# A STUDY TO INVESTIGATE THE EFFECT OF DIFFERENT TIME-SERIES SCALES TOWARDS FLOOD FORECASTING USING MACHINE LEARNING

**NAZLI MOHD KHAIRUDIN [1], NORWATI MUSTAPHA [2], TEH NORANIS MOHD ARIS [3], MASLINA ZOLKEPLI [4]**

[1,2,3,4]Universiti Putra Malaysia, Faculty of Computer Sciences and Information Technology, Malaysia

Email: [1]nazmkhair@gmail.com, [2]norwati@upm.edu.my, [3]nuranis@upm.edu.my, [4]masz@upm.edu.my

**ABSTRACT**

Machine learning has been deemed to be a powerful approach in forecasting hydrological events such as flood using time-series historical data. A flood can be forecast in a manner of lead time whereby short-term forecast is up to 2 days, the medium forecast is between 2 to 10 days, and the long-term forecast is more than 10 days and several months of forecasts will have a seasonal lead time. Even though the determination of forecast lead time is normally bound with the purpose of operation i.e., daily operations or strategical, but the determination of time-series scale pattern to be input into the forecast model still impose a challenging task as it involves availability and variability of the data. Commonly, the hydrological data has a dynamic nature with non-stationary and non-linear characteristics. Therefore, it is important to choose dominant input to provide an accurate forecast. The objective of this study is to investigate the effects of different time-series scales of rainfall data from eight rainfall stations in Kelantan River towards the accuracy of forecasting water level at Kuala Krai station. Pre-processing techniques based on Mutual Information (MI) are also introduced to cater the variability of the data in finding the most dominant features as input to the forecast model. There are four scale patterns that have been investigated which consist of 7 days, 10 days, 14 days, and monthly. The forecasting analysis of all scale patterns were run against three machine learning models which are Artificial Neural Networks (ANN), Long-Short Term Memory (LSTM), and Adaptive Neuro-Fuzzy Inferences System (ANFIS) model. The results show that monthly scale pattern achieve the best performance compared to other scale patterns and LSTM is the best model for forecasting monthly water level. It indicates that longer time-series of scaled pattern may provide better forecasting accuracy and able to capture more information of the seasonal characteristics of the rainfall. Thus, it will largely benefit the flood management in reducing the flood risk and controlling its resources.

**Keywords**: *Flood Forecasting. Machine Learning, Rainfall, Time-Series Scale, Water Level*

## 1. INTRODUCTION

Flood is one of natural disaster that frequently occurred around the world. The impact can be recognized such as physical contact with the flood water, destruction of infrastructure, communication disruption, and levelling the crops. Consequently, it is necessary for the local, regional, and even a national to have a reliable and sustainable flood risk management. To facilitate a strategic decision by authority, an optimized forecast should be developed to help in mitigating the flood risk. It will also help in minimizing the loss of the socio-economic sectors.

The classification of hydrological forecasting such as floods would be according to their lead-time order [1]. A short-term hydrological forecast can have a lead time of up to 2 days while medium hydrological forecasts can have a lead time around 2 to 10 days. Long-term hydrological forecasts can have more than 10 days while several months of forecasts will have a seasonal lead time.

Different time-scale pattern has been used by works of literature to accommodate forecast of different lead time. In short-term hydrological

forecasting, daily time-series scale is widely adopted such as in [2] where the daily streamflow data from three stations in Iran namely Siira, Bilghan, and Gachsar is employed to forecast the daily streamflow using Random Forest Regression (RFR) and Gene Expression Programming (GEP) with decomposition method. This study has indicated the high performance achieved in forecasting high and low points of daily streamflow but did poorly in forecasting the extreme high and low flow events.

In forecasting daily runoff, [3] has used daily time-series scale runoff data of Hongjiadu reservoir as the input. The forecast model is developed using Artificial Neural Network based on Quantum-behaved Particle Swarm Optimization model. It is found that the proposed model provides better accuracy than the traditional ANN model. Alternatively, the daily time-series scale of rainfall data in two stations in Turkey is used to forecast the daily rainfall up to five days [4]. The proposed hybrid wavelet-season-neuro technique has proved that it encompassed great reliability to forecast the rainfall with 1 day lead time. However, it went faltered when the lead time surpassed 2 days.

Three different time-series scale pattern were applied by [5] namely daily, mean weekly and mean monthly to forecast the monthly streamflow in Johor River Basin using Artificial Neural Network (ANN) and Extreme Learning Machine (ELM). Although it is found that the ELM has outperformed ANN in this study, but the daily time-series scale performance in both model much more recommended than the mean weekly and mean monthly scales. This advised that although the forecast is developed using the same machine learning method, the performance can be varied when fit in with different time-series scale. Daily time-series scale data seem to be a promising predictors of machine learning forecast model. Yet, it yields an uncertain or inaccurate performances in some of the machine learning model.

Other than daily and weekly scaled time-series data, the monthly scaled is also being observed in literatures especially when monthly hydrological forecasting is needed. For instance, monthly river flow data is used to forecast 1 month ahead river flow of Tigris river with various antecedents values [6]. The proposed model of Wavelet-Extreme Machine Learning (ELM) has the best performances with consecutive three-month antecedents' inputs.

In forecasting monthly stream flow of Hurman River in Turkey and Diyalah and Lesser Zab Rivers in Iraq, [7] has used monthly flow of the rivers as predictors. The study has presented a machine learning model using stepwise approach that produce good forecasting results for 1 month ahead, but the forecast performance of six months ahead has exacerbated. Similarly, in another monthly time-series scaled application found in [8], the monthly river flow of Zarrinehrud River is forecasted using monthly river flow of Safakhaneh, Santeh and Polanian stations. The result is that when the monthly data is executed towards Multilayer Perceptron (MLP) and Radial Basis Function (RBF) model, it produces better results than Support Vector Regression (SVR) model in terms of forecasting.

A longer period of time-series scaled that beyond monthly period has been beneficial in providing forecast for strategical grounds. Seasonal time-series scaled data in [9] is utilized for seasonal precipitation in Iran using a large climate signals. The proposed model was run with various input combination up to 4 months ahead. It is found that by incorporating seasonal data, the MLP has shown better accuracy than the other models. Yearly scaled time-series data was also once used in hydrological forecast. In the study by [10], they utilizes the annual runoff from Biuliuhe and Mopanshan in China to forecast the long-term runoff. Subsequently, it indicates that machine learning model such as ANN can significantly improve the accuracy of the forecast when the annual runoff data is optimized.

Fundamentally, flood forecasting using various timescales with different lead time are challenging and complex. Historical data used in many research such as rainfall, water level, precipitation, discharge, and ground water level [11] may be imperfect or in scarce as the climate condition of the nature is very dynamic, non-linear, and non-stationary. Consequently, machine learning model would have poor generalization and weak overall performances. Hence, to refine the accuracy of the machine learning model, pre-processing of the historical data has been introduced [12].

