Scientific
Research
Publishing

# A Water Level Forecast of Pattani River in the Southern of Thailand by Deep Learning

**Prattana Deeprasertkul, Kanoksri Sarinnapakorn**

Hydro-Informatics Institute (Public Organization), Ministry of Higher Education, Science, Research and Innovation, Bangkok, Thailand
Email: prattana@hii.or.th

## Abstract

Nowadays, the deep learning methods are widely applied to analyze and predict the trend of various disaster events and offer the alternatives to make the appropriate decisions. These support the water resource management and the short-term planning. In this paper, the water levels of the Pattani River in the Southern of Thailand have been predicted every hour of 7 days forecast. Time Series Transformer and Linear Regression were applied in this work. The results of both were the water levels forecast that had the high accuracy. Moreover, the water levels forecasting dashboard was developed for using to monitor the water levels at the Pattani River as well.

## Keywords

Time Series Transformer, Linear Regression, Water Level Prediction, Data Cleansing

## 1. Introduction

Thailand faces flood events every year partly affects by continuous climate variability and change. Climate change can increase the disaster risk and will become more frequent and more severe. These flood damages affected the people's lives and properties and the country's economy. After the huge flood event in 2011, the National Hydro Data Center system was developed by Hydro-Informatics Institute (Public Organization) that is a central system for collecting water resources data, including spatial data, statistical data, water situation, and forecast data.

Nowadays, deep learning is widely applied to analyze and predict the trend of various disaster situations and offer the alternatives to make the appropriate decisions. Moreover, data visualization is used to display the data analytics which

can create innovations to support decision-making in water management planning and determine related policies. Deep learning techniques have been used for the water level forecasting. One of the deep learning models applied in the floods prediction and management is the Artificial Neural Network or ANN model. The ANN model developed by [1] uses water level and meteorological data as input data and estimates the water flow. The ANN model was proved that had good performance and be able to use for water flow prediction [2]. However, most of the river water level prediction works utilize hydrological data and meteorological data, which are categories of time series data, as input data [3] [4]. Therefore, applying ANN model, there is a problem of insufficient memory when operating on time series data [5] [6]. The work in [7] proposed an LSTM model for flood forecasting. In this study, daily runoff and rainfall data were used as input datasets [2]. Nevertheless, the work in [8] applied the transformer-based model to produce video event proposals. Transformer-based model was also applied to detect and recognize the actions of the person of interest. Moreover, the structure of LSTM model is different from transformer-based model that is transformer-based model which has an Encoder and a Decoder processes but they are not available in LSTM model. The works in [9] [10] also present a transformer-based model applied to their works and achieve a high accuracy of 97%.

Therefore in this work, a system for forecasting the water levels in the Pattani River in the southern of Thailand has been developed to predict the hourly water levels in the river for 7 days forecast. The water level datasets are measured by 5 telemetered stations in the Pattani River and sent to collection in the National Hydro Data Center data warehouse. Time Series Transformer or TST and Linear Regression were applied in this work. TST is deep learning model based on the transformer-based model. Transformer has a self-attention process that introduces the useful data which are non-contiguous ranges to the model. This model is one of the most applied in time series classification, regression, and forecasting works. In addition, the Linear Regression model is a supervised learning model that learns the data pattern from the training datasets. It is also the statistical analysis used to predict the relationship between the variables. Therefore, the water level forecasting results were with a quite good accuracy in this work.

The remainder of this paper is organized as follows: Section 2 presents the technology backgrounds of this work; Section 3 presents the water level forecasting methodologies; Section 4 presents the water levels forecasting results and dashboard. Finally is the conclusion of this work.

## 2. Backgrounds

### 2.1. Inverse Distance Weighting or IDW

IDW is one of the data estimation methods that the assigned values to unknown points are calculated with a weighted average of the values available at the known points [11]. For example, there are 5 know points, and an unknown point is

needed to predict the value as shown in **Figure 1**. Inverse Distance Weighting equation is shown as equation (1).

$$w(x) = \frac{A}{B}$$

$$A = \sum_{i=1}^{n} \frac{1}{d(x, x_i)^p} u_i \tag{1}$$

$$B = \sum_{i=1}^{n} \frac{1}{d(x, x_i)^p}$$

where: $w$ is the predicted value; $d$ is the distance; $x$ is the unknown point; $x_i$ is the nth know point; $u_i$ is the value of the know point; and $p$ is the power.

## 2.2. Spline Interpolation

A spline is a type of polynomial function. In mathematics, splines are often used in a type of interpolation known as spline interpolation. Interpolation is used when there is a set of discrete data points and it is necessary to estimate other points of the same type of data from the given points. Polynomial interpolation is commonly used for small numbers of data points; this is a method that fits an order polynomial function to $n + 1$ data points [12].

