



# Sensitivity of hydrological machine learning prediction accuracy to information quantity and quality

Minhyuk Jeung<sup>1</sup>, Younggu Her<sup>2</sup>, Sang-Soo Baek<sup>3</sup>, Kwangsik Yoon<sup>1</sup>

<sup>1</sup> Department of Rural & Biosystems Engineering (Brain Korea 21), Chonnam National University, Gwangju 61186, Republic of Korea

<sup>2</sup> Department of Agricultural and Biological Engineering / Tropical Research and Education Center, University of Florida, Homestead, Florida 33186, USA

<sup>3</sup> Department of Environmental Engineering, Yeungnam University, Gyeongsan 38541, Republic of Korea

Correspondence to: Kwangsik Yoon (ksyoon@chonnam.ac.kr)

10 **Abstract.** Machine learning (ML) is now commonly employed as a tool for hydrological prediction due to recent advances  
in computing resources and increases in data volume. The prediction accuracy of ML (or data-driven) modeling is known to  
be improved through training with additional data; however, the improvement mechanism needs to be better understood and  
documented. This study explores the connection between the amount of information contained in the data used to train an  
ML model and the model's prediction accuracy. The amount of information was quantified using Shannon's information  
15 theory, including marginal and transfer entropy. Three ML models were trained to predict the flow discharge, sediment, total  
nitrogen, and total phosphorus loads of four watersheds. The amount of information contained in the training data was  
increased by sequentially adding weather data and the simulation outputs of uncalibrated and/or calibrated mechanistic (or  
theory-driven) models. The reliability of training data was considered a surrogate of information quality, and accuracy  
statistics were used to measure the quality (or reliability) of the uncalibrated and calibrated theory-driven modeling outputs  
20 to be provided as training data for ML modeling. The results demonstrated that the prediction accuracy of hydrological ML  
modeling depends on the quality and quantity of information contained in the training data. The use of all types of training  
data provided the best hydrological ML prediction accuracy. ML models trained only with weather data and calibrated  
theory-driven modeling outputs could most efficiently improve accuracy in terms of information use. This study thus  
illustrates how a theory-driven approach can help improve the accuracy of data-driven modeling by providing quality  
25 information about a system of interest.

## 1 Introduction

Machine learning (ML) techniques have become commonly employed for hydrological prediction due to the availability  
of large hydrological data repositories and advances in computing resources and techniques (Sun et al., 2020; Xu and Liang,  
2021). Studies have demonstrated that ML techniques can predict hydrological variables as accurately or even better than  
30 other statistical methods and mechanistic (or theory-driven) modeling (Panidhapu et al., 2020). The prediction accuracy of



ML modeling is known to increase with the volume of data used to train the models (Jha et al., 2018); as such, the accuracy is expected to improve further as hydrological observations and records accumulate over time. However, it remains unclear how prediction accuracy is associated with the characteristics of training data: can any data added to a training set improve the accuracy?

35 Information theory has served as a mathematical tool to measure the amount of information contained in data and its transfer to another set of data (Shannon, 1948a; Shannon, 1948b). This tool can help us understand the correlations or dependencies among multiple interconnected data sets (Pechlivanidis et al., 2018), which helps determine whether the training data contains information that could improve the accuracy of the model (Nearing et al., 2020). Shannon's entropy, often called marginal entropy (ME), is one of the most commonly used information theories that can quantify information content in a set of data (Silva et al., 2017). The concept of transfer entropy (TE) was proposed to measure the amount of information transferred from one variable to another (Schreiber, 2000). Previous studies have employed ME to quantify the amount of information in hydrological datasets (Silva et al., 2017) and TE to qualify the interactions between input and output data in hydrological analyses (Bennett et al., 2019; Konapala et al., 2020). Both ME and TE have great potential as concepts and methods to evaluate the informatic characteristics of training data and their impacts on hydrological ML model performance.

45 Data-driven methods, including ML modeling, rely on historical records and estimates from other analyses, while theory-driven approaches employ existing hydrological concepts and knowledge for prediction. Mechanistic modeling can be classified as a theory-driven method even when its parameter calibration has the nature of a data-driven approach. Mechanistic models employ different assumptions, knowledge, and methods to conceptualize a hydrological system of interest, which is why they provide unique predictions. For example, streamflow hydrographs predicted using Hortonian and Dunne's concepts might be substantially different from each other even after parameter calibration (Loague et al., 2010). Information embedded in hydrological theories and models can help improve the performance of data-driven modeling, and the information is considered in the predictions of mechanistic modeling. Weather records are one of the data sets commonly used to train hydrological ML models (Chen et al., 2020). Previous studies have demonstrated that ML models trained only with meteorological data provide limited accuracy; this is unsurprising given that hydrological processes are usually complicated by many other factors, including topography, soil, land use and cover, geological features, and management practices (Srinivasan et al., 2010; Srivastava et al., 2020). Hence, mechanistic model predictions can be an alternative source of data for the training of hydrological ML models.

60 Mechanistic modeling often or always requires parameter calibration to consider the hydrological characteristics of an area of interest. In a technical sense, parameter calibration is an effort to improve the statistical similarity between observed and predicted variables of interest. The prediction accuracy of mechanistic modeling is usually improved through the calibration process. As a result, the amount and/or quality of information in a relatively accurate prediction may be greater and/or higher than that of information in a relatively inaccurate one. When the prediction accuracy of mechanistic modeling is improved by calibrating its parameters, the calibrated model may have more and/or better-quality information than the uncalibrated



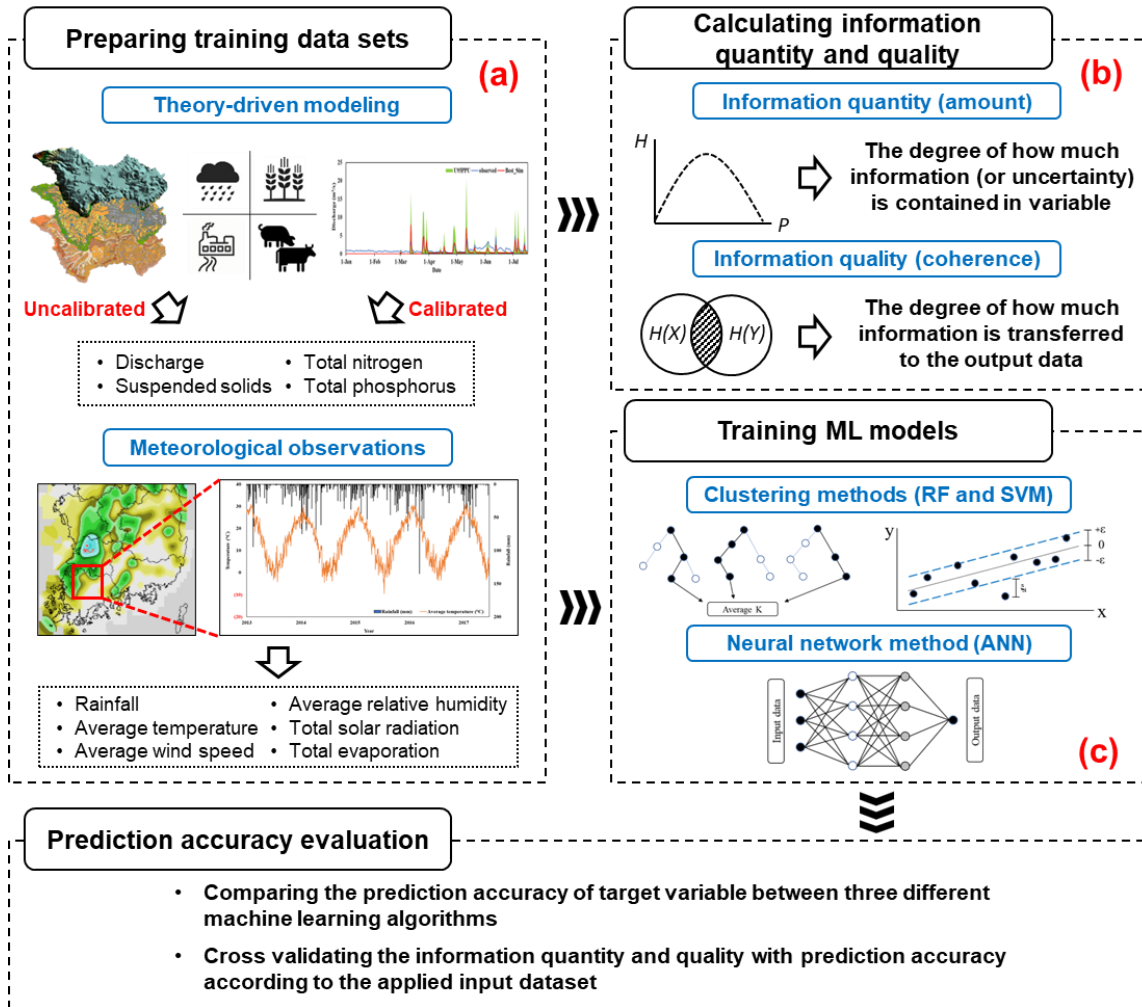
65 model. Thus, a pair of uncalibrated and calibrated mechanistic models for a watershed can be a useful tool to create training  
data sets with different amounts and/or qualities of information for hydrological ML modeling.

This study attempted to relate the quantity and quality of information contained in sets of training data to the prediction  
accuracy of hydrological ML models, with the goal of understanding how to improve the accuracy efficiently. Information  
quantity and quality were quantified using information theory, including ME and TE statistics. Three different ML  
70 algorithms (or models) were tested in the evaluation. The quantity of information was systematically increased by adding  
weather records and uncalibrated and calibrated theory-driven (or mechanistic) model outputs to training data sets. This  
study employed a mechanistic model commonly used to predict flow and water quality to represent a theory-driven approach.  
The data-driven (i.e., ML) and theory-driven (i.e., mechanistic) modeling approaches were applied to predict the flow,  
sediment, total nitrogen (TN), and total phosphorus (TP) loads of four watersheds. Then, the implications of the evaluation  
75 results and the limitations of this study were discussed, and the future direction of hydrological ML modeling was suggested  
based on the findings.