Essentially, this research study will investigate the effects of different time-series scale pattern for certain lead time in water level forecasting. It may be included as the estimator of flood occurrences. To induce a better performance of the forecasting machine learning model, new pre-processing techniques will be manoeuvred to optimize the input. The forecast model is developed using Artificial Neural Networks (ANN), Long-Short Term Memory (LSTM) and Adaptive Neuro-Fuzzy Inferences System (ANFIS). Finally, the performances of each model are put into assessment together with various

input scale.

## 2. MATERIALS AND METHOD

### 2.1 Area of Study and Datasets

In this study, cumulative rainfall data are gathered from eight rainfall stations in Kelantan River Basin. Six of the rainfall stations are located along the Lebir River known as Gunung Gagau, Kuala Koh, Kampung Aring, Kampung Lalok, Kampung Tualang, and Kuala Krai, while the others are located along the Galas River known as Dabong and Limau Kasturi. These rainfall stations are all located in the upper stream. Lebir and Galas rivers are the main tributaries of Kelantan River. The data gained from all these stations will become predictors in forecasting the water level of Kuala Krai station that located at the downstream.

The times-series data of the cumulative rainfall is collected from the period of 1/4/2011 to 30/11/2019 for all stations. These data are scaled into 7 days, 10 days, 14 days, and monthly pattern respectively. These scales are chosen based on widely usage in previous literatures. For each scaled pattern, the dataset will be divided into 75% for training and 25% for testing the model. Table 1 shows the scaled pattern with their respective lead time for water level forecasting.

*Table 1: Scaled Pattern with Lead Times*

| Scaled Pattern | Lead Time |
|---|---|
| **7 days** | 7 days |
| **10 days** | 10 days |
| **14 days** | 14 days |
| **Monthly** | Monthly |

The forecasting of water level in Kuala Krai is deemed as necessary because they are prone to flood. The forecasting can help in minimizing the impact of flood towards the cities and citizens, thus reducing the damages of the infrastructures and crops. The information of forecasted water level will be an important indication to assist authorities in triggering action in managing the flood disaster.

### 2.2 Pre-Processing of The Time-Series Data

Historical data from eight stations along the Galas River and Lebir River are selected as inputs for the water level forecasting model. Using data against different scale pattern will give the challenging task to recognize which data of these stations or combination of stations will produce

better and accurate forecast. Therefore, this study will introduce a pre-processing technique that harnessed the power of signal decomposition together with the measurement of non-linear relationship between input and output using Mutual Information. This integration technique can enhance the model performance and the best selection of input could increase the precision and accuracy of the model [13, 14].

In conjunction with the utilized historical data mentioned in this study, it will be pre-processed by using three decomposition method Empirical Model Decomposition (EMD), Ensemble Empirical Mode Decomposition (EEMD) and Discrete Wavelet Transform (DWT). These decomposition methods are adopted in reference to the widely used of optimized input data of hydrological forecasting. This is also to find which pre-processing method will yield the best performance with the time-series scaled patterns.

DWT is optimizing the input data by decomposing the time-series into shifted and scaled version of the wavelet called a mother wavelet [14]. It can analyse the variation of time-series and will produce the time and frequency information of the signal. In hydrological forecasting, decomposing time-series using DWT is very useful in improving forecast performance [15][16]. EMD introduced by [17] to decompose the input signal into Intrinsic Mode Function (IMF) and residual. It is suitable for non-stationary and non-linear time-series data such as rainfall [18]. EMD is fully self-adaptive and there is no predetermined basis function needed [19]. The use of EMD with forecast model such as Autoregressive Moving Average Model (ARMA) [20] and Support Vector Regression (SVR) [13] has proved to increase the accuracy of the streamflow forecast. EEMD is developed by [21] to overcome the disadvantages of EMD in which it tends to frequently have mode mixing problem [22]. Finite noise is added to the signal to provides a uniform references background of the time-frequency space. EEMD can be such an effective method to extract signals from non-stationary and nonlinear data that are noisy [23]. The different length time series data used in EEMD could produce various performance of the model in which the decomposition and model have to be updated whenever new information is plugged in [19].

In the purpose of getting the optimized input, every decomposition method will be applied to each dataset with four scale patterns. Then, the Spearman correlation coefficient (*p-value*) is estimated to

determine the degree of the correlation by evaluating the level of relationship between the decomposed data and the original data [24]. Formula of *p-value* is defined as the following:

$$p - value = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

where $d_i$ is the difference ranked of each observation, and $n$ is the number of observations. As the result, the decomposed dataset is produced by every decomposition method with 8 features of which are the most correlated data from 8 stations and 1 output which is Kuala Krai water level station. Hence, 12 number of the decomposed datasets shown as in Table 2 is the output when using three decomposition method for each time-series scaled pattern. Each of the dataset is identified by the dataset ID based on decomposition method and scale pattern used. For example, emd_7 represents the dataset that decomposed by EMD for the scale pattern of 7 days.

In gaining superior input to forecast model, the integration of the decompositions methods with Mutual Information (MI) has been done. MI is suitable to sync with non-linear time-series data [25] that can handle a strong non-linear relationship between input and output. The superior input is determined by measuring dependency between random variables and the information dispersion. MI can be calculated using the following formula [26] :

$$MI(A, B) = H(A) + H(B) - H(A, B)$$

where H(A) and H(B) are the entropy of A and B, while the joint entropy of H (A, B) would be:

$$H(A, B) = -\sum_{a \in A} \sum_{b \in B} p_{AB}(a, b) \log p_{AB}(a, b)$$

where $a$ and $b$ is the specific value of $A$ and $B$, respectively $p(a,b)$ is the joint probability of these values occurring together. By integrating MI, each of the decomposed dataset will be ranked according to MI score. The score will determine which station provide dominant input to bring forward into forecast model.

*Table 2: Decomposed Datasets*

| Decomposition Method | Dataset ID | Remarks |
|---|---|---|
| EMD | emd_7 | Dataset for time-series scale of 7 days |
| | emd_10 | Dataset for time-series scale of 10 days |
| | emd_14 | Dataset for time-series scale of 14 days |
| | emd_m | Dataset for time-series scale of monthly |
| EEMD | eemd_7 | Dataset for time-series scale of 7 days |
| | eemd_10 | Dataset for time-series scale of 10 days |
| | eemd_14 | Dataset for time-series scale of 14 days |
| | eemd_m | Dataset for time-series scale of monthly |
| DWT | dwt_7 | Dataset for time-series scale of 7 days |
| | dwt_10 | Dataset for time-series scale of 10 days |
| | dwt_14 | Dataset for time-series scale of 14 days |
| | dwt_m | Dataset for time-series scale of monthly |

## 3. MODEL DEVELOPMENT

Water level forecast model is developed using three machine learning models known as Artificial Neural Networks (ANN), Adaptive Neuro Fuzzy Inference System (ANFIS) and Long-Short Term Memory (LSTM). These machine learning models have proven to provide an accurate forecast to the hydrological data with non-stationary and non-linear characteristics [27][28][29]. Each of those models is fed with optimized time-series scale pattern datasets as input which is produced in pre-processing phase. The general flow of the model development is given in Figure 1.