## 2.3. MICE Interpolation

MICE or Multiple Imputations by Chained Equations is a popular approach to replace the missing values by predicting them using other features from the dataset [13].

## 2.4. Time Series Transformer or TST

TST model is a deep learning model that is based on the transformer model. The kind of work for the transformer model is related to sequence transformed to
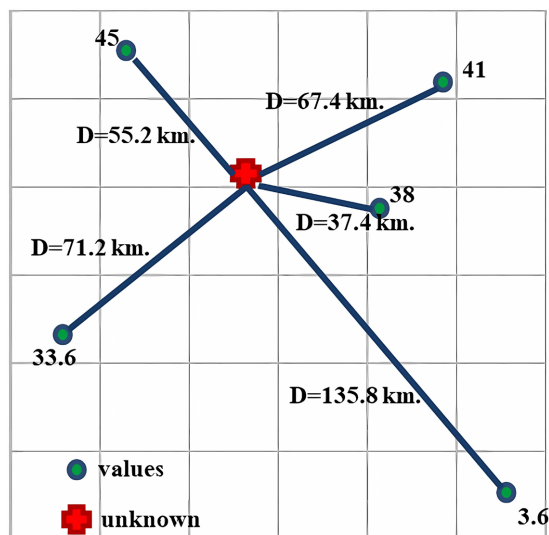


**Figure 1.** The concept of inverse distance weighting.

sequence, such as text translation, speech recognition, time series classification, etc. In general, the transformer model consists of 2 main parts, encoder and decoder, which has different functions. The encoder is a model input data processing to learn the initial information and extract important information from the data and then send it to the decoder. The decoder receives information from the encoder and other information that there may be more to process and generate results [14].

The encoder is only used in this work because it needs to reduce the restriction in the decoder side for a variety of methods and results as shown in **Figure 2**. In addition, the number of model parameters is less than a half of the model using both the encoder and decoder. Not reducing the number of parameters can only reduce the processing time but also solves the overfitting problems. Moreover, the TST model also uses a learnable positional encoding that can be learned by itself for supporting a variety of data types [14].

- Linear Regression is a supervised learning model that learns the data pattern from the training datasets [15]. It is also the statistical analysis used to predict the relationship between the variables. The equation as shown in **Figure 3** used for prediction in this work is in the form of multiple linear regression that is more than one independent variable can be added to predict Y or the water level values in the future. This is more complicated than one variable.

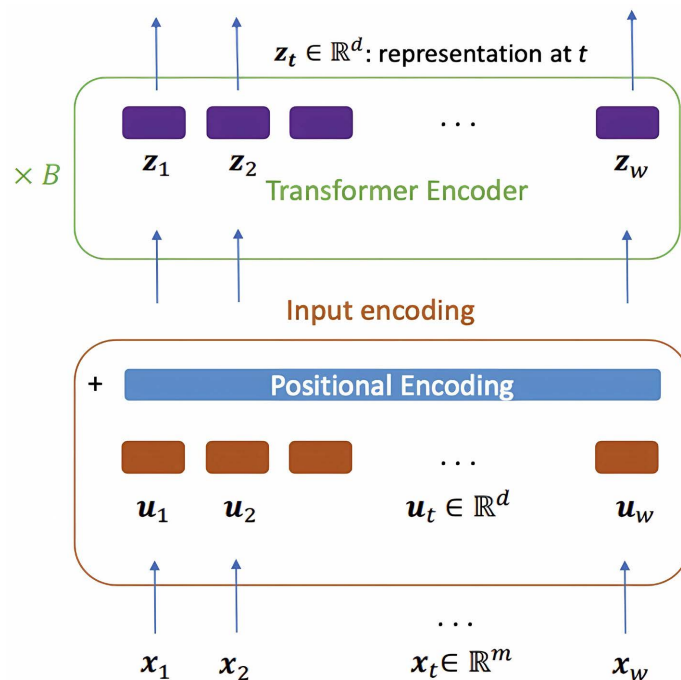## 3. Water Level Forecasting Methodologies

### 3.1. Data

- Water Levels Data



**Figure 2.** The TST model process.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \varepsilon$$

**Figure 3.** Linear regression equation.

The raw water levels data are stored in CSV files that a separate file is collected for each telemetering station data. The datasets consist of 6 stations that means we have 6 datasets, named BLGTD03, BLGTD05, FOP045, X40A, X10A and X275.

