



# Ensemble Model Development for the Prediction of a Disaster Index in Water Treatment Systems

# Jungsu Park <sup>1,\*</sup>, Jae-Hyeoung Park <sup>2</sup>, June-Seok Choi <sup>3</sup>, Jin Chul Joo <sup>1</sup>, Kihak Park <sup>4</sup>, Hyeon Cheol Yoon <sup>5</sup>, Cheol Young Park <sup>6</sup>, Woo Hyoung Lee <sup>7</sup> and Tae-Young Heo <sup>8,\*</sup>

- <sup>1</sup> Department of Civil and Environmental Engineering, Hanbat National University, 125, Dongseo-daero, Yuseong-gu, Daejeon 34158, Korea; jincjoo@hanbat.ac.kr
- <sup>2</sup> G&C Environmental Solution, 16-5, Seongmisan-ro 23-gil, Mapo-gu, Seoul 03979, Korea; gncenvsol@naver.com
- <sup>3</sup> Korea Institute of Civil Engineering and Building Technology, 283, Goyang-daero, Ilsanseo-gu, Goyang-si, Gyeonggi-do 10223, Korea; jschoi@kict.re.kr
- <sup>4</sup> Department of Environmental Energy Engineering, University of Suwon, 17, Wau-an-gil, Bongdam-eup, Hwaseong-si, Gyeonggi-do 18323, Korea; parkihak@naver.com
- <sup>5</sup> Disaster Prevention Research Division, National Disaster Management Research Institute, 365, Jongga-ro, Jung-gu, Ulsan 44538, Korea; hcyoon82@korea.kr
- <sup>6</sup> Bayesian AI Laboratory, BAIES, Fairfax, VA 22030, USA; cparkf@gmu.edu
- <sup>7</sup> Department of Civil, Environmental and Construction Engineering, University of Central Florida, 12800 Pegasus Dr., Orlando, FL 32816, USA; woohyoung.lee@ucf.edu
- <sup>8</sup> Department of Information & Statistics, Chungbuk National University, Chungdae-Ro 1, SeoWon-Gu, Cheongju, Chungbuk 28644, Korea
- \* Correspondence: parkjs@hanbat.ac.kr (J.P.); theo@cbnu.ac.kr (T.-Y.H.); Tel.: +82-42-821-1265 (J.P.)

Received: 17 October 2020; Accepted: 13 November 2020; Published: 15 November 2020



**Abstract:** The quantitative analysis of the disaster effect on water supply systems can provide useful information for water supply system management. In this study, a total disaster index (TDI) was developed using open-source public data in 419 water treatment plants in Korea with 23 input variables. The TDI quantifies the possible effects or damage caused by three major disasters (typhoons, heavy rain, and earthquakes) on water supply systems. The four components (regional factor, risk factor, urgency factor, and response and recovery factor) were calculated using input variables to determine the disaster index (DI) of each disaster. The weight of the input variables was determined using principal component analysis (PCA), and the weights of the DI of three natural disasters and four components used to calculate the TDI were determined by the analytical hierarchy process (AHP). Specifically, two ensemble machine learning models, random forest (RF) and XGBoost (XGB), were used to develop models to predict the TDI. Both models predicted the TDI with the coefficient of determination and root-mean-square error-observations standard deviation ratio of 0.8435 and 0.3957 for the RF model and 0.8629 and 0.3703 for the XGB model, respectively. The relative importance analysis suggests that the number of input variables can be minimized, which improves the models' practical applicability.

**Keywords:** disaster management; ensemble model; machine learning; water supply; water treatment system

# 1. Introduction

Various natural disasters, such as floods and earthquakes, cause considerable damage to water supply systems. This damage includes the destruction of plants, intake systems, pipelines, and electric systems, and the consequent interruption of water supply to the public [1]. The assessment of



damage to water supply systems caused by natural disasters is important for proper management and decision-making processes to prevent and restore the damage caused by natural disasters [2,3].

Assessing risk and measuring disaster resilience are the keys for predicting possible events, quantifying contributing factors, and identifying potential consequences. One good example is the lone house that remained standing after Hurricane Ike in 2008. It was rebuilt based on the experience from Hurricane Rita in 2005 on elevated ground, with an appropriate roof pitch and windows that were designed to withstand winds of up to 209 km/h, thus surviving Hurricane Ike with its winds of 177 km/h [4]. Although there have been many efforts to develop quantitative and indicator-based assessments, such as the comprehensive disaster resilience index (CDRI) [5,6], there is no universal standard for the measurement of disaster and related consequences [7]. A reliable disaster resilience framework with unified terminology and its quantitative evaluation would be an important tool in the decision-making processes for both policymakers and engineering professionals [8].

Statistical methods such as principal component analysis (PCA) or analytic hierarchy process (AHP) are often applied for the evaluation of disaster effects on civil infrastructures. For example, Park et al. [9] suggested a disaster risk index for 51 high-speed railroad stations in Korea. The index was calculated from a linear equation of four main indices (hazard, exposure, vulnerability, and emergency response and recovery capability) suggested by Rossi and Gilmartin [10] where the weights of each main index were determined by PCA. Recent studies have also used statistical analysis based on survey data for the assessment of disaster risk on flooding or water security [11,12].

In recent decades, advanced technologies of data mining and machine learning (ML) have been used to manage disasters such as typhoons and earthquakes [13–19], and various emerging remote sensing technologies have also been increasingly used for monitoring and detecting data related to disaster managements [20]. The continuous increase in available data due to advanced data collection technologies, such as remote sensors or unmanned aerial vehicles (UAV), has accelerated the application and accuracy of ML models [21–23]. Ofli et al. [21] used aerial images captured from UAVs for identifying features of interest such as damaged shelters and blocked loads to assist with disaster response. The features of interest in the image were annotated and used for training ML models, including support vector machine (SVM) and random forest (RF), where the overall accuracy of the model's classification results ranges from 0.73 to 0.85. Sheykhmousa et al. [24] analyzed satellite images of land cover and land use using an SVM classifier. The image data during Typhoon Haiyan, which caused massive damage in the Philippines in 2013, was compared with the image data in 2017, four years after Typhoon Haiyan, to assess the post-disaster recovery process. Chen et al. [15] analyzed the impact indices of flood disasters using RF and developed a risk assessment model based on the neural network method. Various data, including rainfall and socioeconomic data in the Yangtze River Delta area between 2008 and 2018, were used for model development. More recently, Kao et al. [19] used an advanced deep learning algorithm, long short-term-memory (LSTM), to forecast flood events.