The performance for each of the model with their respective input data will be assessed using three statistical methods often known as "goodness of fit" [30] which are:

i.     Root Mean Squared Error (RMSE)

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y})^2} \qquad (1)$$

Mean Absolute Error (MAE)

$$MAE = \frac{1}{N} \sum_{i=1}^{N} |y_i - \hat{y}| \qquad (2)$$

ii.    Nash-Sutcliffe Efficiency (NSE)

$$NSE = 1 - \frac{\sum_{t=1}^{T} (Q_m^t - Q_o^t)^2}{\sum_{t=1}^{T} (Q_o^t - \bar{Q}_o)^2} \qquad (3)$$

In equations (1) and (2), the original value in period $i$ is denoted by $y_i$, and the forecasted value at the period of $i$ is denoted by $\hat{y}_i$. The number of samples is denoted by $N$. For both equations (1) and (2), a better forecast is depicted by the smaller values

of it. In equation (3), the original value is denoted by $Q_0$ and the forecasted value denoted by $Q_m$. $Q_o^t$ is the original value at time $t$. Contrary to the equation (1) and (2), the higher value produce by equation (3) denotes a more powerful forecast model. In literatures, RMSE has been used in assessing flow forecast [27][6] and rainfall forecast [31], while MAE is used to asses runoff forecast [19]. NSE in which provide insights on the predictive skill has been used to assess water level forecast [32] and ground water forecast [33]. These three equations are widely used in measuring hydrological forecast and may be applicable to various machine learning model, thus been used in this study.

To get the best performance of the model with the most superior input data, all data in the rank are tested and produced by the Mutual Information. Then, it will go through a process of repeated retesting and the one with lowest rank for each test cycle will be eliminate in the model testing phase.
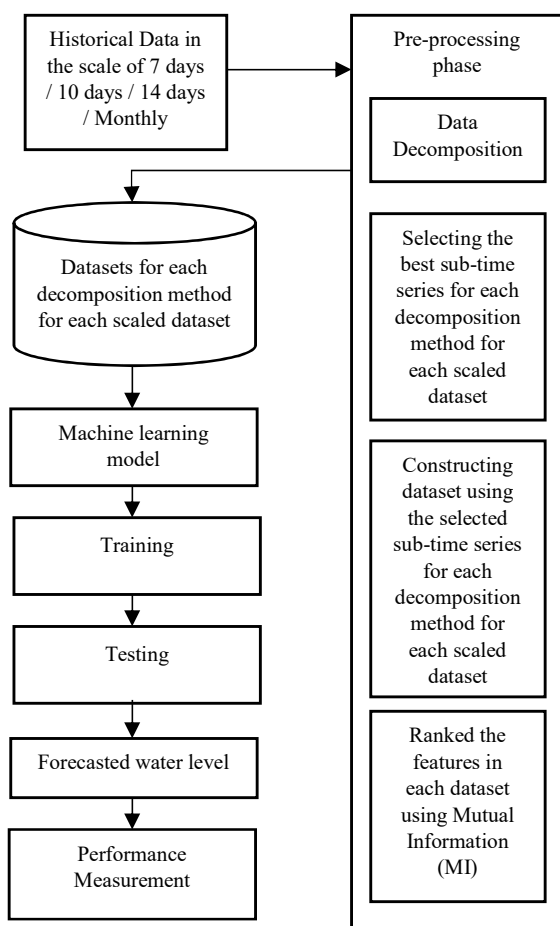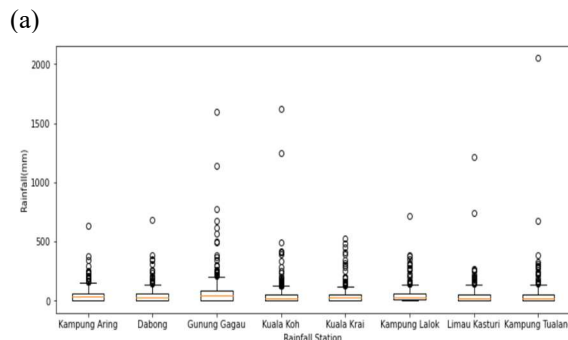
### 4. RESULTS AND DISCUSSION

Every year, Kuala Krai is one of the vulnerable locations to the changes of Northeast Monsoon. It has undergone heavy rainfall almost every single year. Most of the rainfall stations will experience this extreme rainfall in this period. In December 2014, a tremendous flood hit Kuala Krai. The impact is very destructive and caused lots of loss lives. It had estimated that RM1 billion loss was suffered in term of infrastructure and property [34].

Figure 2a-2d below display the box-whiskers plot for all original time-series scaled data of rainfall from 1/4/2011 to 30/11/2019. Box-whiskers plot present the spread and centre of the dataset by five values denote as minimum, first quartile (Q1), median, third quartile (Q3), and maximum. These figures indicates that most of the station receive almost the same amount of rainfall and the outliers which circles in shape represent the high intensity of rainfall of every station. It exhibits that monthly scaled time-series data has the largest variability and spread out by having the largest inter-quartile range and maximum-minimum differences.

Meanwhile, for Figure 3a-3d, they present the box-whiskers plot for water level station with original time-series scaled data of the same period. Department of Irrigation and Drainage Malaysia has introduced 4 classes of water level for Kuala Krai known as Normal Level (17.00m), Alert Level (20.00m), Warning Level (22.5m) and Danger Level (25.00m) [35]. The plot in these figures shows that outliers has exceeded the normal range and foresees the flood occurred in December 2014.
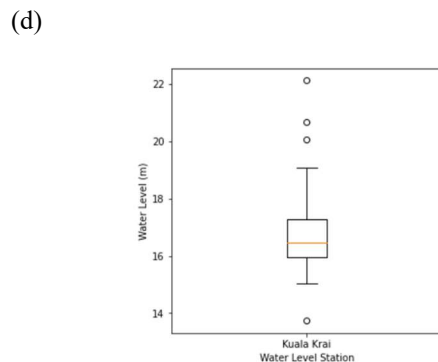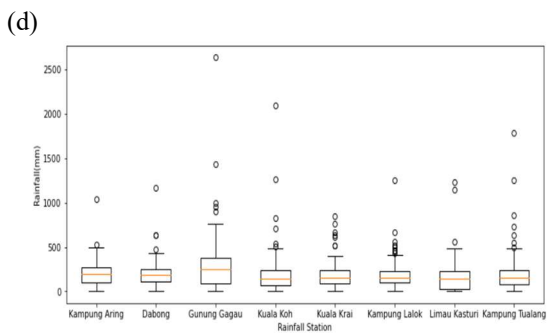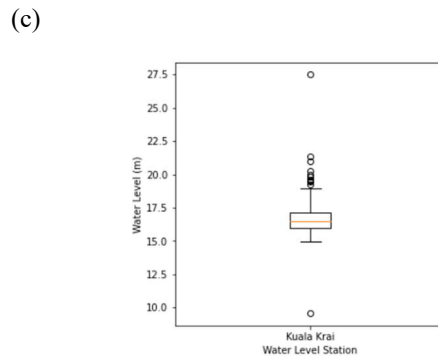
(a)



Figure 1. The Flow Diagram of Model Development

(b)