## 2 Methods and Materials

### 2.1 Overall procedure

Three ML models, including Random Forest (RF), Support Vector Machine (SVM), and Artificial Neural Network (ANN),  
80 were employed in this study (Fig. 1). The ML models were first trained using data sets collected from four study watersheds,  
including weather data observed at weather stations within or close to the watersheds and flow discharge, suspended soils  
(SS), TN, and TP measured at the outlets of the four study watersheds. The Soil and Water Assessment Tool (SWAT) model  
was selected to represent a mechanistic (or theory-driven) model. The SWAT model was used to produce additional data sets  
to which the ML models would be trained. The SWAT models were calibrated to flow (or streamflow discharges), SS, TN,  
85 and TP loads measured at the outlets. Then, the outputs (i.e., flow discharge, SS, TN, and TP loads) of the uncalibrated (i.e.,  
SWAT models with the default parameter values) and calibrated SWAT models were used as additional data sets for the  
training of the three ML models. Here, we assumed there would be a difference in the quality of information contained in the  
uncalibrated and calibrated outputs of the SWAT model, and we employed the marginal and transfer entropy based on  
Shannon's information theory, to quantify the amount and quality of information contained in the training data sets. Finally,  
90 the weather data and uncalibrated and calibrated SWAT model outputs were sequentially fed to the three ML models to  
investigate how the ML models' prediction accuracy reacts to the information quantity and quality of the training data (Fig. 1  
and Table 1). Information use efficiency (IUE) was quantified to evaluate the performance of the three ML models, and the  
four training data sets (Table 1) based on the amount of prediction accuracy improvement made using a unit increase in  
entropy.



95

**Figure 1.** Overall procedure to investigate the contribution of information quantity and quality to the prediction accuracy of hydrological machine learning (ML) modeling.



**Table 1.** Combinations of data sets used to train hydrological ML models.

Training Data Sets	Variables
WDO	P, AT, WS, RH, SR, E
WD+UC	P, AT, WS, RH, SR, E, Q_Uncal*, SS_load_Uncal*, TN_load_Uncal*, TP_load_Uncal*
WD+C	P, AT, WS, RH, SR, E, Q_Cal*, SS_load_Cal*, TN_load_Cal*, TP_load_Cal*
All	P, AT, WS, RH, SR, E, Q_Uncal*, SS_load_Uncal*, TN_load_Uncal*, TP_load_Uncal*, Q_Cal*, SS_load_Cal*, TN_load_Cal*, TP_load_Cal*

\* P, AT, WS, RH, SR, and E represent precipitation, average temperature, wind speed, relative humidity, solar radiation, and evapotranspiration, respectively. The SWAT outputs, including Q, SS load, TN load, and TP load, were used to train ML models separately depending on the target variables. For example, SWAT's SS load simulation results were used only when predicting SS load using ML models.

## 2.2 Data-driven (or machine learning) models

The RF model is based on a regression tree, but it differs as it does not grow with a single tree but rather an entire forest of numerous trees using the bootstrap aggregating technique or bagging technique to help decrease model variance (Breiman et al., 1984; Breiman, 2001). The RF model is known for its ability to be used when there are more variables than observation data, and it does not result in overfitting due to the pruning process (Diaz-Uriate and de Andrés, 2006). The RF model was also reported to offer excellent performance even when predictive variables are irregular (Diaz-Uriate and de Andrés, 2006). In RF modeling, a decision tree grows by splitting a tree node, and it is pruned by removing tree nodes or sections with relatively low explanatory power compared to others (Hasanipanah et al., 2017).

The SVM model divides a high- or infinite-dimensional space using hyperplanes until all data points are separated (Vapnik, 1995; 1998). The SVM model is known to be able to avoid overfitting and produces highly accurate predictions (Aktan, 2011). The goal of the SVM procedures is to identify the optimal hyperplane separating two classes in the high-dimensional space that maximizes the distance between the two data point groups (Ahmed et al., 2017). SVM modeling transforms training data using the kernel function so that a linear hyperplane can separate the data points in high dimensions. Three kernel functions are commonly used: radial basis function (RBF), linear function, and polynomial function. This study employed the RBF, which is the most widely used kernel function (Tao et al., 2008).

The ANN model has been widely used to solve various modeling problems (Khashei and Bijari, 2010). The structure of the ANN model was inspired by the biological structure of the human brain, which is composed of many interconnected



processing elements called neurons (Tosun et al., 2016). The structure is characterized by a network of three layers: input, hidden, and output. The number of input and hidden layers is determined by the number of input variables and the complexity of the problem (Yilmazkaya et al., 2018). Neurons are a critical parameter used in interconnected processing, which is characterized by weights (Tosun et al., 2016). The weights of individual neurons determine how input values are transferred to other values on the output nodes. The weights of connections between layers are calculated by the backpropagation process, which calculates the gradient of prediction error with respect to weights (Siddique and Tokhi, 2001).

### 2.3 Data normalization and accuracy evaluation

ML modeling is known to have low learning rates when some types of training data have value ranges substantially different from those of others (Ioffe and Szegedy, 2015). Data normalization techniques are commonly used to rescale the training data from their original ranges into a common value range so that the ML models can be efficiently and quickly trained. Several data normalization methods are available; linear scaling is one of the most widely used, presumably due to its simplicity and efficacy (Raju et al., 2020; Eq. 1).

$$X' = (x - \min(x)) / (\max(x) - \min(x)) \quad (1)$$

where  $X'$  is the normalized value of the data set (ranges from 0 to 1), and  $x$  is an original value.

The prediction accuracy of the three ML models was evaluated using the Kling-Gupta efficiency coefficient (KGE; Gupta et al., 2009). The KGE considers the strength of the correlation between observed and predicted variables while also comparing the variables' biases and variances. Thus, compared to the Nash-Sutcliffe efficiency and the coefficient of determination, the KGE is less sensitive to relatively large values that lead to biases toward such values (Nash and Sutcliffe, 1970; Gupta et al., 2009; Eq. 2).

$$KGE = 1 - \sqrt{(r - 1)^2 + \left(\frac{\sigma_{sim}}{\sigma_{obs}} - 1\right)^2 + \left(\frac{\mu_{sim}}{\mu_{obs}} - 1\right)^2} \quad (2)$$

where  $\sigma_{obs}$  and  $\sigma_{sim}$  are the standard deviations of observations and simulation results, respectively, and  $\mu_{obs}$  and  $\mu_{sim}$  are the averages of observed and simulated variables, respectively.

A KGE of 1 indicates perfect agreement between observations and predictions (Andersson et al., 2017). Knoben et al. (2019) mathematically demonstrated that the KGE value approaches -0.41 when the predicted (or simulated) values of a variable are equal to the average value of its observations. Thus, a KGE value of -0.41 can be interpreted similarly to an NSE value of 0.00, meaning that the predictions may not be a better than the observed mean (Schaeffli and Gupta, 2007). In this study, we assumed that predictions would be acceptable or satisfactory when the differences between observed and simulated averages of a variable (or percentage biases) and the variances of the differences are less than 25% for flow, 55% for SS, and 70% for TN/TP (Moriassi et al., 2007), which correspond to KGEs of 0.54, 0.17, and -0.03 for flow, SS, and TN/TP, respectively, with an arbitrarily selected threshold correlation of 0.30.



## 2.4 Theory-driven (or mechanistic) model

The SWAT model was designed to predict watershed processes based on theories and known mechanisms that control the generation and transport processes of water, sediment, and nutrients (Nietsch et al., 2002). The SWAT model is popularly used to predict water and nutrient loadings at the watershed and basin scales due to its proven applicability to a variety of landscapes and climate zones as well as its simple but defensible modeling strategies. Moreover, the SWAT model can consider various management practices, including application rates and timing of fertilizers and herbicides/pesticides; tillage and low-impact development practices; and agricultural conservation practices such as filter strips, nutrient management plans, terraces, and tile drainage (Her et al., 2017; Her and Jeong, 2018; Li et al., 2021a). Several studies have attempted to improve the prediction accuracy of SWAT modeling by coupling it with ML techniques, for example, to predict peak flow (Senent-Aparicio et al., 2019), water quality (Noori et al., 2020), and aquifer vulnerability (Jang et al., 2020).

Two versions of the SWAT model, namely uncalibrated and calibrated mechanistic modeling outputs, were prepared to generate two sets of training data for the ML models. The agricultural management practices that were compiled from the study watersheds were incorporated into both models (RDA, 2014). The values of all parameters of the uncalibrated SWAT model remained unchanged; thus, the uncalibrated SWAT models do not necessarily represent the hydrological processes of the study watersheds, and they are not likely to reproduce the observed flow, SS, TN, and TP at an acceptable accuracy level. Accordingly, the quality of information contained in the outputs of the uncalibrated SWAT models may be relatively low compared to that of the calibrated SWAT models. The parameter values of the SWAT models were calibrated to flow, SS, TN, and TP observations made at the study watersheds' outlets. The flow, SS, TN, and TP loads predicted using the calibrated SWAT models were assumed to have relatively high-quality information compared to those of the uncalibrated SWAT model. The quantity and quality of information were quantified using the marginal and transfer entropies described in the following section.

The SUFI-2 algorithm, widely used for SWAT model calibration, was used to explore the multi-dimensional parameter spaces of the SWAT models and locate a solution (or a parameter set) close to the global optimum in this study (Sao et al., 2020). The simulation period was split into three: a warm-up period from January 1, 2008, to July 11, 2013; a calibration period from July 12, 2013, to December 31, 2015; and a validation period from January 1, 2016, to December 31, 2017. The types and value ranges of the calibration parameters were determined based on the previous SWAT modeling experience, the understanding of the calibration parameters, and the literature (Tobin and Bennett, 2017; Tang et al., 2021).

## 2.5 Marginal and transfer entropy

This study measured the quantity and quality of information contained in the training data using ME and TE. In general, a data set that is spread out has relatively high entropy, while another data set that is concentrated on a small range of values has relatively small entropy. The ME is defined as the information content of a variable and used to calculate randomness in time series using Eq. 3 (Shannon, 1948; Cover and Thomas, 2006; Silva et al., 2017):



$$H(X) = -\sum_{i=1}^n p(x_i) \log_2 P(x_i) \quad (3)$$

185 where  $H(X)$  is a measure of information of a discrete random variable  $X$ , and  $P(x)$  is the probability mass function of variable  $x$  in the  $i^{\text{th}}$  step.