Recent studies also used social platforms with text information about disasters to analyze the characteristics of disasters such as typhoons and earthquakes [14,25,26]. Resch et al. [25] analyzed earthquake characteristics from social media information during an earthquake in Napa, California, USA, in 2014 using latent Dirichlet allocation (LDA), which is widely used for topic analysis. The spatial hot spots of the earthquake were determined from the LDA model with 86.45% accuracy compared with the United States Geological Survey earthquake footprint report. More recently, Yu et al. [14] analyzed text information in social media during Typhoon Anemone along the coast of China in August 2012 to develop a typhoon disaster classification system using a model based on a convolutional neural network.

ML is also increasingly used in environmental management. Zhang et al. [27] predicted air pollution by PM 2.0 with a fusion model based on three gradient boosted decision tree (GBDT) algorithms. The root-mean-squared error (RMSE) of the fusion model was 32.300. Bi et al. [28] also used a GBDT-based model, light gradient boost, coupled with a fast Fourier transform for the assessment of a liquefaction disaster. However, even with substantial efforts on the classification and analysis of

disaster, its quantitative and indicator-based assessments on the water infrastructure have not been thoroughly conducted.

In this study, the effects of various disasters on water supply systems from the perspective of management are quantified by statistical data analysis methods, PCA, and analytic hierarchy process (AHP). From the statistical approach, a total disaster index (TDI) was developed. In the second part, tree-based ensemble models (i.e., RF and GBDT) were used to predict TDI, which provides valuable information for the safety management of water supply systems.

#### 2. Methods

#### 2.1. Data Sources

Total 23 input variables of facility specification and operational data in 419 water treatment plants in Korea were used to develop a TDI. The data were obtained from statistical yearbooks and open-source public data (Table 1). The 23 input variables provide information about the water supply systems, including water supply capacity, pipeline density, number of customers, management labor, and regional characteristics of natural conditions where the water treatment plants are located (Table 2). The local peak ground acceleration by an earthquake at each water treatment plant was estimated from the Korea Seismicity Map program developed by Cao et al. [29]. The data for regional natural conditions were obtained from meteorological data available from the national meteorological administration information portal [30]. The financial status of a local government that manages the water treatment plant was collected from the public data portal of the Ministry of the Interior and Safety in Korea [31].

Table 1. Data sources	3.
-----------------------	----

Data	Reference
Water treatment plant operational information and facility specification	Statistical yearbook of water treatment system [32]
Meteorological data	Meteorological administration information portal [30]
Financial status in local government	Ministry of the interior and safety information portal [31]
Design standard for wind speed	Korean Design Standard [33]
Local peak ground acceleration by earthquake	Korea Seismicity Map [29]

Table 2. Input variables.

Variables	Description		
CUSTOMER	Population that receives drinking water from the water treatment plant		
EMPLOYEES_AREA	Employees per management area of the authority * (person/km <sup>2</sup> )		
EMPLOYEES_SITE	Employees per number of water treatment plants of the authority * (person/ea)		
EQ1	Seismic design application status (applied: 1, not applied: 0)		
LINE_DENSE	Total pipeline length per management area of the authority * (m/km <sup>2</sup> )		
LOCAL_EMPLOYEES	Number of employees in the water supply plant		
MONEY	Financial independence of the local government (%)		
PGA_500	500 years frequency peak ground acceleration (%)		
PGA_1000	1000 years frequency peak ground acceleration (%)		
PGA_2400	2400 years frequency peak ground acceleration (%)		
PUMP	Number of water supply pumps in the water treatment plant (ea)		
PUMP_EP	Sum of electrical capacity of all water supply pumps in the water supply plant (kW)		

### Table 2. Cont.

Variables	Description		
Q	Water supply capacity of the water treatment plant (m <sup>3</sup> /day)		
Q_DAILY	Daily average water production capacity of the water treatment plant (m <sup>3</sup> /day)		
Q_MAX	Daily maximum water production capacity of the water treatment plant (m <sup>3</sup> /day)		
Q_PRO	Total annual water production capacity of the water treatment plant (m <sup>3</sup> )		
QT	Maximum water supply capacity per hour of the water treatment plant (m <sup>3</sup> /h)		
QW	Total annual electric power usage of the water treatment plant (kWh)		
QY	Total annual amount of water treated by the water treatment plant (m <sup>3</sup> /year)		
RAIN	Number of flood warning advisories between 2015 and 2019 in the region where the water treatment plant is located (times/km <sup>2</sup> )		
SWIND	Number of strong wind advisories between 2015 and 2019 in the region where the water treatment plant is located (times/km <sup>2</sup> )		
TYPHOON	Number of typhoon warning advisories between 2015 and 2019 in the region where the water treatment plant is located (times/km <sup>2</sup> )		
WIND_R	Regional standard wind speed (m/s) in the region where the water treatment plant is located		

\* authority: the owner of the water supply plant (i.e., the local government) and one authority may manage multiple plants.

#### 2.2. Disaster Index

#### 2.2.1. Type of Disaster

Typhoons and heavy rains are among the most frequent disasters in Korea, while Korea has been known to be relatively safe from earthquakes. However, interest in earthquakes in Korea has increased since the two earthquakes with magnitudes of 5.8 and 5.4 on the Richter scale in 2016 and 2017, respectively. In this study, three natural disasters, typhoons, heavy rains, and earthquakes, were selected as the most influential disasters on the water supply system and used for the TDI development considering natural characteristics in Korea.

#### 2.2.2. Component of Disaster Index

The four components (i.e., regional factor, risk factor, urgency factor, and response and recovery factor), describing the level of the damage caused by each type of disaster, were used to determine the disaster index (DI) of three natural disasters as follows.

- 1. Regional factor (RE) represents regional characteristics such as the frequency of natural disaster occurrence in the selected areas;
- 2. Risk factor (RI) represents the quantity of possible damage caused by natural disasters. For example, the RI increases as the capacity of water treatment plants or the length of water supply pipelines increases;
- 3. Urgency factor (UR) represents the urgency of recovery after a disaster. For example, the UR increases with a larger population in the area receiving drinking water; and,
- 4. Response and recovery factor (RR) represents the recovery ability during and after a disaster, which is estimated by the financial status or manpower of the authority of a water treatment plant, such as the local government.

A total of 23 input variables obtained from open-source public data were used to determine the four components (RE, RI, UR, and RR), as summarized in Table 2.

The weights of each variable for the DI of three natural disasters (typhoon, heavy rain, and earthquake) were determined using PCA. PCA is a statistical method that reduces the dimension of variables and determines each variable's relative importance using an eigenvector. The input variables were standardized as an average of zero and standard deviation of one for PCA analysis [34–36].

