6 4115-28. [https://doi.org/10.](https://doi.org/10.1002/ece3.2102)
1032 [1002/ece3.2](https://doi.org/10.1002/ece3.2102)
1033 154

1034 Zeng, C. N., et al., 2020. Modeling Water Allocation under Extreme Drought of South-to-North
1035 Water Diversion Project in Jiangsu Province, Eastern China. *Frontiers in Earth Science*.
1036 [8.https://doi.org/ARTN_54166410.3389/feart.2020.541664](https://doi.org/ARTN_54166410.3389/feart.2020.541664)

1037 Zhang, L., et al., 2004. A rational function approach for estimating mean annual
1038 evapotranspiration. *Water Resources Research*.
1039 [40.https://doi.org/ArtnW0250210.1029/2003wr002710](https://doi.org/ArtnW0250210.1029/2003wr002710)

1040 Zhang, Y. G., et al., 2018. A High-Resolution Global Map of Soil Hydraulic Properties Produced
1041 by a Hierarchical Parameterization of a Physically Based Water Retention Model.
1042 *Water Resources Research*. 54 9774-9790. <https://doi.org/10.1029/2018wr023539>

1043 Zhang, Z., et al., 2020. Coupled hydrological and biogeochemical modelling of nitrogen
1044 transport in the karst critical zone. *Sci Total Environ*. 732
1045 138902. <https://doi.org/10.1016/j.scitotenv>
1046 .2020.138902

1047 Zou, X.-Y., et al., 2019. A Novel Event Detection Model for Water Distribution Systems Based
1048 on Data-Driven Estimation and Support Vector Machine Classification. *Water*
1049 *Resources Management*. 33 4569-4581 <https://doi.org/10.1007/s11269-019-02317-5>

1050