While the amount of information contained in a variable can be calculated using the ME, we can also calculate the amount of information shared between two variables based on mutual information theory using Eq. 4 (Cover and Thomas, 2006):

$$I(X; Y) = H(X) + H(Y) - H(X, Y) \quad (4)$$

190 where  $I(X, Y)$  is the quantified value between  $X$  and  $Y$ . The mutual information  $I(X, Y)$  represents the expected information gained in  $Y$  from measuring  $X$ , or vice versa. From these definitions, we can calculate the conditional entropy by subtracting the amount of information shared between  $X$  and  $Y$  from  $H(X)$ , which indicates how much information remains about the entire time series  $X$  in case we already know the information content of  $Y$ .

$$H(X|Y) = H(X) - I(X; Y) \quad (5)$$

195 These quantities are all symmetrical and do not explain the amount of information exchanged between variables (Bennett et al., 2019). The TE was devised to consider the asymmetric transfer of information between any two-time series  $X$  and  $Y$  (the information flow from one to another variable), and can be defined as conditional mutual information (Schreiber, 2000):

$$T_{X \rightarrow Y} = I(Y_t; X_t | Y_t) \quad (6)$$

200 where  $T_{X \rightarrow Y}$  is the transfer entropy from  $X$  to  $Y$ , and  $X_t$  or  $Y_t$  denotes the variables  $X$  and  $Y$  in time  $t$ . Once the ME and TE were calculated for the modeling experiments with the unique combinations of the ML models and the training data sets (Fig. 1 and Table 1), the prediction accuracy gain was divided by the increases in the quantity (ME) and quality (TE) of information contained in the training data to calculate information use efficiency (IUE):

$$IUE_{ME} = \frac{P_{WD \rightarrow ID}}{\sum H(x_{WD} \rightarrow x_{ID})} \quad (7)$$

$$IUE_{TE} = \frac{P_{WD \rightarrow ID}}{\sum T_{WD_i \rightarrow Y} \rightarrow \sum T_{ID_i \rightarrow Y}} \quad (8)$$

205 where  $P_{WD \rightarrow ID}$  is the prediction accuracy gain or increase from using additional straining data sets, as compared to the case of only using weather data for the training. The  $H(x_{WD \rightarrow x_{ID}})$  and  $\sum T_{WD_i \rightarrow Y} \rightarrow \sum T_{ID_i \rightarrow Y}$  denotes the marginal entropy and transfer entropy gain or increase from using additional straining data sets, compared to the case of only using weather data for the training.

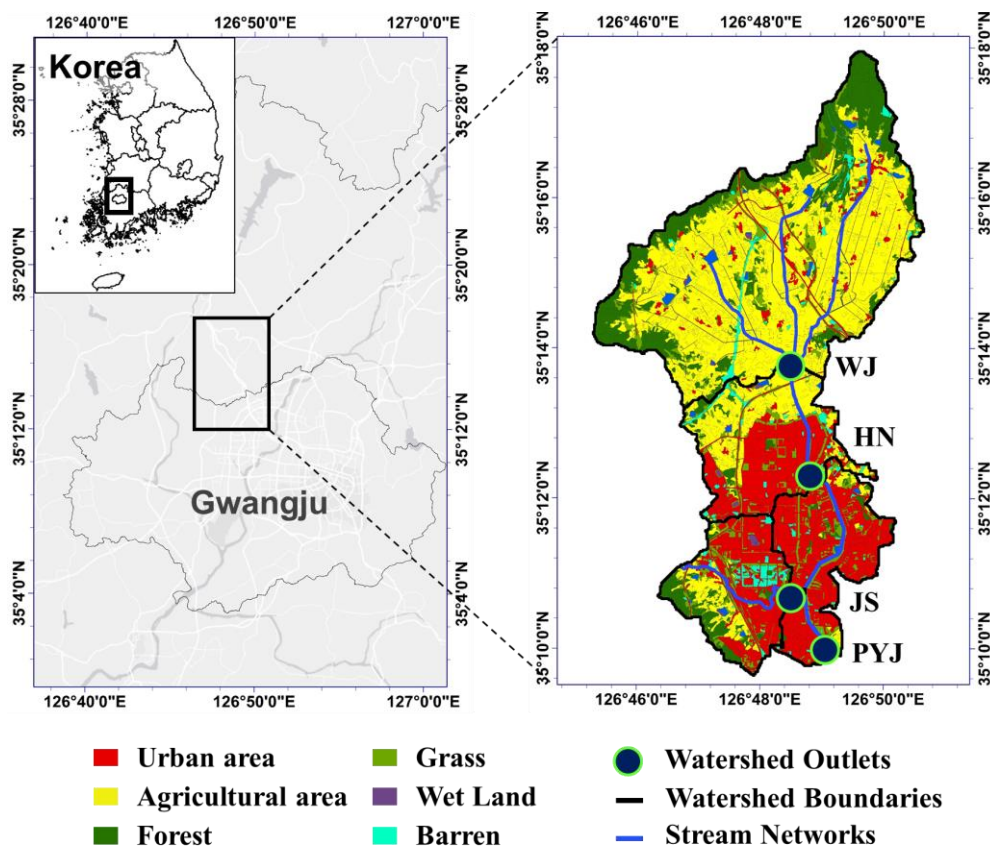
## 2.6 Study watersheds and training data acquisition

210 The Pung-Yeong-Jung (PYJ) river watershed was selected for the modeling experiment of this study. The PYJ watershed can be divided into three sub-watersheds from upstream to downstream: the Wall-Jeong (WJ), Ha-Nam (HN), and Jang-Su (JS) watersheds (Fig. 2 and Table S1). The WJ watershed is nested by the HN watershed, and the HN and JS watersheds are nested by the PYJ watershed. Thus, all direct runoff drained from the three nested watersheds passes the outlet (35°09'58.87" N, 126°49'08.93" E) of the PYJ watershed. The streamflow, SS, TN, and TP concentrations were monitored at the outlets of





215 the four study watersheds for four years and six months, from July 12, 2013, to December 31, 2017. Most of the drainage areas were covered by agricultural land uses, including upland and rice paddy fields (covering 41% of the JS watershed and 62% of the WJ watershed) and forest. Urbanized areas cover 5% (WJ watershed) to 31% (JS watershed) of the watersheds.



220 **Figure 2.** Location of the study watersheds and their land uses and covers.

The Korean Meteorological Administration monitors weather variables, including daily AT, P, E, WS, RH, and SR, at a weather station located approximately 7 km away from the study watersheds. Water pressure sensors and data loggers (OTT Orpheus Mini, Germany) were deployed at the monitoring sites close to the watershed outlets. The cross-sections of the streams were surveyed at the monitoring sites, and the velocity of streamflow was then measured using a flow meter (VALEPORT model 002, UK) across the sections to estimate flow discharge rates. Water quality samples were manually collected every eight days. During a rainfall event, stream water was collected using an automatic sampler (ISCO portable sampler 6712, USA), and the sampling interval was reduced to 1 hour to catch the expected large variations of flow rates and the corresponding water quality concentrations for improved observation accuracy. During the monitoring period, a total of 230 17 large rainfall events were sampled.

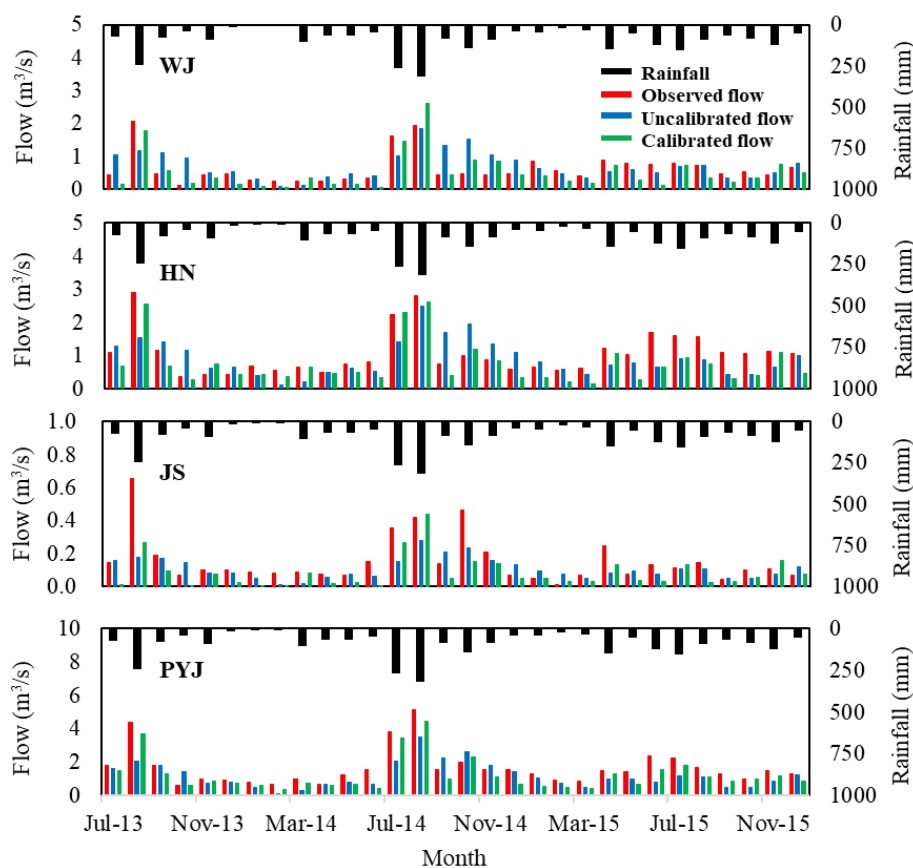


### 3 Results

#### 3.1 Training data: Weather records and monitoring data

Sets of training data were prepared using the daily weather records, the watershed monitoring data, and SWAT modeling results (uncalibrated and calibrated outputs; Figs. 3 and S1–S4). The watersheds have four seasons, with relatively short springs and falls. The watersheds are fairly wet in the summer and dry in the spring. For example, the watersheds receive precipitation of 831–1333 mm annually, with more than half (59% on average) of the precipitation occurring in the summer (from June to September). In spring, the stream might dry up due to the small amount of precipitation and warm air. In the case of the PYJ watershed, streamflow discharges can be large, with as much as 2.64 m<sup>3</sup>/s on average in summer, but they are limited (e.g., 1.21 m<sup>3</sup>/s) enough to reveal the bottom of the stream in spring (from March to May).

240



**Figure 3.** Comparison of monthly streamflow predicted using the mechanistic models (i.e., uncalibrated and calibrated SWAT models) and observed during the training period (July 12, 2013, to December 31, 2015). The daily-scale comparisons can be found in the supplementary document (Figs. S1–S4).



245

The PYJ and JS watersheds had the largest and smallest average daily discharge of 1.69 and 0.16 m<sup>3</sup>/s, respectively (Table S2). The JS watershed had relatively higher SS concentrations compared to the other watersheds as it includes large construction sites (Mendie, 2005; Pullanikkatil et al., 2015; Adeola-Fashae et al., 2019). In addition, the first flush effects of the urbanized watersheds (e.g., the JS watershed; Table S1) led to higher peak SS, TN, and TP concentrations (Chaudhary et al., 2022). The WJ and HN watersheds had relatively higher TN and TP concentrations, presumably due to agricultural management activities such as fertilizer application and livestock farming in their large agricultural areas (Liu et al., 2012; Table S2).

250

### 3.2 Training data: Outputs of the mechanistic modeling

The calibrated SWAT model provided acceptable performance in all watersheds (e.g., KGEs equal to or greater than 0.54 for flow, 0.17 for SS, and -0.03 for TN/TP). The average KGE values for all watersheds were 0.68 for flow, 0.45 for SS, 0.40 for TN, and 0.44 for TP (Table 2). However, as expected, the uncalibrated model could not accurately predict the variables; average KGEs for all watersheds were less than 0.41 for flow, 0.02 for SS, -0.20 for TN, and -0.35 for TP. As such, the information quality of the outputs of the calibrated SWAT modeling may be greater than that of the uncalibrated modeling. The quantity and quality of information were evaluated with marginal and transfer entropies.