#### 2.2.4. AHP Analysis

The DI of three natural disasters and four components were used for the calculation of TDI in 419 water treatment plants. However, there was limited data available for the statistical determination of the relative weight of each natural disaster and four components for the TDI calculation. In addition, although the effect of earthquakes on the water supply system is expected to be extremely large, there were only two significant earthquakes in Korea that occurred in 2016 and 2017. Thus, it should be noted that quantitative data for the analysis of the effect of earthquakes was limited.

The weights of three natural disasters and four components used to calculate the TDI were determined by the AHP suggested by Saaty [37,38]. The AHP is a structured data analysis method for complex decision-making, which is also widely used to analyze disaster data [9,23,39]. In AHP, a pairwise comparison matrix of each element for the decision-making process is structured. This structure relates to the matrix's eigenvector, which represents the weight of each element in the decision-making process [37,40].

The survey results from 62 experts or engineers currently working in water treatment plants were used for AHP analysis. The survey data with a consistency ratio (CR) of less than 0.2 was used to calculate the weight of each input variable to maintain the consistency of the AHP analysis result [40–42].

$$CI = \frac{\lambda_{max} - n_f}{n_f - 1} \text{ and } CR = \frac{CI}{RI}$$
(1)

where

 $\lambda_{max}$ : principal eigenvalue in the pairwise comparison matrix,  $n_f$ : number of features, CI: consistency index, RI: random consistency index (RI = 0.90 for n = 4 and RI = 0.58 for n = 3), and CR: consistency ratio.

#### 2.2.5. Disaster Index Model

The TDI is determined by the weighted sum of the DI for three natural disasters using the following equations (Equations (2)–(5)).

$$TDI = a(TI) + b(HI) + c(EI)$$
<sup>(2)</sup>

$$TI = a_t(RE_t) + b_t(RI_t) + c_t(UR_t) - d_t(RR_t)$$
(3)

$$HI = a_h(RE_h) + b_h(RI_h) + c_h(UR_h) - d_h(RR_h)$$
(4)

$$EI = a_e(RE_e) + b_e(RI_e) + c_e(UR_e) - d_e(RR_e)$$
(5)

where

TI: DI for typhoon; HI: DI for heavy rain; EI: DI for earthquake; a, b, and c: weight of each natural DI; and a<sub>t</sub>, b<sub>t</sub>, c<sub>t</sub>, d<sub>t</sub>, a<sub>h</sub>, b<sub>h</sub>, c<sub>h</sub>, d<sub>h</sub>, a<sub>e</sub>, b<sub>e</sub>, c<sub>e</sub>, and d<sub>e</sub>: weight of each component. Subscripts (i.e., t, h and e) from Equations (3)–(5) represents typhoon, heavy rain, and earthquake respectively.

#### 2.3. Disaster Prediction Model

#### 2.3.1. Model Selection

Two ensemble models, RF and GBDT, have been increasingly used as ML models to manage the water environment. Both models show good performance, even for nonlinear relationship analysis, and data with outliers are also applicable for both classification and regression [43,44].

RF is a tree-based ensemble model in which a random data selection approach generates multiple decision trees. RF randomly selects several sets of input features from the original input features by a bagging method before generating the decision trees, which increases the independence and variability of each decision tree. The final RF prediction is determined by averaging the predictive results from individual decision trees in RF [45]. Consequently, the prediction performance of RF can be dramatically improved [46–48] and outperforms other ML models [49]. RF has shown high performance in various domains and has also been continuously applied to environmental research, such as water quality prediction [50,51].

GBDT is an ensemble model based on a gradient boosting method (GBM), called a sequential tree-based calculation process [45,52,53], and a set of decision trees. Unlike RF which determines the final prediction by voting (for classification) or averaging (for regression), GBDT uses the decision tree, called a weak learning model, from a previous stage in the ML process to improve model performance in the following stage. Residual errors of the prior stage are included in developing the decision tree in the current stage to reduce the residual errors by optimizing a specified loss function [45,52]. This optimization process is sequentially performed until the predefined number of decision trees is reached, which is a major difference with RF, where the calculation of each tree is independent.

GBDT is optimized by minimizing an objective function, J, for a training data set with n samples. The regularization term can be added to avoid overfitting of the model [44,54]. Equation (6) shows an illustrative example of the objective function of GBDT [44,54].

$$\mathbf{J} = \sum_{i=n}^{n} L(y_i, \hat{y}_i) + \sum_{k=1}^{K} \Omega(f_k)$$
(6)

where

*f<sub>k</sub>*: function of the *k*th decision tree,

*L*: loss function that calculates the difference between an observation ( $y_i$ ) and model prediction ( $\hat{y}_i$ ) in each decision tree,

 $\Omega$ : regularization function that penalizes the complexity of the model, and n: number of data samples.

The schematics of RF and GBDT are compared in Figure 1, where *X* denotes input features as  $X = x_1, x_2, ..., x_n, h(X, \theta_k), (k = 1, 2, ..., K)$  is a collection of decision trees, and the  $\theta_k$  are independent and identically distributed random vectors [44,45,54].



Figure 1. Schematics of the random forest (RF) and gradient boosted decision tree (GBDT) models.

In this study, both RF and GBDT models were used for the TDI estimation of 419 drinking water treatment plants. The Python open-source libraries of Scikit-learn (for RF) and XGBoost (for GBDT) were used for regression model development [55,56]. XGBoost (XGB) is one of the most popular GBDT implementations developed by Chen and Guestrin [45,54]. Scikit-learn is also a popular Python-based ML library developed by Pedregosa et al. [55].

#### 2.3.2. Model Optimization

The hyperparameters of RF and XGB were optimized by a trial and error method with ten-fold cross-validation using the grid search library in Scikit-learn [57]. The models were developed with 23 input variables of 419 water treatment plants, where the ratio of data used for training and testing of the models was 8:2.

#### 2.3.3. Feature Importance (FI) of Input Variables

The relative importance of input variables on RF and XGB model performance was calculated using the feature importance (FI) algorithm in Scikit-learn [57]. The FI in the tree-based model was computed as the total impurity reduction of the model brought by that feature [55,58,59].