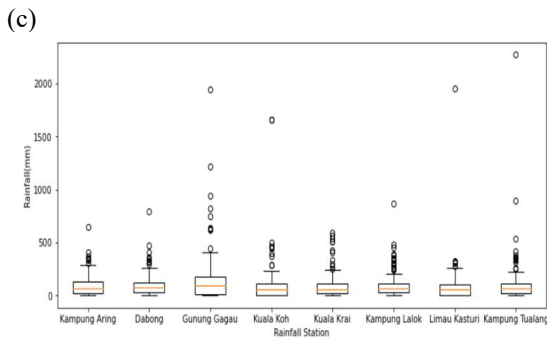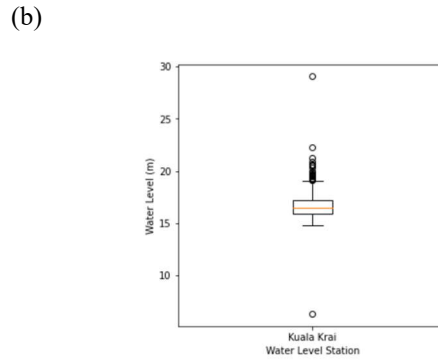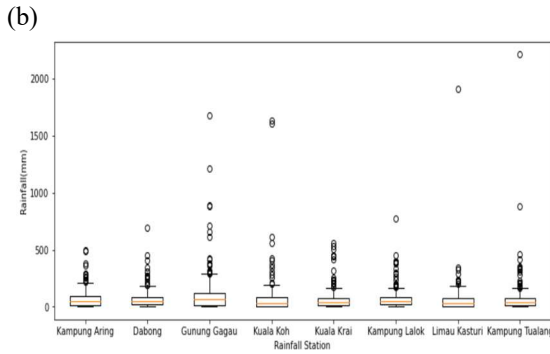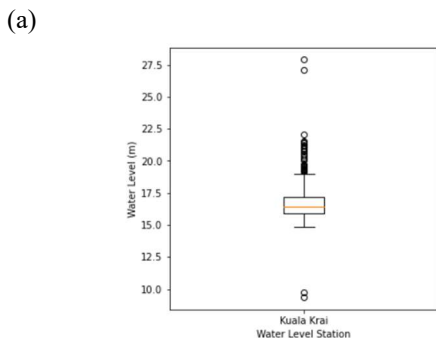(b)



(c)



(c)



(d)



(d)



*Figure 2: (a) Box-Whiskers Plot for 7 Days Rainfall Time-Series Scale Pattern, (b) Box-Whiskers Plot for 10 Days Rainfall Time-Series Scale Pattern, (c) Box-Whiskers Plot for 14 Days Rainfall Time-Series Scale Pattern, (d) Box-Whiskers Plot for Monthly Rainfall Time-Series Scale Pattern.*

*Figure 3: (a) Box-Whiskers Plot for 7 Days Water Level Time-Series Scale Pattern, (b) Box-Whiskers Plot for 10 Days Water Level Time-Series Scale Pattern, (c) Box-Whiskers Plot for 14 Days Water Level Time-Series Scale Pattern, (d) Box-Whiskers Plot for Monthly Water Level Time-Series Scale Pattern.*

(a)



As original time-series scaled data have various variability, it becomes an obstacle for the rainfall station to provide better model performance. The introduced of pre-processing method could decomposed the original data using several decomposition methods before ranking it using entropy called Mutual Information. The purpose of decomposing the original signal is to optimize the input. This method provides an insight of which station has the most information and strong non-linear relationship between input and output. Recent

studies indicate that the behaviour of the model is improving when pre-processing is applied thus increases the forecast accuracy [22][16][19][6]. In Table 3, it shows the Mutual Information rank for each of the times-series scaled pattern.

*Table 3: Mutual Information Rank*

| Scale | Decompositi-on Method | More to Less Dominant | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| **7 days** | EMD | 2 | 7 | 1 | 4 | 3 | 6 | 8 | 5 |
| | EEMD | 1 | 2 | 7 | 5 | 6 | 8 | 4 | 3 |
| | DWT | 1 | 6 | 2 | 7 | 5 | 4 | 8 | 3 |
| **10 days** | EMD | 2 | 7 | 6 | 5 | 1 | 8 | 3 | 4 |
| | EEMD | 1 | 5 | 4 | 3 | 7 | 2 | 6 | 8 |
| | DWT | 4 | 5 | 8 | 6 | 2 | 7 | 1 | 3 |
| **14 days** | EMD | 1 | 4 | 6 | 7 | 2 | 3 | 8 | 5 |
| | EEMD | 2 | 4 | 1 | 7 | 3 | 8 | 5 | 6 |
| | DWT | 1 | 4 | 2 | 8 | 6 | 7 | 3 | 5 |
| **Monthly** | EMD | 1 | 4 | 6 | 7 | 2 | 3 | 8 | 5 |
| | EEMD | 2 | 4 | 1 | 7 | 3 | 8 | 5 | 6 |
| | DWT | 1 | 4 | 2 | 8 | 6 | 7 | 3 | 5 |

*1=Gunung Gagau, 2=Kuala Koh, 3=Kampung Aring, 4=Tualang, 5=Kampung Lalok, 6=Kuala Krai, 7=Limau Kasturi, 8=Dabong*

Afterwards, we investigate the effect of water level forecasting performance of Kuala Krai station with different time-series scaled pattern in various lead time using Artificial Neural Networks (ANN), Adaptive Neuro Fuzzy Inference System (ANFIS) and Long-Short Term Memory (LSTM) models. For each time-series scaled pattern, the data are divided into training datasets of 75% and testing datasets of 25%. To evaluate the best performance of the model, all models are run repeatedly using the decomposed scaled data in the Mutual Information rank with the lowest rank is removed in each running cycle. Furthermore, the performance of every model with various time-series scaled pattern are assessed using the statistical measurement of RMSE, MAE and NSE. Table 4 presents the results of the model evaluation of the testing dataset when using these measurements.