260

**Table 2.** Accuracy statistics (KGEs) of a theory-driven (or SWAT) model in the training period. The KGE scores that satisfy the acceptable accuracy criteria (i.e., 0.54 for flow, 0.17 for SS, -0.03 for TN/TP) are in bold.

Watershed	Flow		SS		TN		TP	
	Uncal	Cal	Uncal	Cal	Uncal	Cal	Uncal	Cal
WJ	0.49	<b>0.71</b>	<b>0.28</b>	<b>0.52</b>	-0.28	<b>0.41</b>	-0.39	<b>0.43</b>
HN	0.50	<b>0.70</b>	-0.06	<b>0.36</b>	-0.09	<b>0.43</b>	-0.33	<b>0.47</b>
JS	0.18	<b>0.57</b>	-0.35	<b>0.45</b>	-0.44	<b>0.37</b>	-0.41	<b>0.27</b>
PYJ	0.46	<b>0.72</b>	<b>0.22</b>	<b>0.48</b>	<b>0.01</b>	<b>0.40</b>	-0.27	<b>0.57</b>

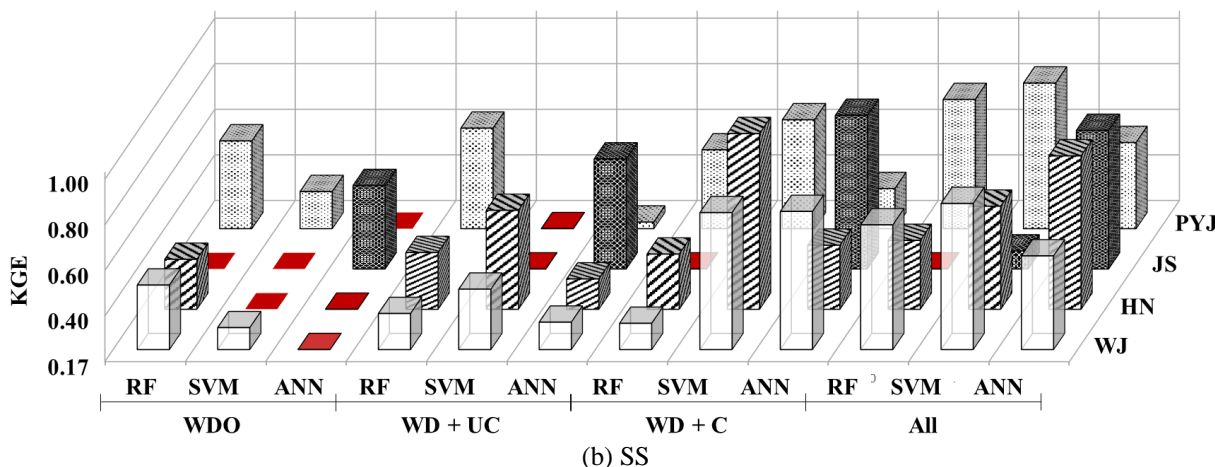
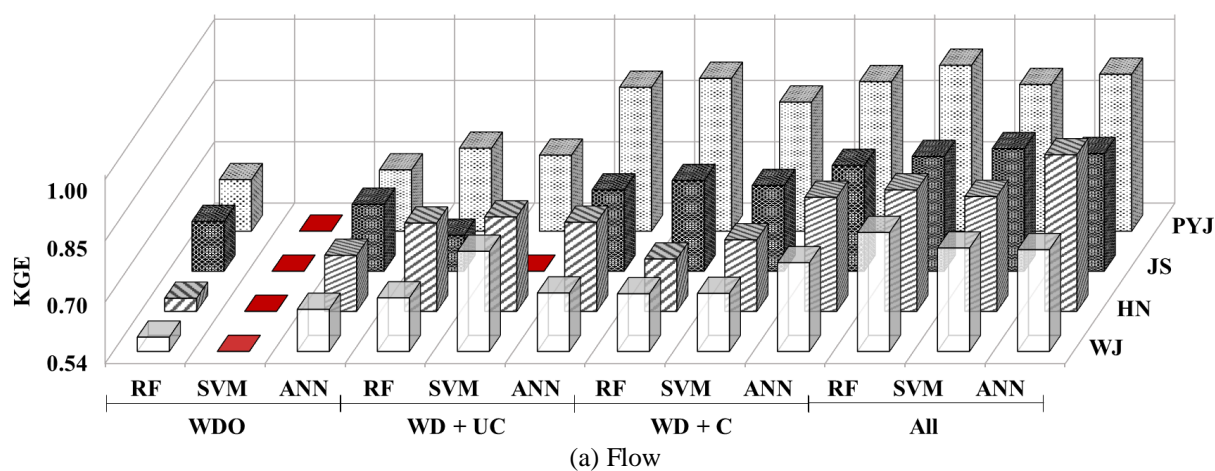
### 3.3 Prediction accuracy of machine learning modeling

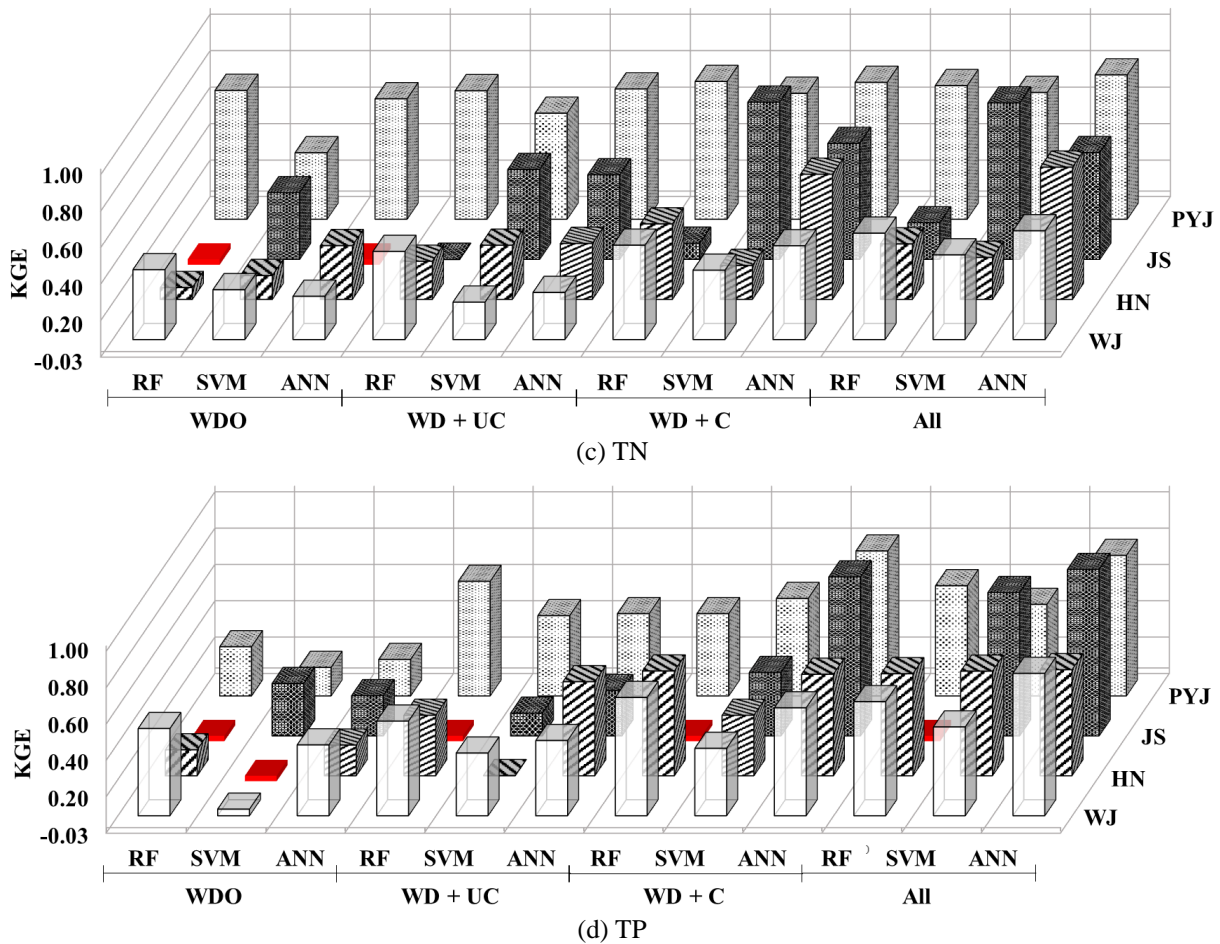
The four ML models were trained with different sets of training data: weather data only (WDO), the uncalibrated SWAT modeling outputs added to WDO (WD+UC), the calibrated SWAT modeling outputs added to WDO (WD+C), and all training data (All or WD+UC+C). The trained ML models yielded unique performances in the predictions depending on the

265



training data set types (Fig. 4). Overall, the ML models' flow prediction accuracy consistently improved as additional data sets were added to the training data, including WDO to WD+UC, WDO+C, and All. For example, the WDO case provided acceptable accuracy (KGE of 0.67 greater than the threshold of 0.54) in the prediction of flow using the RF algorithm at the outlet of the PYJ watershed. When the outputs of the uncalibrated and/or calibrated SWAT modeling were added to the training data, the accuracy of the ML modeling was increased to KGEs of 0.74 (11.6% increase with WD+UC) and 0.91 (37.2% increase with WD+C) in the case of using the RF model. The additional training data sets also improved the accuracy of the water quality ML modeling. However, ML models trained only using the weather data and uncalibrated mechanistic modeling outputs failed to meet the acceptable accuracy levels (i.e., 0.17 for SS and -0.03 for TN/TP; Fig. 4). In Fig. 4, the KGE scores overall increase from left to right. Negative KGE scores are frequently found in the JS watershed, indicating the models relatively poorly performed for the watershed.