#### 2.3.4. Model Evaluation

The model performance was evaluated by three evaluation indexes (Equations (7)–(9)), RMSE, coefficient of determination ( $\mathbb{R}^2$ ), and RMSE-observation standard deviation ratio (RSR). RSR ranges from 0 to 1 and approaches 0 when the model shows a good fit with observation. The model is considered to predict the observation when RSR < 0.70 [60,61].

$$R^{2} = 1 - \frac{\sum_{i=1}^{n} (M_{i,obs} - M_{i,model})^{2}}{\sum_{i=1}^{n} (M_{i,obs} - \overline{M_{i,obs}})^{2}}$$
(7)

$$RMSE = \sqrt{\frac{\sum_{i=1}^{n} \left(M_{i,obs} - M_{i,model}\right)^{2}}{n}}$$
(8)

$$RSR = \frac{\sqrt{\sum_{i=1}^{n} \left(M_{i,obs} - M_{i,model}\right)^2}}{\sqrt{\sum_{i=1}^{n} \left(M_{i,obs} - \overline{M_{i,obs}}\right)^2}}$$
(9)

where

 $M_{i,obs}$ : observed values,  $\overline{M}_{i,obs}$ : mean of observed values, and  $M_{i,model}$ : model predicted value.

#### 3. Results and Discussion

#### 3.1. Characteristics of Input Variables

Total 23 input variables for the development of DI were identified from open-source public statistical data. The characteristics of the input variables are summarized in Table 3. The frequency of warning advisories of natural disasters was calculated at each water treatment plant from the sum of the three variables in Table 3 (i.e., RAIN, SWIND, and TYPHOON). The frequency of warning advisories ranged from 0.017 to 1.29 times/km<sup>2</sup> and tended to be higher in areas near the ocean as shown in Figure 2 using ArcGIS pro.



**Figure 2.** A spatial distribution of water treatment plants and frequency of natural disasters determined from the disaster warning advisories in Korea.

Variables	Average	Max	Min	Standard Deviation
CUSTOMER	89,902.749	3,030,917.000	0.000	276,768.491
EMPLOYEES_AREA	0.195	3.017	0.006	0.449
EMPLOYEES_SITE	25.819	304.333	0.833	50.505
EQ1	0.222	1.000	0.000	0.416
LINE_DENSE	2718.442	25,817.066	274.779	3812.635
LOCAL_EMPLOYEES	8.043	134.000	0.000	16.030
MONEY	18.972	78.450	4.020	15.198
PGA_500	8.532	11.000	3.199	1.792
PGA_1000	11.515	14.000	5.431	1.884
PGA_2400	16.695	19.000	9.222	1.959
PUMP	2.852	176.000	0.000	9.093
PUMP_EP	444.527	14,400.000	0.000	1637.036
Q	48,083.905	1,600,000.000	30.000	141,017.670
Q_DAILY	30,932.308	1,081,369.000	6.000	90,686.216
Q_MAX	37,629.160	1,221,400.000	20.000	106,454.278
Q_PRO	11,278,843.986	394,699,507.000	2288.000	33,103,181.424
QT	1818.253	69,120.000	0.000	6465.596
QW	2,406,152.687	78,414,779.000	0.000	7,988,910.560
QY	11,623,330.642	402,072,337.000	3276.000	33,936,442.393
RAIN	0.084	0.900	0.009	0.123
SWIND	0.077	0.813	0.002	0.108
TYPHOON	0.016	0.060	0.001	0.013
WIND_R	29.084	44.000	24.000	5.143

Table 3. Characteristics of input variables.

# 3.2.1. PCA Analysis

The weights for each natural disaster index were determined from PCA with 23 input variables (Table 4). The eigenvectors were calculated from PCA and normalized to make the sum of weights of each component to be 1.

Disaster (Index)	Component	Input Variable (Symbol)	Weight
	RE <sub>t</sub>	WIND_R (RE <sub>t1</sub> ) TYPHOON (RE <sub>t2</sub> ) SWIND (RE <sub>t3</sub> )	0.309 0.345 0.346
		sum	1.000
		$Q(RI_{t1})$ $QY(RI_{t2})$	0.143 0.143
		$Q_{PRO}(RI_{t3})$	0.144
	RIt	PUMP ( $RI_{t4}$ )	0.052
	· ·	$PUMP_EP(RI_{t5})$	0.140
Typhoon (TI)		$QI(RI_{6})$	0.132
		$QW(\mathbf{M}_{t7})$	0.130
		Sum	1 000
-			1.000
		$Q_{\text{DAILY}}(\text{UR}_{1})$	0.334
	URt	$Q_{MAX}(UR_{t2})$	0.334
		COSTOMER (UKt3)	0.552
-		sum	1.000
		LOCAL_EMPLOYEES (RR <sub>t1</sub> )	0.248
		MONEY (RR <sub>t2</sub> )	0.235
	RRt	EMPLOYEES_SITE (RR <sub>t3</sub> )	0.263
		EMPLOYEES_AREA (RR <sub>t4</sub> )	0.254
		sum	1.000
		WIND_R (RE <sub>h1</sub> )	0.500
	RE <sub>h</sub>	RAIN (RE <sub>h2</sub> )	0.500
_		sum	1.000
		$Q(RI_{h1})$	0.143
		QY (RI <sub>h2</sub> )	0.143
		Q_PRO (RI <sub>h3</sub> )	0.144
		PUMP (RI <sub>h4</sub> )	0.052
	RI <sub>h</sub>	PUMP_EP (RI <sub>h5</sub> )	0.140
		QT (RI <sub>h6</sub> )	0.132
Heavy rain (HI)		QW (RI <sub>h7</sub> )	0.136
		LINE_DENSE (RI <sub>h8</sub> )	0.110
		sum	1.000
		O DAILY (UR <sub>b1</sub> )	0.334
	I ID	O MAX (UR <sub>b2</sub> )	0.334
	UR <sub>h</sub>	$COSTUMER (UR_{h3})$	0.332
		sum	1.000
-		LOCAL_EMPLOYEES (RRh1)	0.248
		MONEY (RR <sub>h2</sub> )	0.235
	RR <sub>h</sub>	EMPLOYEES_SITE (RRh3)	0.263
		EMPLOYEES_AREA (RRh4)	0.254
		sum	1.000

Table 4. PCA analysis for weight of each component.