● Rainfall Data

The raw rainfall data are also stored in CSV files that the separate file is collected for each telemetering station data. The datasets consist of 21 stations which means we have 21 datasets, named FOP042, FOP044, FOP047, KABG, KTUM, MUNG, STH001, STH002, STH003, STH004, STH006, STH008, STH012, STH015, STH017, STH018, STH020, STH024, STH027, STH029, and STH030. All rainfall telemetering stations are around the Pattani River area.

● Reservoir Data

The raw data of the reservoir discharge volume is saved in CSV file. There is only one station of Bang Lang Reservoir at the Pattani River upstream.

## 3.2. Data Pre-Processing

This section describes the steps of data preparation in order to be ready for use in the next section. In this work, there are the following preprocesses that are the data cleansing, the data interpolation, and the data augmentation. The workflow of the data preparation is shown in **Figure 4**. The details are as follows.

● Data Cleansing

This section describes the steps and the function coding of the anomaly data cleansing. The source codes were divided into two parts. First, the source codes used in water levels data cleansing. Second, the source codes were used to clean the rainfall data. The workflows of water levels data cleansing and rainfall data cleansing are shown in **Figure 5** and **Figure 6**, respectively.

● Remove_waterlevel_anomaly(df_waterlevel)

It is the work function for obtaining the water level data tables that were previously processed for preparing the raw data. The details of the functions are as follows.

Step 1: Obtain water level data in tabular form via the df_waterlevel variable.

Step 2: Replace the water level values that are equal to −999 with the Nan value. −999 is the missing value of a water level censor in telemetering station.
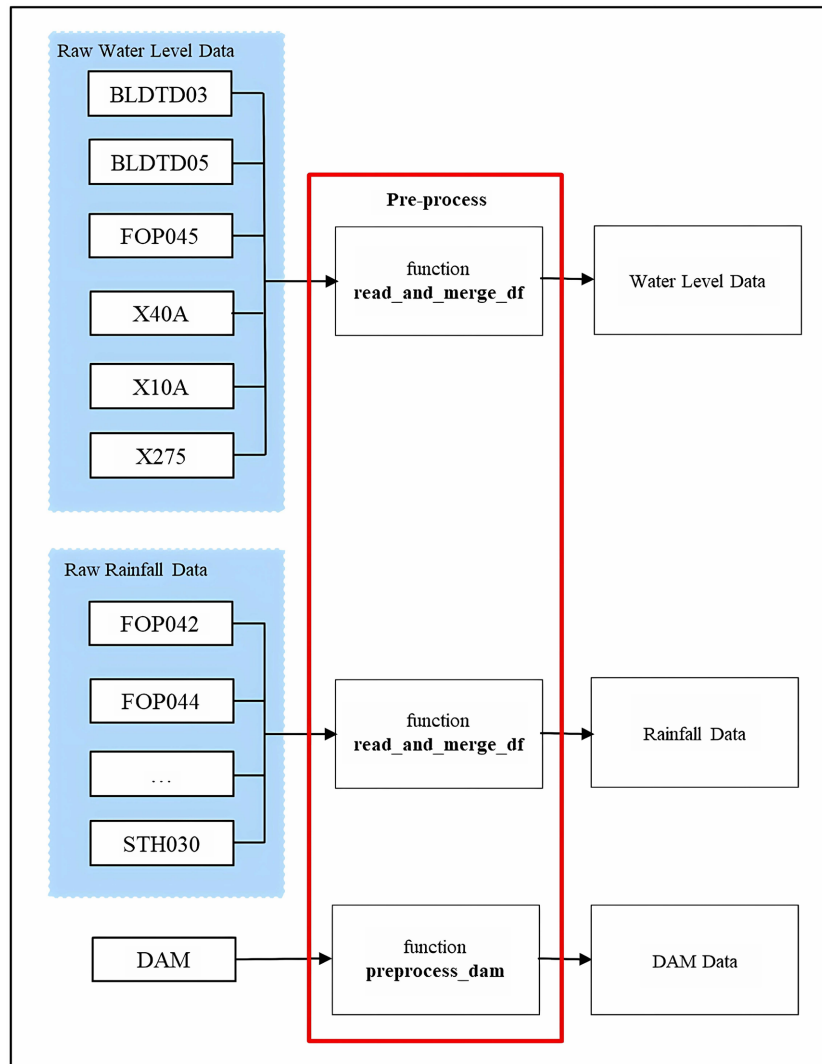
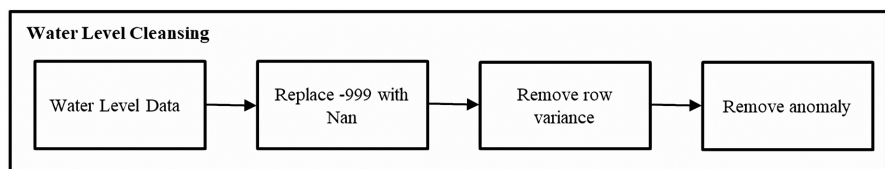**Figure 4.** Overview of the raw data preparation process.