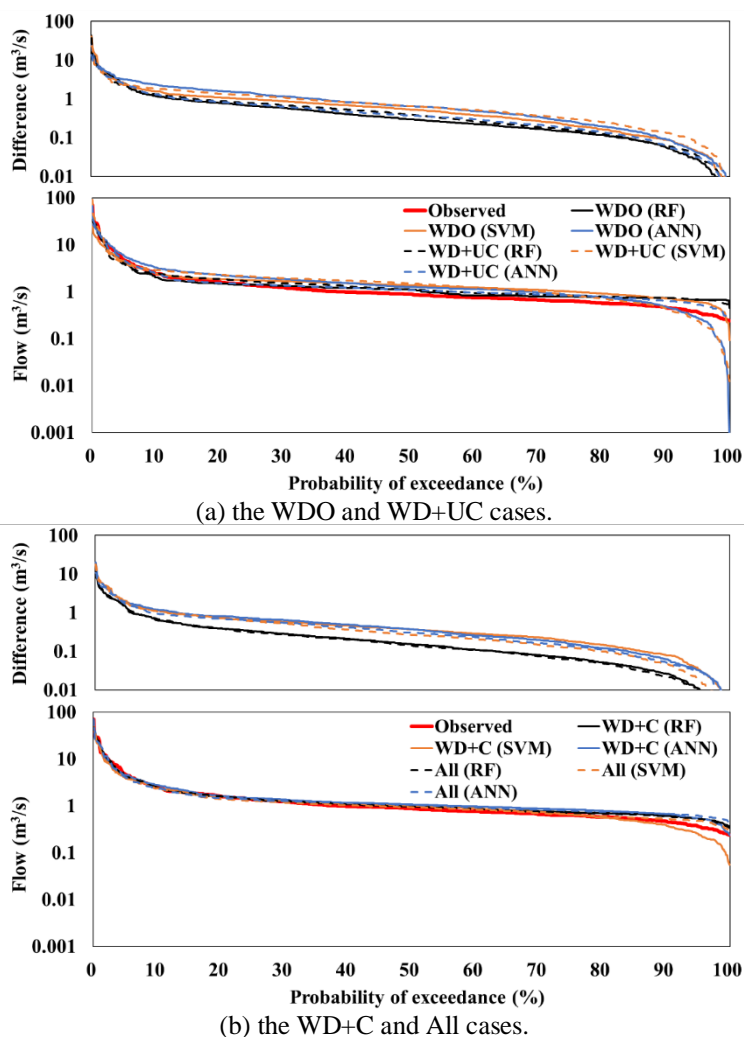


**Figure 4.** Prediction accuracy (KGE) of hydrological ML models trained with the different training data set combinations. The KGE values that do not satisfy the acceptable accuracy levels (e.g., i.e., 0.54 for flow, 0.17 for SS, and -0.03 for TN/TP) are marked with solid red rectangles on the x-y plane.

A flow duration curve (FDC) provides a graphical way of investigating the frequency of extreme events, such as floods and droughts. The FDCs were derived from the observed and predicted flow hydrographs and compared to evaluate prediction accuracy in the frequency domain (Figs. 5 and S5–S7). The FDCs created from flow predictions made using the ML models trained with all training data (the All case) were the closest to the observed FDC in both high (e.g., flooding) and low (e.g., drought) exceedance probability regions. The WDO and WD+UC cases created relatively large differences (under- and over-estimations) between the predicted and observed FDCs, especially for extreme events (i.e., flooding and drought). For example, the differences between the RF predictions and observations for the 5% (flooding) and 95% (drought) exceedance probabilities of the PYJ watershed were 12.1% and 49.9% in the All case respectively, and they increased to 23.0% and 108.6% in the WD+UC case. The findings indicate that the ML models trained with all available training data



sets (the All case) can more accurately predict the extremes than the relatively less trained ML models (the WDO and WD+UC cases).



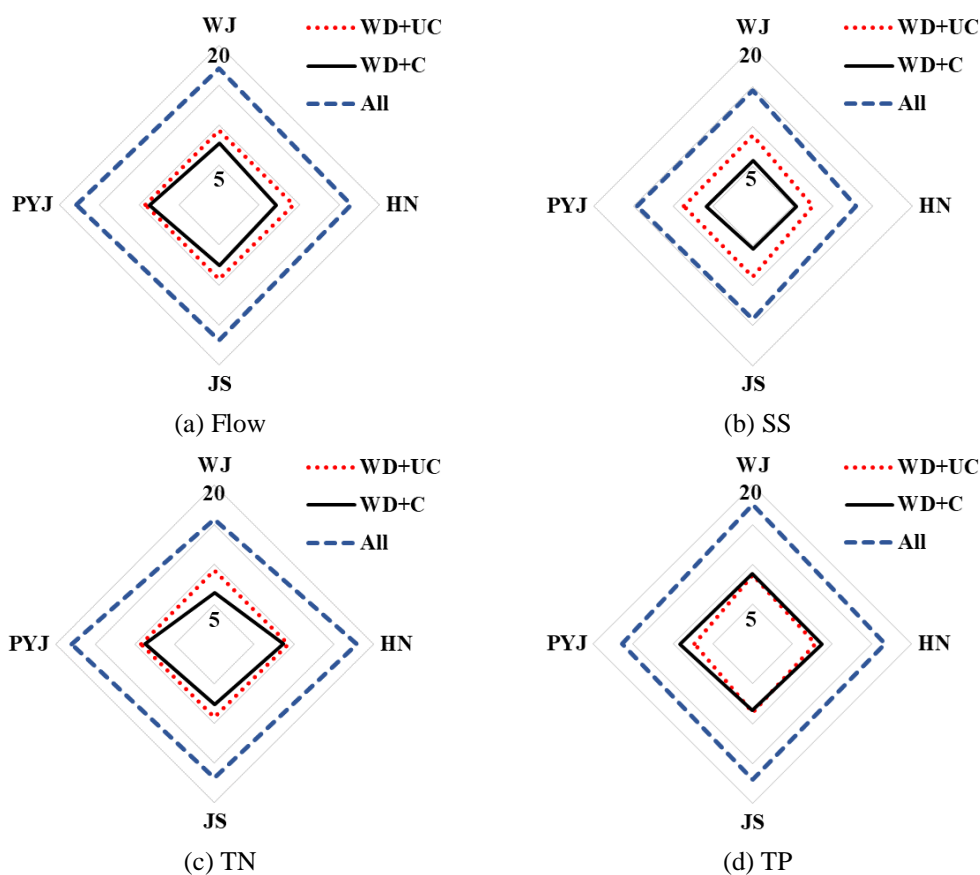
**Figure 5.** Comparison between observed and ML-predicted FDCs at the outlet of the PYJ watershed.

### 295 3.4 Information quantity and quality

The amount and quality of information contained in the training data sets for the ML training were quantified using the ME and TE concepts, respectively. Then, IUE was calculated to understand how efficiently information quantity and quality can improve the prediction accuracy of hydrological ML modeling. ME of the training data sets generally increased as additional data (i.e., the outputs of the uncalibrated and calibrated mechanistic or SWAT modeling) were added to the weather data (i.e., WDO case; Fig. 6). In the case of predicting the flow of the PYJ watershed, for example, ME increased by



8.7 to 17.9 bits when the uncalibrated and/or calibrated SWAT modeling outputs were added to the training data respectively, compared to the WDO case (Fig. 6[a]). The All cases increased ME more substantially than the WD+UC and WD+C cases, regardless of the watersheds and variables. The WD+C cases did not always increase ME more (or efficiently) than the WD+UC cases, and the ME increases were negligible even when they did occur (e.g., in the case of predicting TP loads); this confirms that ME does not consider the association between two variables (i.e., watershed responses observed and simulated using the calibrated mechanistic models) in the training data sets, which is one of the features that ME has. Thus, ME does not change depending on the types of ML models as ME only counts information in the training data.

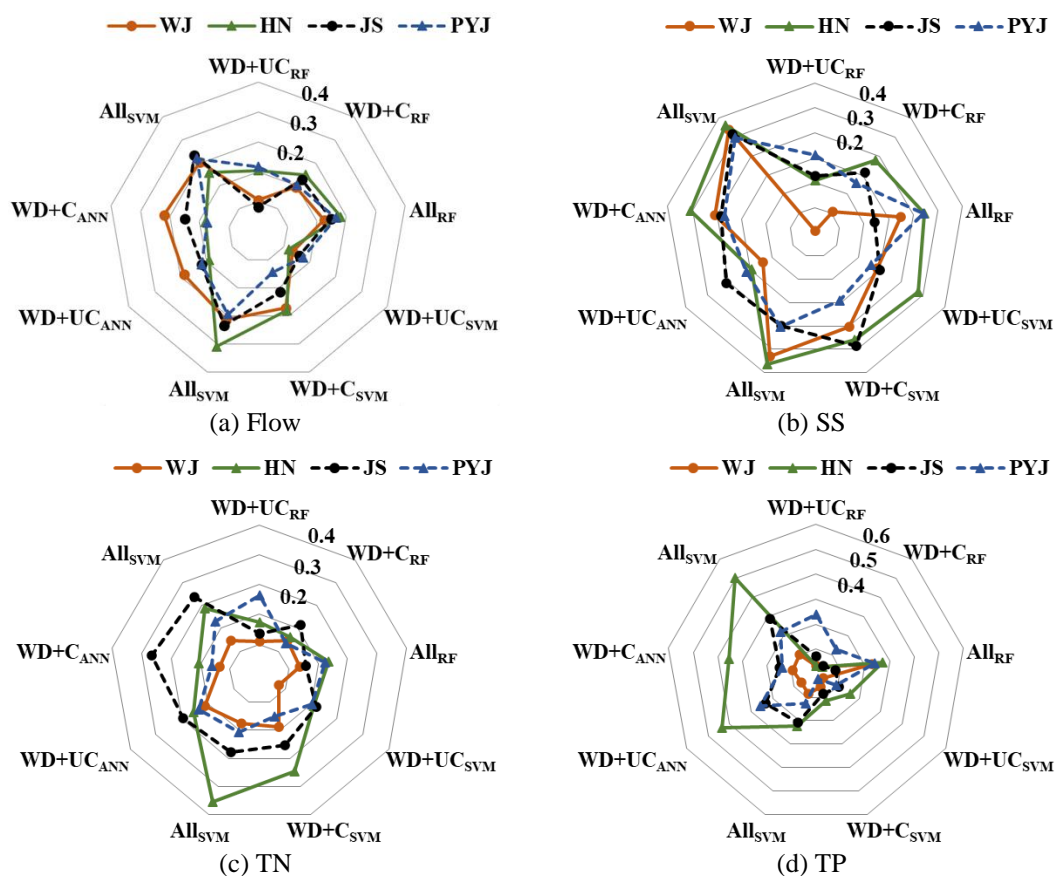


**Figure 6.** Increases in ME due to the addition of additional training data sets. The WDO training data set serves as the baseline for this comparison.