Disaster (Index)	Component	Input Variable (Symbol)	Weight
		PGA_500 (RE <sub>e1</sub> )	0.333
	DE	PGA_1000 (RE <sub>e2</sub> )	0.336
	КLе	PGA_2400 (RE <sub>e3</sub> )	0.331
		sum	1.000
-		Q (RI <sub>e1</sub> )	0.143
		QY (RI <sub>e2</sub> )	0.143
		Q_PRO (RI <sub>e3</sub> )	0.144
		PUMP (RI <sub>e4</sub> )	0.052
Earthquake (EI)	RIe	PUMP_EP (RI <sub>e5</sub> )	0.140
		QT (RI <sub>e6</sub> )	0.132
		QW (RI <sub>e7</sub> )	0.136
		LINE_DENSE (RI <sub>e8</sub> )	0.110
		sum	1.000
	UR <sub>e</sub>	Q_DAILY (UR <sub>e1</sub> )	0.334
		Q_MAX (UR <sub>e2</sub> )	0.334
		COSTUMER (UR <sub>e3</sub> )	0.332
		sum	1.000
-		LOCAL_EMPLOYEES (RR <sub>e1</sub> )	0.234
		EQ1 (RR <sub>e2</sub> )	0.056
	RR	MONEY (RR <sub>e3</sub> )	0.222
	KΛ <sub>e</sub>	EMPLOYEES_SITE (RR <sub>e4</sub> )	0.249
		EMPLOYEES_AREA (RR <sub>e5</sub> )	0.239
		sum	1.000

Table 4. Cont.

# 3.2.2. AHP Analysis

The weights for each disaster type were determined from the AHP analysis using the survey data (CR < 0.2) (Table 5). The response rate of the survey was in the range between 52 and 69% for each item. The weights of each disaster are in the order of typhoons, earthquakes, and heavy rain.

(a) Weights for Disaster Type.					
Disaster Weight					
0.481					
0.198					
0.321					
Sum 1.000					
CR 0.054					
(b) Weights for Each Component.					
Disaster Component Weight					
REt	0.275				
RIt	0.265				
URt	0.216				
RRt	0.244				
Sum	1.000				
CR	0.017				
	Weight           0.481           0.198           0.321           1.000           0.054           for Each Compo           Component           RE <sub>t</sub> RIt           URt           RRt           Sum           CR				

Table 5. Analytical hierarchy process (AHP) analysis results.

(b) Weights for Each Component.			
Disaster	Disaster Component		
	RE <sub>h</sub>	0.279	
	RI <sub>h</sub>	0.247	
Howarain	UR <sub>h</sub>	0.221	
Tleavy failt	RR <sub>h</sub>	0.253	
	Sum	1.000	
	CR	0.004	
Earthquake	RE <sub>e</sub>	0.215	
	RI <sub>e</sub>	0.370	
	URe	0.235	
	RRe	0.180	
	Sum	1.000	
	CR	0.040	

Table 5. Cont.

#### 3.2.3. Disaster Index (DI)

The TDI was determined using the following model (Equations (10)–(13)) which were developed from PCA and AHP analysis (Tables 4 and 5).

$$TDI = 0.481(TI) + 0.198(HI) + 0.321(EI)$$
(10)

$$TI = 0.275(RE_t) + 0.265(RI_t) + 0.216(UR_t) - 0.244(RR_t)$$
(11)

where

 $RE_{t} = 0.309(RE_{t1}) + 0.345(RE_{t2}) + 0.346(RE_{t3}),$ 

$$\label{eq:RIt} \begin{split} RI_t &= 0.143(RI_{t1}) + 0.143(RI_{t2}) + 0.144(RI_{t3}) + 0.052(RI_{t4}) + 0.140(RI_{t5}) + 0.132(RI_{t6}) + 0.136(RI_{t7}) + 0.110(RI_{t8}), \end{split}$$

$$\begin{split} &UR_t = 0.334(UR_{t1}) + 0.334(UR_{t2}) + 0.332(UR_{t3}) \text{, and} \\ &RR_t = 0.248(RR_{t1}) + 0.235(RR_{t2}) + 0.263(RR_{t3}) + 0.254(RR_{t4}). \end{split}$$

$$HI = 0.279(RE_h) + 0.247(RI_h) + 0.221(UR_h) - 0.253(RR_h)$$
(12)

where

 $RE_h = 0.500(RE_{h1}) + 0.500(RE_{h2}),$ 

$$\begin{split} \text{RI}_{h} &= 0.143(\text{RI}_{h1}) + 0.143(\text{RI}_{h2}) + 0.144(\text{RI}_{h3}) + 0.052(\text{RI}_{h4}) + 0.140(\text{RI}_{h5}) + 0.132(\text{RI}_{h6}) + 0.136(\text{RI}_{h7}) \\ &+ 0.110~(\text{RI}_{h8}), \end{split}$$

 $UR_h = 0.334(UR_{h1}) + 0.334(UR_{h2}) + 0.332(UR_{h3})$ , and

 $RR_{h} = 0.248(RR_{h1}) + 0.235(RR_{h2}) + 0.263(RR_{h3}) + 0.254(RR_{h4}).$ 

$$EI = 0.215(RE_e) + 0.370(RI_e) + 0.235(UR_e) - 0.180(RR_e)$$
(13)

where

 $RE_e = 0.333(RE_{e1}) + 0.336(RE_{e2}) + 0.331(RE_{e3}),$ 

$$\begin{split} RI_{e} &= 0.143(RI_{e1}) + 0.143(RI_{e2}) + 0.144(RI_{e3}) + 0.052(RI_{e4}) + 0.140(RI_{e5}) + 0.132(RI_{e6}) + 0.136(RI_{e7}) \\ &+ 0.110(RI_{e8}), \end{split}$$

 $UR_e = 0.334(UR_{e1}) + 0.334(UR_{e2}) + 0.332(UR_{e3})$ , and

 $RR_{e} = 0.234(RR_{e1}) + 0.056(RR_{e2}) + 0.222(RR_{e3}) + 0.249(RR_{e4}) + 0.239(RR_{e5}).$ 

Using the developed models, TDI values of 419 water treatment plants were determined with the range between -0.526 and 3.813 with an average of 0 and a standard deviation of 0.343. A higher TDI represents a higher potential of effect or damage by a disaster in water treatment systems. The TDI tends to be higher in water treatment plants near metropolitan cities as well as the areas near ocean.

The TDI was developed considering the natural status of Korea. For example, there were only two earthquakes in 2016 and 2017, which were considered to have caused actual damage to water treatment plants in Korea. As the data available for the quantification of damage by earthquakes is minimal, the AHP based on survey data was used for the DI calculation.

Although there were not many cases of damage in water treatment systems from earthquakes, the weight of the earthquake was larger than that of heavy rain. The AHP results represent that, although earthquakes have been rare in Korea, the damage and consequences by an earthquake would not be negligible when it occurs, indicating that a preventive plan against earthquakes should be prepared in advance. In addition, given that most of the facilities already experience heavy rain and are relatively well prepared for these instances, it is expected that the actual damage caused by heavy rain is relatively small compared to other disasters.