*Table 4: Performance Measurement Results*

| Model | Time series scale | RMSE | MAE | NSE |
|---|---|---|---|---|
| LSTM | LSTM | 1.7995 | 1.2744 | -0.0335 |
| | eemd_7 | 1.7377 | 1.1712 | 0.0363 |
| | dwt_7 | 1.5554 | 1.0811 | 0.2278 |
| | emd_10 | 1.9360 | 1.2237 | -0.0719 |
| | eemd_10 | 2.0596 | 1.2404 | -0.2132 |
| | dwt_10 | 1.7180 | 1.1698 | 0.1560 |
| | emd_14 | 1.7586 | 1.1725 | 0.0324 |
| | eemd_14 | 1.7982 | 1.1604 | -0.0117 |
| | dwt_14 | 1.7642 | 1.0895 | 0.0263 |
| | emd_m | 1.0447 | 0.7519 | 0.4943 |
| | eemd_m | 1.2259 | 0.8253 | 0.3037 |
| | dwt_m | 0.9356 | 0.6742 | 0.5945 |
| ANN | emd_7 | 1.8006 | 1.3049 | -0.0348 |
| | eemd_7 | 1.6704 | 1.1342 | 0.1094 |
| | dwt_7 | 1.5833 | 1.1188 | 0.1999 |
| | emd_10 | 1.7859 | 1.1499 | 0.0879 |
| | eemd_10 | 1.9615 | 1.2190 | -0.1003 |
| | dwt_10 | 1.6745 | 1.1225 | 0.1981 |
| | emd_14 | 1.7720 | 1.2468 | 0.0176 |
| | eemd_14 | 1.7017 | 1.1902 | 0.0940 |
| | dwt_14 | 1.6751 | 1.1411 | 0.1221 |
| | emd_m | 1.2862 | 0.9524 | 0.2335 |
| | eemd_m | 1.3478 | 1.0288 | 0.1584 |
| | dwt_m | 1.1302 | 0.9181 | 0.4082 |
| ANFIS | emd_7 | 1.4898 | 0.9654 | 0.2916 |
| | eemd_7 | 1.7701 | 1.2344 | -0.0001 |
| | dwt_7 | 1.5953 | 1.0521 | 0.1878 |
| | emd_10 | 1.9942 | 1.1720 | -0.1372 |
| | eemd_10 | 1.8262 | 1.1226 | 0.0464 |
| | dwt_10 | 1.8509 | 1.0773 | 0.0203 |
| | emd_14 | 1.5742 | 1.0593 | 0.2247 |
| | eemd_14 | 1.6954 | 1.0846 | 0.1007 |
| | dwt_14 | 1.6967 | 1.0277 | 0.0993 |
| | emd_m | 1.3058 | 0.8688 | 0.1850 |
| | eemd_m | 1.2184 | 0.7769 | -0.2919 |
| | dwt_m | 1.0046 | 0.6885 | -0.1619 |

Figure 4, Figure 5, and Figure 6 display the trend for all the performance measurement of RMSE, MAE and NSE.
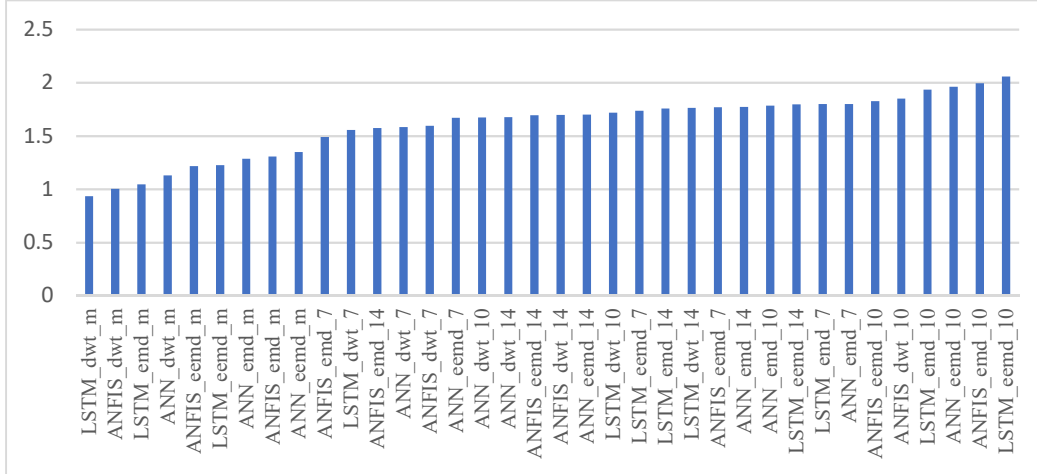
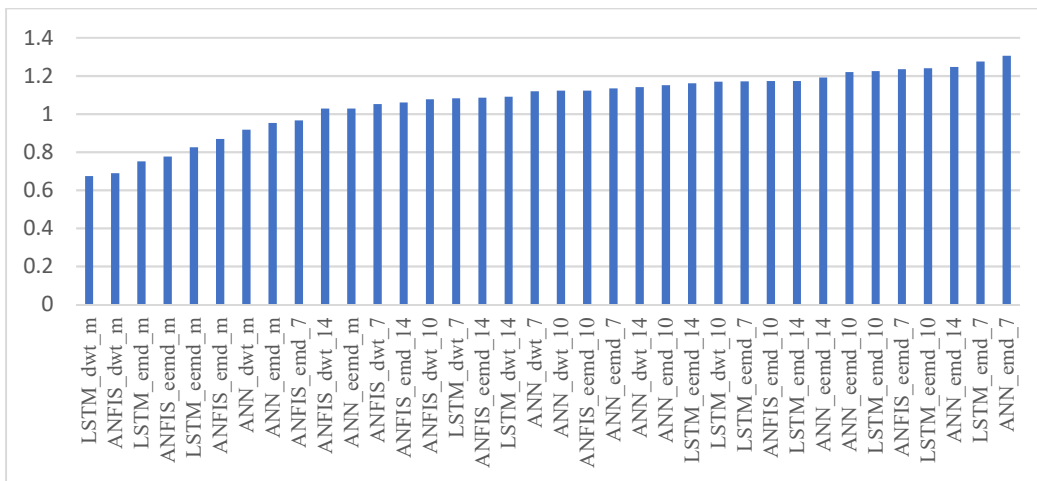*Figure 4: Performance of RMSE for All Models*



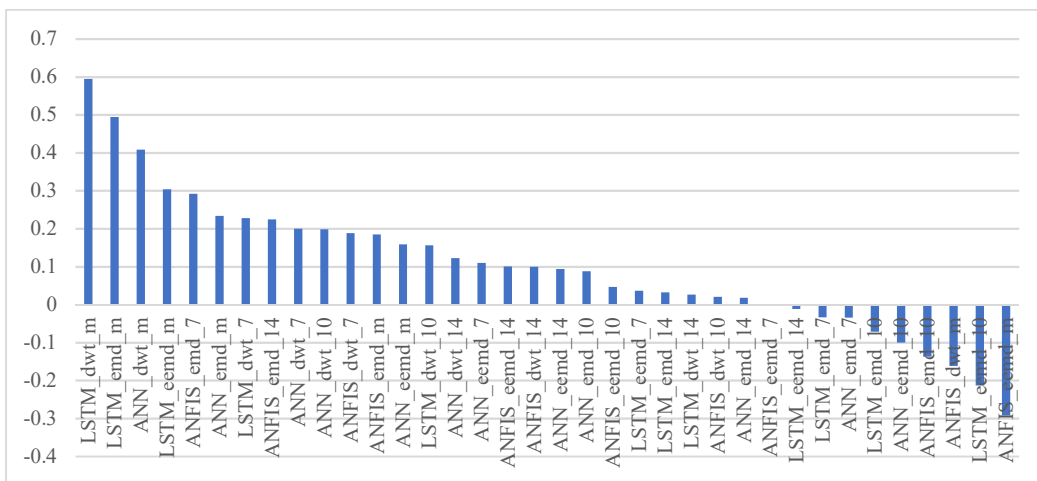*Figure 5: Performance of MAE for All Models*



*Figure 6: Performance of NSE for All Models*

The performance of RMSE pattern indicates that the lowest point achieves when water level is forecasted using monthly time-series scaled data for the monthly lead time. All model that forecasted monthly water level using the monthly time-series data has outperformed the other models. The evaluation of the MAE performance, although it is a noticeable result that most of the model in monthly time-series scaled data recorded the lowest value, but the use of time-series with 7 days scaled with EMD decomposition in ANFIS model has surpassed the performance of ANN model that using monthly time-series data scaled with EEMD decomposition. The result of NSE is much more varied at the point of highest and lowest achieved by the model that used monthly time-series scaled dataset. The highest value of NSE is recorded for the LSTM model using monthly time-series scaled data decomposed by DWT, while the lowest NSE is for ANFIS model using monthly time-series scaled data decomposed by EEMD.