TE did not always increase with additional training data. For example, in the case of the RF modeling trained with the WDO data set for the WJ watershed, the TE of SS loads decreased from 0.385 to 0.190 and 0.294 when adding uncalibrated



and calibrated mechanistic modeling outputs, respectively (Fig. 7); this indicates that a loss of information was commonly  
 315 found in the target and training data sets when adding additional data, such as uncalibrated and calibrated modeling outputs,  
 to the training data set. The amount of precipitation information (0.174 bits) was transferred to the SS prediction for the WJ  
 watershed in the case of WDO. However, when adding the uncalibrated mechanistic (i.e., SWAT) modeling output to the  
 training data set, the amount of transferred precipitation information decreased to 0.066 bits, whereas only 0.044 bits were  
 320 transferred from the uncalibrated SWAT modeling output. Here, the information loss of 0.064 bits can be calculated by  
 subtracting 0.110 bits (amount of information on precipitation and uncalibrated mechanistic modeling output when applying  
 WD+UC training data set) from 0.174 bits. TE considers the amount of information contained in the training data sets and  
 then transfers it into the predictions made using the trained ML models. TE considers the amount of information commonly  
 found in input and output variables and the direction of information flow from variable x to another variable y (Schreiber,  
 2000). Thus, TE depends on the types of training data sets, prediction variables, and ML models (Figs. 6 vs. 7).  
 325



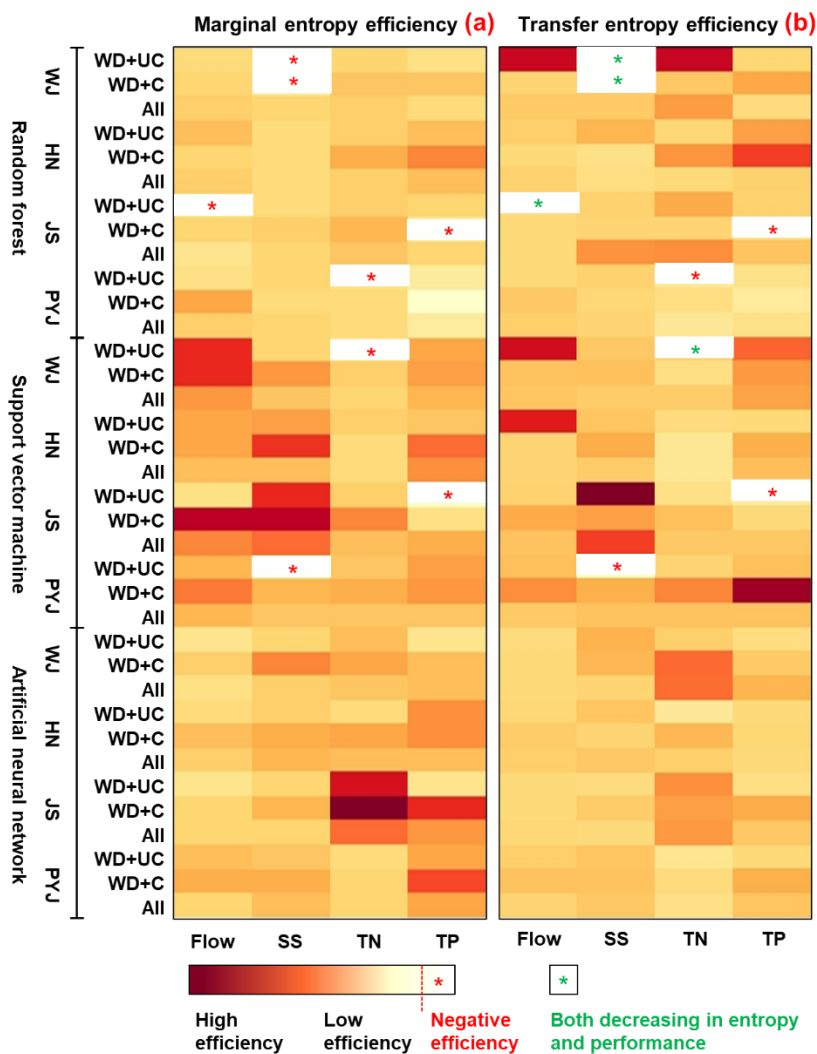
**Figure 7.** Increases in TE due to the addition of training data sets. The WDO training data set serves as the baseline for this comparison.





### 3.5 Information use efficiency

IUE represents the relative improvement of prediction accuracy compared to the baseline per unit change of information quantity (Eqs. 7 and 8). IUE was calculated by dividing the increases in KGEs (the WDO training data set serves as the baseline) by the differences between the amount of information quantified using ME (IUE-ME) or TE (IUE-TE) contained in the training data sets (Fig. 8).



335 **Figure 8.** Comparison of information use efficiency calculated from the entropy (ME and TE) and accuracy (KGE) statistics provided by using the different training sets. “Negative efficiency” describes the case where prediction accuracy decreased with increases in entropy, which is presented with a red symbol. In addition, decreases in both entropy and prediction accuracy are presented with a green symbol.



The case of WD+C provided a relatively higher IUE-ME compared to the other training data cases (Table S3). This means  
340 that the prediction accuracy of ML modeling can be most efficiently improved when the outputs of the calibrated  
mechanistic modeling are added to the training data sets (i.e., WD+C). Interestingly, WD+C may be more efficient than the  
All case, which added the uncalibrated theory-driven modeling outputs to WD+C. This finding implies that information  
quality can more efficiently improve the prediction accuracy of hydrological ML modeling than information quantity.  
However, it is worth noting that the All case still provided the best prediction accuracy (or the highest KGE), but its  
345 efficiency in increasing KGE scores was lower than that of WD+C when considering the relative accuracy improvement to  
the amount of added information.

IUE-ME were often negative, especially in the WD+UC case, indicating that prediction accuracy decreased even when  
entropy increased (red star in Fig. 8[a]); this is because ME always increases with additional input variables, regardless of  
their quality or association with the target variables. IUE-TE also showed negative efficiency, which means the KGE  
350 decreased with increases in the TE. Model performance might not necessarily relate to TE because of complicated  
associations among weather forcings, management practices, watershed features, and responses (Konapala et al., 2020).  
KGEs also decreased when TE decreased (green star in Fig. 8[b]), which implies that TE can capture the decrease in  
information flow between independent (i.e., weather data, uncalibrated and calibrated modeling outputs) and dependent (or  
target) variables that may lead to decreased prediction accuracy (or decreased KGEs). This inverse relationship was  
355 primarily detected when adding uncalibrated mechanistic modeling outputs to the training data set, demonstrating the role of  
information quality in ML modeling training.

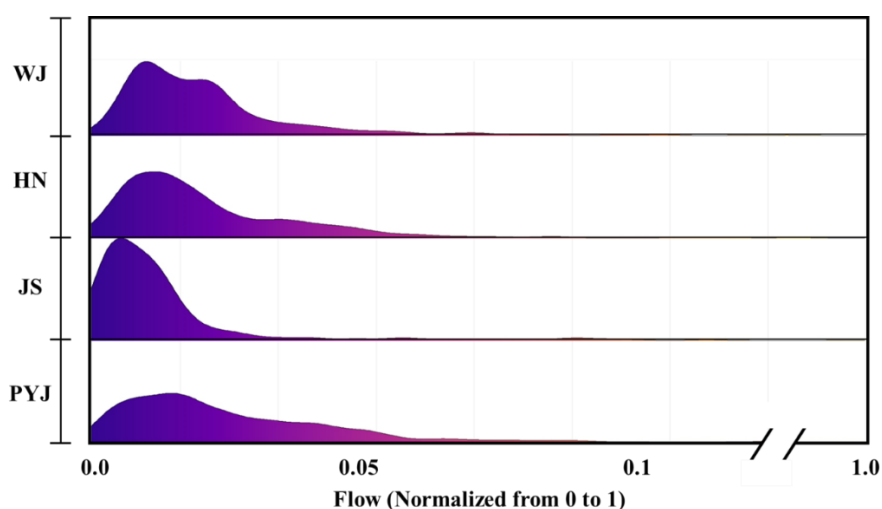
#### 4 Discussions

This study investigated how the prediction accuracy of hydrological ML modeling is associated with the quantity and  
quality of information contained in the training data. The results exhibited that prediction accuracy (KGE scores) generally  
360 increased with the amount and quality of information contained in the training data sets (all cases except the cases with stars  
in Fig. 8). Hence, access to both a large quantity and high-quality information helps increase hydrological ML modeling  
accuracy. However, the prediction accuracy of hydrological ML modeling and its association with entropy scores were found  
to be dependent on the study watersheds, target variables, and the ML models.

The accuracy of the ML modeling varied by the watersheds. Regardless of the training data sets, the ML models provided  
365 the best prediction accuracy for the PYJ watershed, which has the largest drainage area, while they did the worst for the JS  
watershed, which has the smallest drainage area; this implies the potential impact of watershed features and responses (flow,  
SS, TN, and TP) on ML prediction accuracy. For example, entropy in the watershed responses of the PYJ watershed was  
consistently higher than that of the JS watershed (Table 3). In the case of WDO, the amount of information contained in the  
independent variables (i.e., only weather records observed at a single station) of the training data set should be the same for  
370 the PYJ and JS watersheds. However, their responses (dependent or target variables) differ and thus have different entropy



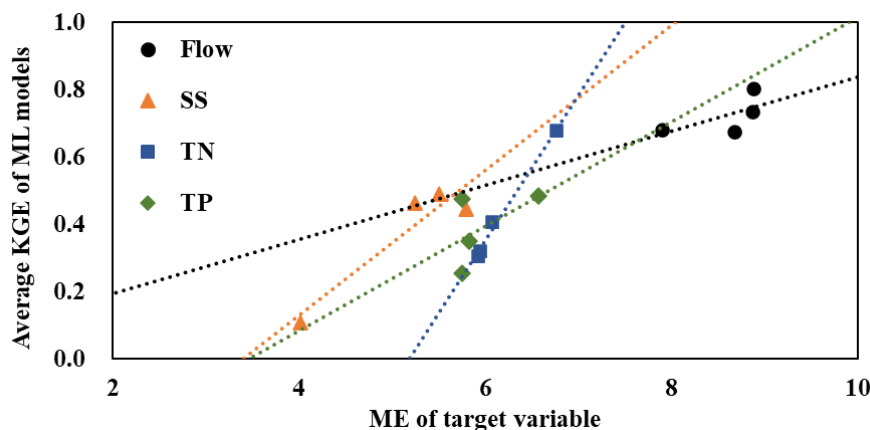
(or information) scores. The responses of the PYJ watershed are spread out over wide value ranges, which means relatively high entropies compared to those of the other watersheds, especially the JS watershed (Figs. 9 and S8–S10). The flow observed at the outlet of the JS watershed are relatively highly biased toward low flow ranges. ML model prediction accuracy was found to be associated with the entropy in the watershed responses (Fig. 10). The results indicated that the KGE (i.e., prediction accuracy) scores of the ML models generally increased with increases in the amount of information contained in the target variables (Fig. 10).



**Figure 9.** Density (or frequency) distributions of observed flow data (i.e., target variable) during the training period. The flow was normalized from 0 to 1 for each watershed.

**Table 3.** ME quantified for target variables observed at the watershed outlets in the training period.

Watershed	Flow	SS	TN	TP
WJ	8.680	5.507	6.080	5.757
HN	8.868	5.252	5.946	5.828
JS	7.896	4.014	5.924	5.760
PYJ	8.884	5.801	6.770	6.578
Average	8.582	5.144	6.180	5.981



385 **Figure 10.** Linear relationship between the average KGE scores of three ML models trained using the four training data sets  
and the ME of the target variables.