#### 3.3. Ensemble Model Simulation

#### 3.3.1. Total Disaster Index (TDI) Prediction using Ensemble Models

Two ensemble ML models, RF and XGB, were used to develop a model to predict TDI. The model performance with the test data set was evaluated by three indices, as summarized in Table 6. The R<sup>2</sup> and RSR were 0.8435 and 0.3957 for the RF model and 0.8629 and 0.3703 for the XGB model, respectively.

Table 6. Summary of model evaluation result	lts.
---	------

Model	RMSE	<b>R</b> <sup>2</sup>	RSR
RF	0.100	0.8435	0.3957
XGB	0.093	0.8629	0.3703

The observed data and model predictions are compared in Figure 3. The model prediction shows a similar good fit with observations both in the RF and XGB models, while XGB showed a slightly better performance for all three evaluation indexes (Table 6 and Figure 3).



Figure 3. Comparison of model prediction.

The FI of 23 input variables for both RF and XGB models to predict DI are shown in Figure 4. The FI was different between RF and XGB, while the variables that represent the scale of water treatment plants such as PUMP\_EP and Q tend to have a higher effect on model performance for both models. For RF, the sum of FI in the highest nine input variables was more than 80%, while for XGB, the sum of FI in the highest four variables was more than 80% of the total FI for XGB.



Figure 4. Feature importance (FI) of (a) RF and (b) XGBoost (XGB).

The performance of the models was compared between RF and XGB using fewer input variables, starting with 1 and adding up to 10 input variables with the order from the highest FI (Figure 5). The RF model showed a tendency to improve the performance of the model as the number of input variables increased from one to ten, and even when using three input variables, the RSR was 0.6954, indicating that the model accurately predicted the observation. XGB shows better performance when using fewer input variables. The RSR is 0.5323 when only three input variables were applied, which reduces to 0.3937 when using ten input variables. The FI analysis shows that several input variables with higher feature importance have a considerable effect on model performance. The analysis results show that both the RF and XGB models show similar performance when using five or more input variables with higher FI. The FI is one of the factors and not an absolute standard considered for model structure. The necessary input variables are not always obtainable from the actual operation and management of water treatment systems. Thus, the practical applicability of the model would be improved as fewer input variables are used. The FI analysis suggests that the model shows acceptable performance if only part of the input variables with the highest FI would increase the practical applicability of the model.



Figure 5. Model sensitivity to the number of input variables included in a model (RF or XGB).

#### 4. Summary and Conclusions

In this study, a disaster index (DI) for predicting the effect or damage caused by three major natural disasters in Korea (i.e., typhoons, heavy rain, and earthquakes) was newly developed to quantify each natural disaster's effect on water utilities.

Although the operational data in water utilities provided a good understanding regarding the effect of disasters, the data is usually collected in an individually specified format often site-specific, making it difficult to collect, organize, and analyze the data. In addition, the operational data for water utilities was not easily accessible, limiting the comprehensive development of the DI. Therefore, in this study, the DI of natural disasters in water treatment systems was developed using statistical open-source public data. Two well-defined statistical data analysis methods (i.e., AHP and PCA) were used for the determination of DI.

The open-source public data have greater accessibility and are updated regularly, so the DI can also be updated considering the current status, which is also a significant benefit of using open-source

public data. The DI developed in this study may be site-specific at a given location and conditions of water utilities, but the developed framework would be applicable for quantifying the effect of disasters on water treatment systems in other regions with different natural status.

In the second part, two ensemble models (i.e., RF and XGB) were used to develop models to predict TDI. Both RF and XGB showed similar satisfactory performance for prediction of the DI, while the XGB showed a slightly better performance in general. The FI analysis also suggested that the models have sufficient performance for practical use with only several input variables of the highest FI, which can improve the practical applicability of the models.

Quantitative assessment of disaster effects on water treatment systems is essential for better management of the water treatment systems and stable supply of drinking water to the public. However, data related to disaster analysis are often limited and even hardly quantifiable. One of the possible solutions would be to keep collecting data, analyze them statistically, while facilitating frequent discussions from experts experiencing the disasters in their utilities [11,35]. The recent advance of information and communication technologies, such as sensor-based real-time monitoring methods, can provide various continuous monitoring data about the operational condition of water treatment plants and related infrastructure, which can improve the pre- and post-management planning processes [20,22]. However, the quantification and assessment of disasters on water treatment systems are still in an early stage, and the use of field operational data and responses, in particular during disaster events, is currently limited at this time.

This study provided quantified information on the impact of various natural disasters on water treatment systems with open-source public data, which would be useful for creating a plan to reduce damage to water supply systems caused by natural disasters. Further study is warranted to use high-frequency real-time data to improve the model performance and practical applicability.

Author Contributions: Data curation and software, J.P.; conceptualization: J.P., J.-H.P., J.-S.C., J.C.J., K.P., W.H.L. and T.-Y.H.; investigation, J.P., J.-H.P., J.-S.C., J.C.J., K.P., H.C.Y., C.Y.P., W.H.L. and T.-Y.H.; writing-original draft, J.P.; writing—review and editing, J.-H.P., J.-S.C., J.C.J., K.P., H.C.Y., C.Y.P. and T.-Y.H.; project administration, J.P. and J.-H.P.; supervision, J.P. and T.-Y.H.; funding acquisition, J.P. and J.-H.P. All authors have read and agreed to the published version of the manuscript.

**Funding:** This work was supported by Korea Environment Industry and Technology Institute (KEITI) through Environmental R&D Project on the Disaster Prevention of Environmental Facilities Project, funded by Korea Ministry of Environment (MOE) (2019002870001).

Conflicts of Interest: The authors declare no conflict of interest.