On the other hand, for 7 days' time-series scaled pattern ANFIS model with EMD decomposition method, has achieved the best performance for all the measurement. For 10 days' time-series scaled pattern, best performance of RMSE and NSE is achieved by ANN model with DWT decomposition. Comparably, the lowest MAE is achieved by ANFIS model with DWT decomposition. For 14 days' time-series scaled pattern, the best performance is achieved by the ANFIS model, but the best RMSE and NSE is achieved by coupling it with EMD decomposition and lowest MAE is achieved by coupling it with DWT decomposition method. With monthly time-series scaled pattern, best performance of all the measurement achieved by the LSTM model with DWT decomposition.

Figure 7 presents the observed and forecasted water level best model performance for 7 days' time-series scaled pattern for ANFIS model with EMD decomposition. Meanwhile, the Figure 8a-8b shows the observed and forecasted water level best model performance for 10 days' time-series scaled pattern. Figure 8(a) has the values for ANN model with DWT decomposition that achieved best performance in RMSE and NSE while 8(b) has the value for ANFIS model with DWT decomposition that achieved the lowest MAE. Next, Figure 9a-9b present the observed and forecasted water level best model performance for 14 days' time-series scaled pattern. Figure 9 (a) has the value of observed and forecasted value for ANFIS model with EMD decomposition with best performance for RMSE and NSE, while Figure 9(b) has the value for ANFIS

model with DWT decomposition that gain lowest MAE. Finally, Figure 10 present the observed and forecasted value for monthly time-series scaled pattern using LSTM model with DWT decomposition. This final model has the best performance in all the performance measurement.
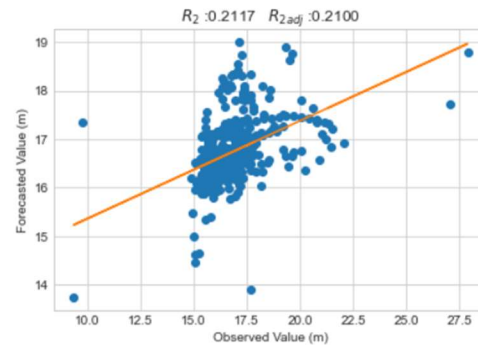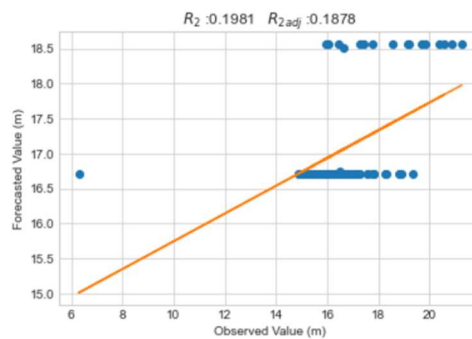


*Figure 7: Observed and Forecasted Water Level for ANFIS Model with EMD Decomposition (7 Days' Time-Series Scaled Pattern)*
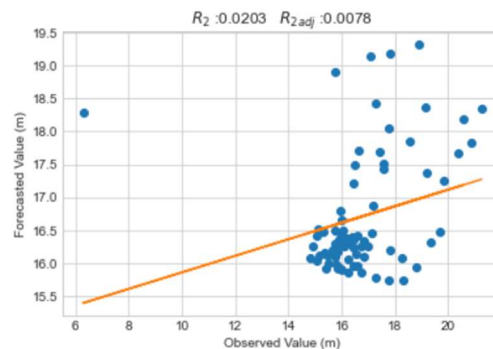
(a)



(b)



*Figure 8: (a) Observed and Forecasted Water Level for ANN Model with DWT Decomposition (b) Observed and Forecasted Water Level for ANFIS Model with DWT Decomposition (10 Days' Time-Series Scaled Pattern)*
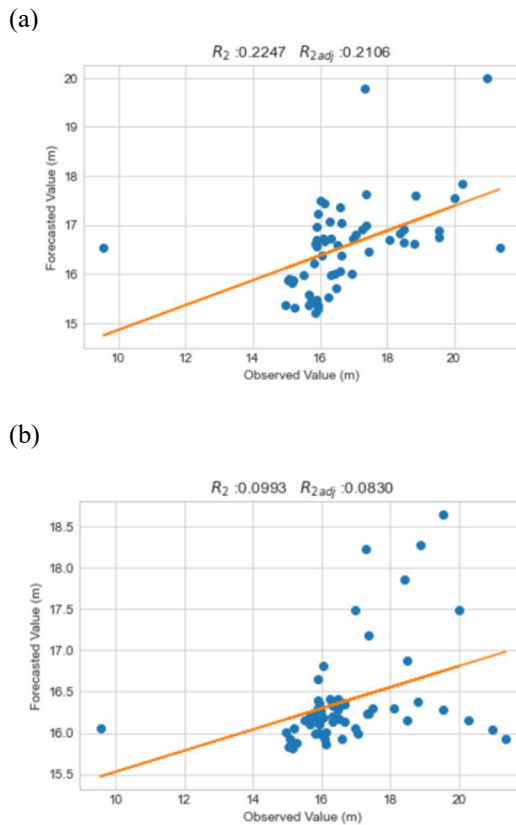
(a)



(b)



*Figure 9: (a) Observed and Forecasted Water Level for ANFIS Model with EMD Decomposition (b) Observed and Forecasted Water Level for ANFIS Model with DWT Decomposition (14 Days' Time-Series Scaled Pattern)*
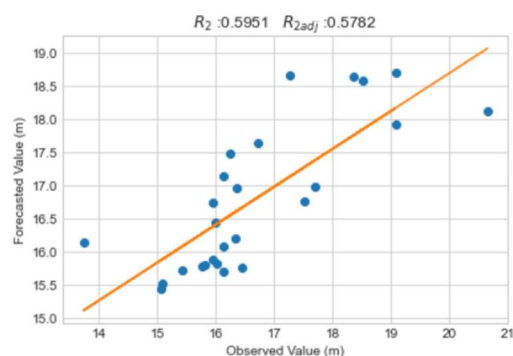


*Figure 10:  Observed and Forecasted Water Level for LSTM Model with DWT Decomposition (Monthly Time-Series Scaled Pattern)*

The above scatterplots show the distribution of data resulted from the best performance model for each time-series scale pattern. Although all models seem to have a positive relationship between the variables, yet the spread of the data seems to be varied. As a result, it is found that ANFIS model with EMD decomposition that used 7 days' time-series scaled pattern seem to centre its result in 17.5:17 line and less scattered. For model with 10 days' time-series scaled pattern, it is shown that ANN model with DWT decomposition seem to give monotonous forecast while the ANFIS model with DWT decomposition is more scattered and have a weak relation. In model using 14 days' time-series pattern, it indicates a moderate relation in ANFIS model with EMD decomposition while showing more scattered data with DWT decomposition.