The prediction accuracy of the ML models also varied according to the variables of interest (flow, SS, TN, and TP). The RF, SVM, and ANN ML models were best at predicting the flow of the study watersheds compared to the other variables  
390 (Fig. 4). For example, the ML models provided KGEs of 0.557 (WDO) to 0.854 (All) when predicting flow versus KGEs of 0.093 (WDO) to 0.607 (All) for SS loads. This variance is presumably because of the previously described differences in the amount of information contained in the watershed responses or target variables, which is also why prediction accuracy varied by watershed. For example, the flow hydrographs commonly have relatively higher entropies (8.582 on average) than the other variables' hydrographs (5.144 for SS, 6.180 for TN, and 5.981 for TP on average) for all study watersheds (Table 3). In  
395 the frequency domain, normalized SS load data have the most biased distributions toward low values (small SS loads) and the highest frequencies among the watershed responses or target variables, leading to relatively low entropy in the SS data (Figs. 9 and S8–S10). The SS, TN, and TP concentrations observed at the watershed outlets have relatively small variations compared to the flow (Table S2), which might be attributed to the fact that water quality variables were much less frequently measured (or sampled) than flow (Table S2, the number of observations); thus, potentially large concentration variations  
400 might not be apparent in the observations. These comparison results imply that the frequency of water quality sampling can affect the amount of information in training data and the accuracy of hydrological ML model prediction.

None of the ML models consistently provided more accurate predictions than the others (Fig. 4). This finding aligns with other studies that identified no ML model that is universally applicable to all data sets or problems (Alzubi et al., 2018). Some study has demonstrated that the RF model is more accurate compared to the SVM and ANN models (Al-Mukhtar,  
405 2019). Conversely, other study has determined that the SVM or ANN model outperform the RF model (Ahmad et al., 2018). In this study, the RF model provided relatively better accuracy than other ML models when predicting the streamflow of the PYJ watershed using all training data sets (the All case). The ANN model was the only one that could provide acceptable accuracy (KGE of 0.84, which is greater than the threshold of 0.17 for SS) in the prediction of SS loads in the JS watershed



with the WD+C training data. The SVM model provided a relatively greater KGE score than the other models when  
410 predicting the SS loads of the HN watersheds using the WD+UC training data.

The ANN and SVM models could improve their predictions more efficiently in terms of the amount of information (ME)  
added to the training data compared to the RF model (Fig. 8[a]). The RF model uses a random sampling method to select the  
feature subspace for each node in growing the trees (Breiman, 2001), which is a model parameter called the “number of  
variables.” Previous studies (Wang and Xia., 2016; Ye et al., 2013) have argued that when applying the random sampling  
415 method to a high-dimension data set, model may select many subspaces that do not include informative features and will  
increase error bounds for the RF model. This study agrees with previous studies: RF performed relatively poorly when  
dimension of a training data set was higher (i.e., the large number of independent variables) than SVM and ANN. A  
traditional ANN model with one or two hidden layers is known to suffer performance degradation due to its rapid growth in  
the number of connection weights (Krenker et al., 2011). However, one study demonstrated that a deep neural network that  
420 employs numerous hidden layers, such as the one used in this study, could yield promising performance with high-  
dimensional training data (Liu et al., 2017).

Negative IUE-TE values were found when watershed responses were predicted using the RF and SVM models (red star in  
Fig. 8[b]), especially in the WD+UC case. The RF and SVM models occasionally failed to utilize additional information  
contained in the training data, presumably due to the “curse of dimensionality” (Bellman, 1961) and/or the complicated  
425 nonlinear relationship between independent and dependent (or target variables) data. Both RF and SVM models use  
“piecewise” linear decision boundaries or hyperplanes to partition the input space, but the decision boundaries can be  
nonlinear overall at the end of the decisions or partitioning. To handle nonlinear cases, SVM models employ a kernel  
function to transform the nonlinear decision space into a linear one, and RF models use nonlinear decision boundaries  
(Kirasich et al., 2018). However, studies have supported that nonlinear decision boundaries might not always be able to help  
430 solve the high-dimensionality issue, mainly because random sampling could sample less informative features (or noises)  
when growing trees (Wang and Xia., 2016; Ye et al., 2013). On the other hand, the kernel function of the SVM model has  
been found to handle high-dimensional non-linear data well (Huang et al., 2018), which contrasts with our finding. In this  
study, we used a radial basis kernel function and optimized the ML models using a Bayesian optimization method, which  
may affect predictive accuracy with feature selection (Shawe-Taylor and Sun, 2011). The ANN model did not provide any  
435 negative IUE scores in this study, and its performance is less affected by the relatively low-quality information included in  
the training data set (the WD+UC case in Table S3). These findings suggest that the ANN model can more efficiently utilize  
quality information than the other two models.

## 5 Conclusions

From the modeling experiment, this study revealed that the accuracy of hydrological ML prediction is closely associated  
440 with the quantity and quality of data used to train the models. Prediction accuracy was maximized when all available data



(i.e., the All case) were employed for training, and it was most efficiently improved in terms of information use when relatively high-quality data (i.e., the WD+C case) were added to the training data set. Information use efficiency was affected by the amount of information contained in the dependent (or target) variables of the training data, which varied by watershed and variable type. ML model performance was case-dependent, and the ANN model could more efficiently utilize the quality information contained in the training data set than the SVM and RF models. Relatively low-quality information (i.e., the WD+UC) case sometimes did not improve prediction accuracy, demonstrating the significance of the quality as well as the quantity of training data. These findings are expected to elucidate the relationship between information and ML modeling accuracy, highlighting the importance of data quality and information in ML model training.

**Author contributions.** **MJ:** Conceptualization, Software, Validation, Formal analysis, Writing - Original Draft; **YH:** Conceptualization, Methodology, Supervision, Writing - Review & Editing; **SB:** Validation, Formal analysis, Data Curation; **KY:** Conceptualization, Supervision, Writing - Review & Editing.

**Competing interests.** The contact author has declared that none of the authors has any competing interests

455

**Acknowledgements.** This research was supported by a project titled “A Long-term Monitoring for the Nonpoint Sources Discharge” (Yeongsan and Seomjin River Water Management Committee).

## References

- Adeola-Fashae, O., Abiola-Ayorinde, H., Oludapo-Olusola, A., Oluseyi-Obateru, R., 2019. Landuse and surface water quality in an emerging urban city. *Appl. Water Sci.* 9(25), Doi: 10.1007/s13201-019-0903-2.
- Ahmad, I., Basher, M., Iqbal, M.J., Rahim, A., 2018. Performance Comparison of Support Vector Machine, Random Forest, and Extreme Learning Machine for Intrusion Detection. *IEEE Access* 6, 33789-33795.
- Ahmed, S., Khalid, M., Akram, U., 2017. A Method of Short-Term Wind Speed Time Series Forecasting Using Support Vector Machine Regression Model. “2017 6th International Conference on Clean Electrical Power (ICCEP), 190-195. Doi: 10.1109/ICCEP.2017.8004814.
- Aktan, S., 2011. Application of machine learning algorithm for business failure prediction. *Invest. Manage. And Financial Inno.* 8(2), 52-65.
- Al-Mukhtar, M., 2019. Random forest, support vector machine, and neural networks to modelling suspended sediment in Tigris River-Baghdad. *Environ. Monit. Assess.* 191, 673.
- Alzubi, J., Nayyar, A., Kumar, A., 2018. Machine Learning from Theory to Algorithms: An Overview. *Journal of Physics: Conference Series* 1142, 012012.

470



- Andersson, J.C.M., Arheimer, B., Traoré, F., Gustafsson, D., Ali, A., 2017. Process refinements improve a hydrological model concept applied to the Niger River basin. *Hydrol. Process.* 31, 4540-54.
- Bellman, R., 1961. *Adaptive Control Processes: A Guided Tour*. Princeton University Press.
- 475 Bennett, A., Nijssen, B., Ou, G., Clark, M., Nearing, G., 2019. Quantifying Process Connectivity with Transfer Entropy in Hydrologic Models. *Water Resour. Res.* 55, 4613-4629.
- Breiman, L., 2001. Random Forests. *Mach. Learn.* 45, 5-32.
- Breiman, L., Friedman, J., Olshen, R.A., Stone, C.J., 1984. *Classification and Regression Trees*. CRC press, Wadsworth.
- Cabaneros, S.M.S., Calautit, J.K.S., Hughes, B.R., 2017. Hybrid Artificial Neural Network Models for Effective Prediction  
480 and Mitigation of Urban Roadside NO<sub>2</sub> Pollution. *Energy Procedia* 142, 3524-3530.
- Chaudhary, S., Chua, L.H.C., Kansal, A., 2022. Event mean concentration and first flush from residential catchments in different climate zones. *Water Res.* 219, 118594.
- Chen, Z., Zhu, Z., Jiang, H., Sun, S., 2020. Estimating daily reference evapotranspiration based on limited meteorological data using deep learning and classical machine learning methods. *J. Hydrol.* 591, 125286.
- 485 Choi, J.Y., Engel, B.A., Chung, H.W., 2002. Daily streamflow modelling and assessment based on the curve-number technique. *Hydrol. Process.* 16, 3131-3150.
- Cover, T.M., Thomas, J.A., 2006. *Elements of Information Theory*. John Wiley & Sons, Inc., Hoboken.
- Díaz-Uriarte, R., Alvarez de Andrés, S., 2006. Gene selection and classification of microarray data using random forest. *BMC bioinformatics* 7, 1-13.
- 490 Gupta, H.V., Kling, H., Yilmaz, K.K., Martinez, G.F., 2009. Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling. *J. Hydrol.* 377, 80-91.
- Hasanipanah, M., Faradonbeh, R.S., Amnieh, H.B., Armaghani, D.J., Monjezi, M., 2017. Forecasting blast-induced ground vibration developing a CART model. *Eng. Comput.* 33, 307-316.
- Her, Y., Jeong, J., 2018. SWAT+ versus SWAT2012: Comparison of sub-daily urban runoff simulations. *Trans. ASABE*  
495 61(4), 1287-1295.
- Her, Y., Jeong, J., Arnold, J., Gosselink, L., Glick, R., Jaber, F., 2017. A new framework for modeling decentralized low impact developments using Soil and Water Assessment Tool. *Environ. Model. Softw.* 96, 305-322.
- Huang, S., Cai, N., Pacheco, P.P., Narrandes, S., Wang, Y., Xu, W., 2018. Applications of support vector machine (SVM) learning in cancer genomics. *Cancer Genom. Proteom.* 15, 41-51.
- 500 Ioffe, S., Szegedy, C., 2015. Batch Normalization: Acceleration Deep Network Training by Reducing Internal Covariate Shift. arXiv preprint.
- Jang, W.S., Engel, B., Yeum, C.M., 2020. Integrated environmental modeling for efficient aquifer vulnerability assessment using machine learning. *Environ. Model. Softw.* 124, 104602.
- Jha, D., Ward, L., Paul, A., Liao, W.-k., Choudhary, A., Wolverton, C., Agrawal, A., 2018. ElemNet: Deep Learning the  
505 Chemistry of Materials from Only Elemental Composition. *Sci. Rep.* 8, 17593.