# References

- 1. Pan American Health Organization (PAHO). *Emergencies and Disasters in Drinking Water Supply and Sewage Systems: Guidelines for Effective Response;* PAHO: Washington, DC, USA, 2002; pp. 5–12.
- 2. Davis, C.A. Water system service categories, post-earthquake interaction, and restoration strategies. *Earthq. Spectra* **2014**, *30*, 1487–1509. [CrossRef]
- 3. Matthews, J.C. Disaster resilience of critical water infrastructure systems. *J. Struct. Eng.* **2016**, 142, C6015001. [CrossRef]
- 4. World Meteorological Organization (WMO). *Atlas of Mortality and Economic Losses from Weather, Climate and Water Extremes (1970–2012);* WMO-No. 1123; WMO: Geneva, Switzerland, 2014.
- 5. Marzi, S.; Mysiak, J.; Essenfelder, A.H.; Amadio, M.; Giove, S.; Fekete, A. Constructing a comprehensive disaster resilience index: The case of Italy. *PLoS ONE* **2019**, *14*, e0221585. [CrossRef] [PubMed]
- Beccari, B. A comparative analysis of disaster risk, vulnerability and resilience composite indicators. *PLoS Curr.* 2016, 8. [CrossRef] [PubMed]
- Franc, J.M.; Ingrassia, P.L.; Verde, M.; Colombo, D.; Della Corte, F. A simple graphical method for quantification of disaster management surge capacity using computer simulation and process-control tools. *Prehosp. Disast. Med.* 2015, 30, 9. [CrossRef] [PubMed]
- 8. Cimellaro, G.P.; Reinhorn, A.M.; Bruneau, M. Framework for analytical quantification of disaster resilience. *Eng. Struct.* **2010**, *32*, 3639–3649. [CrossRef]

- 9. Park, Y.; Han, S.; Choi, S. Development of Disaster Risk Index for Evaluating the Natural Disaster Hazards of High-speed Railroad Facilities. *J. Korean Soc. Hazard Mitig.* **2019**, *19*, 1–9. [CrossRef]
- 10. Rossi, R.J.; Gilmartin, K.J. *The Handbook of Social Indicators: Sources, Characteristics, and Analysis;* Garland STPM Press: New York, NY, USA, 1980.
- 11. Bruce, A.; Brown, C.; Avello, P.; Beane, G.; Bristow, J.; Ellis, L.; Fisher, S.; Freeman, S.G.; Jiménez, A.; Leten, J.; et al. Human dimensions of urban water resilience: Perspectives from Cape Town, Kingston upon Hull, Mexico City and Miami. *Water Secur.* **2020**, *9*, 100060. [CrossRef]
- 12. Lee, S.; Yoon, H. Development of disaster risk assessment method in river confluence using AHP. J. Korean Soc. Hazard Mitig. 2018, 18, 545–553. [CrossRef]
- 13. Zagorecki, A.T.; Johnson, D.E.; Ristvej, J. Data mining and machine learning in the context of disaster and crisis management. *Int. J. Emerg. Manag.* **2013**, *9*, 351–365. [CrossRef]
- 14. Yu, J.; Zhao, Q.; Chin, C.S. Extracting Typhoon Disaster Information from VGI Based on Machine Learning. *J. Mar. Sci. Eng.* **2019**, *7*, 318. [CrossRef]
- 15. Chen, J.; Li, Q.; Wang, H.; Deng, M. A machine learning ensemble approach based on random forest and radial basis function neural network for risk evaluation of regional flood disaster: A case study of the Yangtze River Delta, China. *Int. J. Environ. Res. Public Health* **2020**, *17*, 49. [CrossRef] [PubMed]
- Khouj, M.; Lopez, C.; Sarkaria, S.; Marti, J. Disaster management in real time simulation using machine learning. In Proceedings of the 2011 24th Canadian Conference on Electrical and Computer Engineering (CCECE), Niagara Falls, ON, Canada, 8–11 May 2011; pp. 001507–001510.
- 17. Chang, F.J.; Hsu, K.; Chang, L.C. (Eds.) *Flood Forecasting Using Machine Learning Methods*; MDPI: Basel, Switzerland, 2019.
- Chang, F.-J.; Guo, S. Advances in hydrologic forecasts and water resources management. *Water* 2020, *12*, 1819. [CrossRef]
- 19. Kao, I.-F.; Zhou, Y.; Chang, L.-C.; Chang, F.-J. Exploring a Long Short-Term Memory based Encoder-Decoder framework for multi-step-ahead flood forecasting. *J. Hydrol.* **2020**, *583*, 124631. [CrossRef]
- 20. Khan, A.; Gupta, S.; Gupta, S.K. Multi-hazard disaster studies: Monitoring, detection, recovery, and management, based on emerging technologies and optimal techniques. *Int. J. Disast. Risk Reduct.* **2020**, *47*, 101642. [CrossRef]
- Ofli, F.; Meier, P.; Imran, M.; Castillo, C.; Tuia, D.; Rey, N.; Briant, J.; Millet, P.; Reinhard, F.; Parkan, M. Combining human computing and machine learning to make sense of big (aerial) data for disaster response. *Big Data* 2016, *4*, 47–59. [CrossRef]
- 22. Park, J.; Kim, K.T.; Lee, W.H. Recent Advances in Information and Communications Technology (ICT) and Sensor Technology for Monitoring Water Quality. *Water* **2020**, *12*, 510. [CrossRef]
- 23. Orencio, P.M.; Fujii, M. A localized disaster-resilience index to assess coastal communities based on an analytic hierarchy process (AHP). *Int. J. Disast. Risk Reduct.* **2013**, *3*, 62–75. [CrossRef]
- 24. Sheykhmousa, M.; Kerle, N.; Kuffer, M.; Ghaffarian, S. Post-disaster recovery assessment with machine learning-derived land cover and land use information. *Remote Sens.* **2019**, *11*, 1174. [CrossRef]
- Resch, B.; Usländer, F.; Havas, C. Combining machine-learning topic models and spatiotemporal analysis of social media data for disaster footprint and damage assessment. *Cartogr. Geogr. Inf. Sci.* 2018, 45, 362–376. [CrossRef]
- 26. Ragini, J.R.; Anand, P.R.; Bhaskar, V. Big data analytics for disaster response and recovery through sentiment analysis. *Int. J. Inf. Manag.* **2018**, *42*, 13–24. [CrossRef]
- 27. Zhang, Y.; Zhang, R.; Ma, Q.; Wang, Y.; Wang, Q.; Huang, Z.; Huang, L. A feature selection and multi-model fusion-based approach of predicting air quality. *ISA Trans.* **2020**, *100*, 210–220. [CrossRef] [PubMed]
- Bi, C.; Fu, B.; Chen, J.; Zhao, Y.; Yang, L.; Duan, Y.; Shi, Y. Machine learning based fast multi-layer liquefaction disaster assessment. *World Wide Web* 2019, 22, 1935–1950. [CrossRef]
- 29. Cao, A.-T.; Tran, T.-T.; Nguyen, T.-H.-X.; Kim, D. Simplified Approach for Seismic Risk Assessment of Cabinet Facility in Nuclear Power Plants Based on Cumulative Absolute Velocity. *Nucl. Technol.* **2020**, 206, 743–757. [CrossRef]
- 30. Korea Meteorological Administration Information Portal. Available online: https://data.kma.go.kr (accessed on 28 March 2020).
- Korea Ministry of the Interior and Safety Information Portal. Available online: http://lofin.mois.go.kr/portal/ main.do (accessed on 15 April 2020).