For LSTM model with DWT decomposition that used monthly time-series scaled pattern, it is found that the data is much more scattered in the concentration of the fit line. It is also appeared that this model has the highest value of coefficient determination ($R_2$) and adjusted coefficient determination ($R_{2adj}$). This indicate that the observed data has a significant agreement with the forecasted data [36]. This model also has achieved the best performance of RMSE, MAE and NSE measurement.

The result of this study has indicated that longer time-series scaled pattern would provide a more accurate forecast for the water level of Kuala Krai in a longer lead time, in this case it would be monthly. The findings of this study are in line with existing literatures in which the use of monthly time-series scaled pattern has been proven to provide better hydrological forecast. As such, forecasting stream flow of Hurman River in Turkey and Diyalah and Lesser Zab Rivers in Iraq has showing better forecast for 1 months ahead [7]. The use of monthly data to forecast river flow of Tigris river also has been proven to achieve the best performances [6]. Other work in providing river flow forecast for Zarrinehrud River [8] using monthly data towards machine learning models has produce better result.

When data is scaled with longer time, it captured more of the seasonal characteristics of the rainfall. Longer lead time is very important when curb with the flood risk and resources. It normally used for strategical purpose of flood management. It helps in reducing the loss and impact of flood to the surrounding especially the infrastructure, crops, and vegetation.

Rainfall in Kelantan River is influenced by the movement of the monsoon. It mostly happens during the end of the year. Hence, 7 days, 10 days, and 14 days scaled data are best used for operational and monitoring purposes. It can be used to give early

warning to the people and triggering an action for evacuation or risk mitigation.

## 5. CONCLUSION

In this study, we investigate the effect of different time-series scaled pattern to forecast the water level of Kuala Krai with various lead time. The scaled time-series data collected and categorized into 7 days, 10 days, 14 days, and monthly to forecast water level with 7 days, 10 days, 14 days, and monthly lead time respectively. Cumulative rainfall data from eight rainfall stations along Lebir River and Galas River are the input to forecast the water level of Kuala Krai. Forecasting model was developed using Artificial Neural Networks (ANN), Adaptive Neuro Fuzzy Inference System (ANFIS) and Long-Short Term Memory (LSTM).

In conclusion, the results of the performance measurement shows that it is proven to better use longer scaled data for a longer lead time. In this case, the monthly scaled time-series data utilized in forecasting water level with monthly lead time has achieved better performance among the other models that using medium scaled time-series data to forecast water level with medium lead time. Nevertheless, the use of scaled time-series data can be influenced by the purpose of forecasting. Longer scaled time-series data in longer lead time water level forecasting provide great information for planning and strategical purpose while medium scaled time-series data for medium lead time water level forecasting is important to provide people with the current information and organizing daily operation.

Extended study in the future can be conducted to investigate other hydrological parameters which has the same non-linear and dynamic characteristics with rainfall such as streamflow and runoff. Other time scales can also be used such as seasonal and annual into the forecast. Future study may provide researcher with a new insights and knowledge of forecasting the water level phenomenon with best lead time and input scale.

## REFERENCES

[1]    World Meteorological Organization, *TECHNICAL REGULATIONS: Volume III - Hydrology*, no. 2. 2006.

[2]    S. Ogale and S. Srivastava, "Modelling and short term forecasting of flash floods in an urban environment," *25th Natl. Conf. Commun. NCC 2019*, pp. 1–6, 2019, doi: 10.1109/NCC.2019.8732193.

[3]    C. T. Cheng, W. J. Niu, Z. K. Feng, J. J. Shen, and K. W. Chau, "Daily reservoir runoff forecasting method using artificial neural network based on quantum-behaved particle swarm optimization," *Water (Switzerland)*, vol. 7, no. 8, pp. 4232–4246, 2015, doi: 10.3390/w7084232.

[4]    A. Altunkaynak and T. A. Nigussie, "Prediction of daily rainfall by a hybrid wavelet-season-neuro technique," *J. Hydrol.*, vol. 529, no. P1, pp. 287–301, 2015, doi: 10.1016/j.jhydrol.2015.07.046.

[5]    Z. M. Yaseen, M. F. Allawi, A. A. Yousif, O. Jaafar, F. M. Hamzah, and A. El-Shafie, "Non-tuned machine learning approach for hydrological time series forecasting," *Neural Comput. Appl.*, vol. 30, no. 5, pp. 1479–1491, 2018, doi: 10.1007/s00521-016-2763-0.

[6]    Z. M. Yaseen, S. M. Awadh, A. Sharafati, and S. Shahid, "Complementary data-intelligence model for river flow simulation," *J. Hydrol.*, vol. 567, no. October, pp. 180–190, 2018, doi: 10.1016/j.jhydrol.2018.10.020.

[7]    A. Mahmood Al-Juboori and A. Guven, "A stepwise model to predict monthly streamflow," *J. Hydrol.*, vol. 543, pp. 283–292, 2016, doi: 10.1016/j.jhydrol.2016.10.006.

[8]    M. A. Ghorbani, H. A. Zadeh, M. Isazadeh, and O. Terzi, "A comparative study of artificial neural network (MLP, RBF) and support vector machine models for river flow prediction," *Environ. Earth Sci.*, vol. 75, no. 6, pp. 1–14, 2016, doi: 10.1007/s12665-015-5096-x.

[9]    B. Choubin, S. Khalighi-Sigaroodi, A. Malekian, and Ö. Kişi, "Multiple linear regression, multi-layer perceptron network and adaptive neuro-fuzzy inference system for forecasting precipitation based on large-scale climate signals," *Hydrol. Sci. J.*, vol. 61, no. 6, pp. 1001–1009, 2016, doi: 10.1080/02626667.2014.966721.

[10]   W. chuan Wang, K. wing Chau, L. Qiu, and Y. bo Chen, "Improving forecasting accuracy of medium and long-term runoff using artificial neural network based on EEMD decomposition," *Environ. Res.*, vol. 139, pp. 46–54, 2015, doi: 10.1016/j.envres.2015.02.002.

[11] A. Mosavi, P. Ozturk, and K. W. Chau, "Flood prediction using machine learning models: Literature review," *Water (Switzerland)*, vol. 10, no. 11, pp. 1–40, 2018, doi: 10.3390/w10111536.

[12] T. Zhou, F. Wang, and Z. Yang, "Comparative analysis of ANN and SVM models combined with wavelet preprocess for groundwater depth prediction," *Water (Switzerland)*, vol. 9, no. 10, 2017, doi: 10.3390/w9100781.

[13] S. Zhu, J. Zhou, L. Ye, and C. Meng, "Streamflow estimation by support vector machine coupled with different methods of time series decomposition in the upper reaches of Yangtze River, China," *Environ. Earth Sci.*, vol. 75, no. 6, 2016, doi: 10.1007/s12665-016-5337-7.