- Khashei, M., Bijari, M., 2010. An artificial neural network (p,d,q) model for timeseries forecasting. *Expert Syst Appl.* 37, 479-489.
- Kim, N.W., Lee, J., 2008. Temporally weighted average curve number method for daily runoff simulation. *Hydrol. Process.* 22, 4936-4948.
- 510 Kirasich, K., Smith, T., Sadler, B., 2018. Random forest vs logistic regression: binary classification for heterogeneous datasets. *SMU Data Science Review* 1, 9.
- Knoben, W.J.M., Freer, J.E., Woods, R.A., 2019. Technical note: Inherent benchmark or not? Comparing Nash–Sutcliffe and Kling–Gupta efficiency scores, *Hydrol. Earth Syst. Sci.* 23, 4323-31.
- Konapala, G., Kao, S.C., Addor, N., 2020. Exploring Hydrologic Model Process Connectivity at the Continental Scale  
515 Through an Information Theory Approach. *Water Resour. Res.* 56, e2020WR027340.
- Krenker, A., Bester, J., Kos, A., 2011. Introduction to the artificial neural networks. In K. Suzuki (Ed.), *Artificial neural networks-methodological advances and biomedical applications*. Rijeka, Croatia: IntechOpen.
- Li, S., Liu, Y., Her, Y., Chen, J., Guo, T., Shao, G., 2021a. Improvement of simulating sub-daily hydrological impacts of rainwater harvesting for landscape irrigation with rain barrels/cisterns in the SWAT model. *Sci. Total Environ.* 798,  
520 149336.
- Li, T.Y., Xiang, H., Yang, Y., Wang, J., Yildiz, G., 2021b. Prediction of char production from slow pyrolysis of lignocellulosic biomass using multiple nonlinear regression and artificial neural network. *J. anal. appl. pyrolysis.* 159, 105286.
- Liu, B., Wei, Y., Zhang, Y., Yang, Q., 2017. Deep Neural Networks for High Dimension, Low Sample Size Data. In  
525 *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI)*, 2287-2293.
- Liu, Z., Yang, J., Yang, Z., Zou, J., 2012. Effects of rainfall and fertilizer types on nitrogen and phosphorus concentrations in surface runoff from subtropical tea fields in Zhejiang, China. *Nutr. Cycl. Agroecosyst.* 93, 297-307. Doi: 10.1007/s10705-012-9517-x.
- Loague, K., Heppner, C.S., Ebel, B.A., VanderKwaak, J.E., 2010. The quixotic search for a comprehensive understanding of  
530 hydrologic response at the surface: Horton, Dunne, Dunton, and the role of concept-development simulation. *Hydrol. Process.* 24, 2499-2505.
- Mendie, U.E., 2005. *The theory and practice of clean water production for domestic and industrial use: Purified and package water*. Lacto-Medal Ltd, Lagos.
- Moriasi, D.N., Arnold, J.G., Van Liew, M.W., Bingner, R.L., Harmel, R.D., Veith, T.L., 2007. Model evaluation guidelines  
535 for systematic quantification of accuracy in watershed simulations. *Trans. ASAE.* 50 (3), 885–900. Doi: 10.13031/2013.23153.
- Nash, J.E., Sutcliffe, J.V., 1970. River flow forecasting through conceptual models part I—A discussion of principles. *J. Hydrol.* 10, 282-290.





- Nearing, G.S., Ruddell, B.L., Bennett, A.R., Prieto, C., Gupta, H.V., 2020. Does Information Theory Provide a New  
540 Paradigm for Earth Science? Hypothesis Testing. *Water Resources Research* 56, e2019WR024918.
- Nietsch, S.L., Arnold, J.G., Kiniry, J.R., Srinivasan, R., Williams, J.R., 2002. SWAT: Soil and water assessment tool user's  
manual. Texas Water Resources Institute, USDA Agricultural Research Service, College Station, TX.
- Noori, N., Kalin, L., Isik, S., 2020. Water quality prediction using SWAT-ANN coupled approach. *J. Hydrol.* 590, 125220.
- Panidhappu, A., Li, Z., Aliashrafi, A., Peleato, N.M., 2020. Integration of weather conditions for predicting microbial water  
545 quality using Bayesian Belief Networks. *Water Res.* 170, 115349.
- Pechlivanidis, I.G., Gupta, H., Bosshard, T., 2018. An Information Theory Approach to Identifying a Representative Subset  
of Hydro-Climatic Simulations for Impact Modeling Studies. *Water Resources Research* 54, 5422-5435.
- Pullanikkatil, D., Palamuleni, L.G., Ruhiga, T.M., 2015. Impact of land use on water quality in the Likangala catchment,  
southern Malawi. *Afr. J. Aquat. Sci.* 40(3), 277-286. Doi: 10.2989/16085914.2015.1077777.
- 550 Qiu, L., Zheng, F., Yin, R.-s., 2012. SWAT-based runoff and sediment simulation in a small watershed, the loessial hilly-  
gullied region of China: Capabilities and challenges. *International J. Sediment. Res.* 27, 226–234.
- Raju, V.N.G., Lakshmi, K.P., Jain, V.M., Kalidindi, A., Padma, V., 2020. Study the Influence of  
Normalization/Transformation process on the Accuracy of Supervised Classification. 2020 Third International Conference  
on Smart Systems and Inventive Technology (ICSSIT), pp. 729-735.
- 555 Rural Development Administration (RDA), 2014. Agricultural work schedule – Machine transplanting cultivation.  
<http://www.nongsaro.go.kr>.
- Sao, D., Kato, T., Tu, L.H., Thouk, P., Fitriyah, A., Oeurng, C., 2020. Evaluation of Different Objective Functions Used in  
the SUFI-2 Calibration Process of SWAT-CUP on Water Balance Analysis: A Case Study of the Pursat River Basin,  
Cambodia. *Water* 12, 2901.
- 560 Schaeffli, B., Gupta, H.V., 2007. Do Nash values have value? *Hydrol. Process.* 21, 2075-2080.
- Schreiber, T., 2000. Measuring information transfer. *Phys. Rev. Lett.* 85, 461.
- Senent-Aparicio, J., Jimeno-Sáez, P., Bueno-Crespo, A., Pérez-Sánchez, J., Pulido-Velázquez, D., 2019. Coupling machine-  
learning techniques with SWAT model for instantaneous peak flow prediction. *Biosyst. Eng.* 177, 67-77.
- Shannon, C.E., 1948a. A mathematical theory of communication. *The Bell system technical journal* 27, 379-423.
- 565 Shannon, C.E., 1948b. A mathematical theory of communication. *The Bell System Technical Journal* 27, 623-656.
- Shawe-Taylor, J., Sun, S., 2011. A review of optimization methodologies in support vector machines. *Neurocomputing*  
74(17), 3609-3618. doi: 10.1016/j.neucom.2011.06.026.
- Siddique, M., Tokhi, M.O., 2001. Training neural networks: backpropagation vs. genetic algorithms. IJCNN'01. International  
Joint Conference on Neural Networks. Proceedings (Cat. No. 01CH37222). IEEE, pp. 2673-2678.
- 570 Silva, V.d.P.R.d., Belo Filho, A.F., Singh, V.P., Almeida, R.S.R., Silva, B.B.d., de Sousa, I. F., Holanda, R.M.d., 2017.  
Entropy theory for analysing water resources in northeastern region of Brazil. *Hydrol. Sci. J.* 62(7), 1029-1038. Doi:  
10.1080/02626667.2015.1099789.



- Srinivasan, R., Zhang, X., Arnold, J., 2010. SWAT ungauged: hydrological budget and crop yield predictions in the Upper Mississippi River Basin. *Trans. ASABE* 53, 1533-1546.
- 575 Srivastava, A., Kumari, N., Maza, M., 2020. Hydrological Response to Agricultural Land Use Heterogeneity Using Variable Infiltration Capacity Model. *Water Resour. Manag.* 34, 3779-3794.
- Sun, W., Lv, Y., Li, G., Chen, Y., 2020. Modeling River Ice Breakup Dates by k-Nearest Neighbor Ensemble. *Water* 12(1), 222.
- Tang, X., Zhang, J., Wang, G., Jin, J., Liu, C., Liu, Y., He, R., Bao, Z., 2021. Uncertainty Analysis of SWAT Modeling in the Lancang River Basin Using Four Different Algorithms. *Water* 13, 341.
- 580 Tao, J., Chen, W., Wang, B., Jiezheng, X., Nianzhi, J., Luo, T., 2008. Real-Time Red Tide Algae Classification Using Naive Bayes Classifier and SVM. 2008 2nd International Conference on Bioinformatics and Biomedical Engineering, 2888-2891.
- Tobin, K.J., Bennett, M.E., 2017. Constraining SWAT Calibration with Remotely Sensed Evapotranspiration Data. *J. Am. Water Resour. Assoc.* 53(3), 594-604.
- 585 Tosun, E., Aydin, K., Bilgili, M., 2016. Comparison of linear regression and artificial neural network model of a diesel engine fueled with biodiesel-alcohol mixture. *Alex. Eng. J.* 55, 3081-3089. Doi: 10.1016/j.aej.2016.08.011.
- Vapnik, V., 1995. *The nature of statistical learning theory*. Berlin: Springer.
- Vapnik, V., 1998. *Statistical learning theory*. New York: John Wiley & Sons.
- Wang, Y., Xia, S.T., 2016. A novel feature subspace selection method in random forests for high dimensional data. 2016 International Joint Conference on Neural Networks (IJCNN), 4383-4389.
- 590 Xu, T., Liang, F., 2021. Machine learning for hydrologic sciences: An introductory overview. *Wiley Interdisciplinary Reviews. Water* 8, e1533.
- Ye, Y., Wu, Q., Zhixue Huang, J., Ng, M.K., Li, X., 2013. Stratified sampling for feature subspace selection in random forests for high dimensional data. *Pattern Recognit.* 46, 769-787.
- 595 Yilmazkaya, E., Dagdelenler, G., Ozcelik, Y., Sonmez, H., 2018. Prediction of mono-wire cutting machine performance parameters using artificial neural network and regression models. *Eng. Geol.* 239, 96-108. Doi: 10.1016/j.enggeo.2018.03.009.
- Zeiger, S., Hubbart, J.A., 2016. Quantifying suspended sediment flux in a mixed-land-use urbanizing watershed using a nested-scale study design. *Sci. Total Environ.* 542, 315-323.