- 32. Korea Ministry of Environment (MOE). 2018 Statics of Waterworks; MOE: Sejong, Korea, 2020.
- 33. Korea Ministry of Land, Infrastructure and Transport (MOLIT). *Korea Design Standard*; MOLIT: Sejong, Korea, 2016; p. 45.
- 34. Razmkhah, H.; Abrishamchi, A.; Torkian, A. Evaluation of spatial and temporal variation in water quality by pattern recognition techniques: A case study on Jajrood River (Tehran, Iran). *J. Environ. Manag.* **2010**, *91*, 852–860. [CrossRef]
- 35. Tripathi, M.; Singal, S.K. Use of Principal Component Analysis for parameter selection for development of a novel Water Quality Index: A case study of river Ganga India. *Ecol. Indic.* **2019**, *96*, 430–436. [CrossRef]
- 36. Sahoo, M.M.; Patra, K.; Khatua, K. Inference of water quality index using ANFIA and PCA. *Aquat. Procedia* **2015**, *4*, 1099–1106. [CrossRef]
- 37. Saaty, T.L. The Analytic Hierarchy Process; Mcgraw Hill: New York, NY, USA, 1980.
- 38. Wind, Y.; Saaty, T.L. Marketing applications of the analytic hierarchy process. *Manag. Sci.* **1980**, *26*, 641–658. [CrossRef]
- Chakraborty, S.; Kumar, R.N. Assessment of groundwater quality at a MSW landfill site using standard and AHP based water quality index: A case study from Ranchi, Jharkhand, India. *Environ. Monit. Assess.* 2016, 188, 335. [CrossRef]
- 40. Saaty, T.L. How to make a decision: The analytic hierarchy process. *Eur. J. Oper. Res.* **1990**, *48*, 9–26. [CrossRef]
- 41. Saaty, R.W. The analytic hierarchy process—What it is and how it is used. *Math. Model.* **1987**, *9*, 161–176. [CrossRef]
- 42. Saaty, T.L. Priority setting in complex problems. IEEE Trans. Eng. Manag. 1983, 3, 140–155. [CrossRef]
- 43. Uddameri, V.; Silva, A.L.B.; Singaraju, S.; Mohammadi, G.; Hernandez, E.A. Tree-Based Modeling Methods to Predict Nitrate Exceedances in the Ogallala Aquifer in Texas. *Water* **2020**, *12*, 1023. [CrossRef]
- 44. Shin, Y.; Kim, T.; Hong, S.; Lee, S.; Lee, E.; Hong, S.; Lee, C.; Kim, T.; Park, M.S.; Park, J. Prediction of Chlorophyll-a Concentrations in the Nakdong River Using Machine Learning Methods. *Water* **2020**, *12*, 1822. [CrossRef]
- 45. Zhang, D.; Qian, L.; Mao, B.; Huang, C.; Huang, B.; Si, Y. A data-driven design for fault detection of wind turbines using random forests and XGboost. *IEEE Access* **2018**, *6*, 21020–21031. [CrossRef]
- 46. Breiman, L. Random forests. Mach. Learn. 2001, 45, 5–32. [CrossRef]
- 47. Genuer, R.; Poggi, J.-M.; Tuleau-Malot, C. Variable selection using random forests. *Pattern Recognit. Lett.* **2010**, *31*, 2225–2236. [CrossRef]
- 48. Hollister, J.W.; Milstead, W.B.; Kreakie, B.J. Modeling lake trophic state: A random forest approach. *Ecosphere* **2016**, 7, e01321. [CrossRef]
- 49. Fernández-Delgado, M.; Cernadas, E.; Barro, S.; Amorim, D. Do we need hundreds of classifiers to solve real world classification problems? *J. Mach. Learn. Res.* **2014**, *15*, 3133–3181.
- 50. Singh, B.; Sihag, P.; Singh, K. Modelling of impact of water quality on infiltration rate of soil by random forest regression. *Model. Earth Syst. Environ.* **2017**, *3*, 999–1004. [CrossRef]
- 51. Read, E.K.; Patil, V.P.; Oliver, S.K.; Hetherington, A.L.; Brentrup, J.A.; Zwart, J.A.; Winters, K.M.; Corman, J.R.; Nodine, E.R.; Woolway, R.I. The importance of lake-specific characteristics for water quality across the continental United States. *Ecol. Appl.* **2015**, *25*, 943–955. [CrossRef] [PubMed]
- 52. Friedman, J.H. Greedy function approximation: A gradient boosting machine. *Ann. Stat.* **2001**, 1189–1232. [CrossRef]
- 53. Ke, G.; Meng, Q.; Finley, T.; Wang, T.; Chen, W.; Ma, W.; Ye, Q.; Liu, T.-Y. Lightgbm: A highly efficient gradient boosting decision tree. In Proceedings of the Advances in Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA, 4–9 December 2017; pp. 3146–3154.
- Chen, T.; Guestrin, C. Xgboost: A scalable tree boosting system. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, 13–17 August 2016; pp. 785–794.
- 55. Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V. Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **2011**, *12*, 2825–2830.
- 56. XGBoost. Available online: https://xgboost.readthedocs.io/en/latest/build.html (accessed on 15 February 2020).
- 57. Scikit-Learn. Available online: https://scikit-learn.org/stable/index.html (accessed on 3 January 2020).

- 58. Fabris, F.; Doherty, A.; Palmer, D.; De Magalhães, J.P.; Freitas, A.A. A new approach for interpreting random forest models and its application to the biology of ageing. *Bioinformatics* **2018**, *34*, 2449–2456. [CrossRef]
- 59. Grömping, U. Variable importance assessment in regression: Linear regression versus random forest. *Am. Stat.* **2009**, *63*, 308–319. [CrossRef]
- Moriasi, D.N.; Arnold, J.G.; Van Liew, M.W.; Bingner, R.L.; Harmel, R.D.; Veith, T.L. Model evaluation guidelines for systematic quantification of accuracy in watershed simulations. *Trans. ASABE* 2007, 50, 885–900. [CrossRef]
- 61. Bennett, N.D.; Croke, B.F.; Guariso, G.; Guillaume, J.H.; Hamilton, S.H.; Jakeman, A.J.; Marsili-Libelli, S.; Newham, L.T.; Norton, J.P.; Perrin, C. Characterising performance of environmental models. *Environ. Model. Softw.* **2013**, *40*, 1–20. [CrossRef]

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



© 2020 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (http://creativecommons.org/licenses/by/4.0/).