[14] M. Tayyab, J. Zhou, X. Dong, I. Ahmad, and N. Sun, "Rainfall-runoff modeling at Jinsha River basin by integrated neural network with discrete wavelet transform," *Meteorol. Atmos. Phys.*, vol. 131, no. 1, pp. 115–125, 2019, doi: 10.1007/s00703-017-0546-5.

[15] M. Tayyab, J. Zhou, R. Adnan, and X. Zeng, "Application of Artificial Intelligence Method Coupled with Discrete Wavelet Transform Method," *Procedia Comput. Sci.*, vol. 107, no. Icict, pp. 212–217, 2017, doi: 10.1016/j.procs.2017.03.081.

[16] M. Ravansalar, T. Rajaee, and O. Kisi, "Wavelet-linear genetic programming: A new approach for modeling monthly streamflow," *J. Hydrol.*, vol. 549, no. April, pp. 461–475, 2017, doi: 10.1016/j.jhydrol.2017.04.018.

[17] N. E. Huang *et al.*, "The empirical mode decomposition and the Hubert spectrum for nonlinear and non-stationary time series analysis," *Proc. R. Soc. A Math. Phys. Eng. Sci.*, vol. 454, no. 1971, pp. 903–995, 1998, doi: 10.1098/rspa.1998.0193.

[18] F. F. Li, Z. Y. Wang, X. Zhao, E. Xie, and J. Qiu, "Decomposition-ANN Methods for Long-Term Discharge Prediction Based on Fisher's Ordered Clustering with MESA," *Water Resour. Manag.*, vol. 33, no. 9, pp. 3095–3110, 2019, doi: 10.1007/s11269-019-02295-8.

[19] Q. F. Tan *et al.*, "An adaptive middle and long-term runoff forecast model using EEMD-ANN hybrid approach," *J. Hydrol.*, vol. 567, pp. 767–780, 2018, doi: 10.1016/j.jhydrol.2018.01.015.

[20] S. Huang, J. Chang, Q. Huang, and Y. Chen, "Monthly streamflow prediction using modified EMD-based support vector machine," *J. Hydrol.*, vol. 511, pp. 764–775, 2014, doi: 10.1016/j.jhydrol.2014.01.062.

[21] N. E. Huang and Z. Wu, "A Review On Hilbert-Huang Transform : Method And Its Applications," *Rev. Geophys.*, vol. 46, no. 2007, pp. 1–23, 2008, doi: 10.1029/2007RG000228.1.INTRODUCTION.

[22] Q. Ouyang, W. Lu, X. Xin, Y. Zhang, W. Cheng, and T. Yu, "Monthly rainfall forecasting using EEMD-SVR based on phase-space reconstruction," *Water Resour. Manag.*, vol. 30, no. 7, pp. 2311–2325, 2016, doi: 10.1007/s11269-016-1288-8.

[23] W. chuan Wang, K. wing Chau, D. mei Xu, and X. Y. Chen, "Improving Forecasting Accuracy of Annual Runoff Time Series Using ARIMA Based on EEMD Decomposition," *Water Resour. Manag.*, vol. 29, no. 8, pp. 2655–2675, 2015, doi: 10.1007/s11269-015-0962-6.

[24] Q. Chen *et al.*, "Empirical mode decomposition based long short-term memory neural network forecasting model for the short-term metro passenger flow," *PLoS One*, vol. 14, no. 9, pp. 1–18, 2019, doi: 10.1371/journal.pone.0222365.

[25] V. Nourani, T. R. Khanghah, and A. H. Baghanam, "Implication of feature extraction methods to improve performance of hybrid wavelet-ann rainfall-runoff model," *Case Stud. Intell. Comput. Achiev. Trends*, pp. 457–498, 2014, doi: 10.1201/b17333.

[26] V. Nourani, M. T. Alami, and F. D. Vousoughi, "Wavelet-entropy data pre-processing approach for ANN-based groundwater level modeling," *J. Hydrol.*, vol. 524, pp. 255–269, 2015, doi: 10.1016/j.jhydrol.2015.02.048.

[27] N. B. M. Khairudin, N. B. Mustapha, T. N. B. M. Aris, and M. B. Zolkepli, "Comparison of Machine Learning Models for Rainfall Forecasting," *2020 Int. Conf. Comput. Sci. Its Appl. Agric. ICOSICA 2020*, 2020, doi: 10.1109/ICOSICA49951.2020.9243275.

[28] F. Mekanik, M. A. Imteaz, and A. Talei, "Seasonal rainfall forecasting by adaptive network-based fuzzy inference system (ANFIS) using large scale climate signals," *Clim. Dyn.*, vol. 46, no. 9–10, pp. 3097–3111, 2015, doi: 10.1007/s00382-015-2755-

2.

[29]  I. R. Widiasari, L. E. Nugoho, Widyawan, and R. Efendi, "Context-based Hydrology Time Series Data for A Flood Prediction Model Using LSTM," *Proc. - 2018 5th Int. Conf. Inf. Technol. Comput. Electr. Eng. ICITACEE 2018*, pp. 385–390, 2018, doi: 10.1109/ICITACEE.2018.8576900.

[30]  M. Perumal and R. K. Price, "A Fully Mass Conservative Variable Parameter Mccarthy-Muskingum Method: Theory and Verification," *J. Hydrol.*, vol. 502, pp. 89–102, 2013, doi: 10.1016/j.jhydrol.2013.08.023.

[31]  J. Abbot and J. Marohasy, "Using lagged and forecast climate indices with artificial intelligence to predict monthly rainfall in the Brisbane catchment, Queensland, Australia," *Int. J. Sustain. Dev. Plan.*, vol. 10, no. 1, pp. 29–41, 2015, doi: 10.2495/SDP-V10-N1-29-41.

[32]  Z. M. Yaseen, R. C. Deo, I. Ebtehaj, and H. Bonakdari, "Hybrid Data Intelligent Models and Applications for Water Level Prediction," no. 1, pp. 121–139, 2018, doi: 10.4018/978-1-5225-4766-2.ch006.

[33]  B. Yadav, S. Ch, S. Mathur, and J. Adamowski, "Assessing the suitability of extreme learning machines (ELM) for groundwater level prediction," *J. Water L. Dev.*, vol. 32, no. 1, pp. 103–112, 2017, doi: 10.1515/jwld-2017-0012.

[34]  "Agensi Kerajaan Tempatan Berdaya Tahan," Accessed: Aug. 13, 2021. [Online]. Available: https://www.preventionweb.net/files/67981_67981resilientlocalgovernmentunitsk.pdf.

[35]  "River Level." http://infobanjir.water.gov.my/waterlevel_page.cfm?state=KEL (accessed Aug. 15, 2021).

[36]  J. Adamowski, H. F. Chan, S. O. Prasher, and V. N. Sharda, "Comparison of multivariate adaptive regression splines with coupled wavelet transform artificial neural networks for runoff forecasting in Himalayan micro-watersheds with limited data," *J. Hydroinformatics*, vol. 14, no. 3, pp. 731–744, 2012, doi: 10.2166/hydro.2011.044.