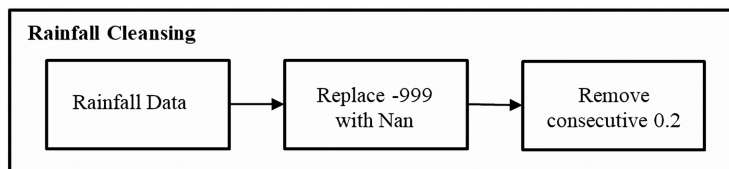**Figure 5.** The workflow of water level cleansing.

**Figure 6.** The workflow of rainfall cleansing.

Step 3: Remove the low variances of water level data. The water level values which had a little change over a continuous time period were removed. If they

were below 9e−5 (0.00009) for 48 hours, these values were considered to be anomalies and then were removed.

- Remove_rainfall_anomaly(df_rainfall, df_geo)

It is the work function for obtaining the rainfall data tables that were previously processed for preparing the raw data. The details of the functions are as follows.

Step 1: Obtain rainfall data in tabular form via the df_rainfall variable and obtain the location of each telemetering station in tabular form via the df_geo variable.

Step 2: Replace the rainfall values that are equal to −999 with the Nan value. −999 is the missing value of a rainfall censor in telemetering station.

Step 3: Select the window_size period for example 4 hours. Therefore, every 4 hours, detect the rainfall values which equal to 0.2, consecutively. Then remove these rainfall data that are anomaly values.

- Data Interpolation
  - Water level data
    - Spline Interpolation is the polynomial function. This method is applied to replace the values for the missing data. However, this method was only used for the data that were lost in less than 6 consecutive hours.
    - MICE Interpolation is assumed that the missing data can be replaced by other data and must not be the same type of missing data. Therefore, in case of the water level data is missing, the values were estimated with nearby station data by using the iterative imputer function from Scikit-Learn.
  - Rainfall data

Inverse Distance Weighting or IDW is a continuous range estimation method. By analyzing which rainfall telemetering stations affected the stations desired to estimate the values. The distance is an important factor in the correlation between two stations.

- Data Augmentation

Data augmentation is the process of adding data by generating data at the interesting time. For model development to be effective, the data must have the sufficient quantity [13]. However, rainfall data and reservoir data were expanded according to that time period. There are three ways to increase the amount of data as follows:

  - Jittering

Jittering is adding Gaussian noise to the original signal. The mean and standard deviation values are used to fill in.

  - Scaling

Scaling is multiplying the original signal by a fixed value. The constant is derived from the Gaussian noise with mean and standard deviation values.

  - Magnitude Warping

The original signal is multiplied by the cubic spline function, the Knot value is equal to 4, and the magnitude of the curve is derived from the Gaussian noise.

The data augmentation uses the signals for a time period which may be ran-

dom or selected specifically. The signals were taken through the 3-step processing as mentioned above. They are arranged in the following order: Scaling > Magnitude Warping > Jittering.

## 3.3. Water Level Forecasting Models

For studying and experimenting with various models, the conclusion that the models chosen to apply to the projection of the water levels was the TST and the Linear Regression models. Both produced the results with higher accuracy than the models that were tested in the previous period.

The accuracies of the TST and Linear Regression models were tested by comparing the results with the baseline model which was the Long Short-Term Memory or LSTM model. The predictions of the baseline model were different from TST model as follows.

- TST Model is the deep learning model based on the transformer model. In general, transformer model consists of two main parts that are encoder and decoder. However, the TST model will only be used in the encoder part of this work. The transformer model has a process of self-attention that introduces useful data ranges to the model. For example, the model may be particularly interested in data from the past 1$^{st}$ and 3$^{rd}$ days for further processing.

All data processed from the previous step were taken and used to develop the TST model. This process consists of two main steps: Feature Engineering is the data extraction from variables and Train is the data training process of the model by Hyperparameters Tuning which is the optimal parameters approach fitting the model. The all operations of the system are shown in Figure 7 that these functions are the main parts of model development. There are the steps of reading the data file, extracting the data from variables as well as finding the optimal parameters for the model. The details of the procedure are as follows.

Step 1: Input the water levels data, the rainfall data and the dam discharge volume data obtained from the data preparation process together.
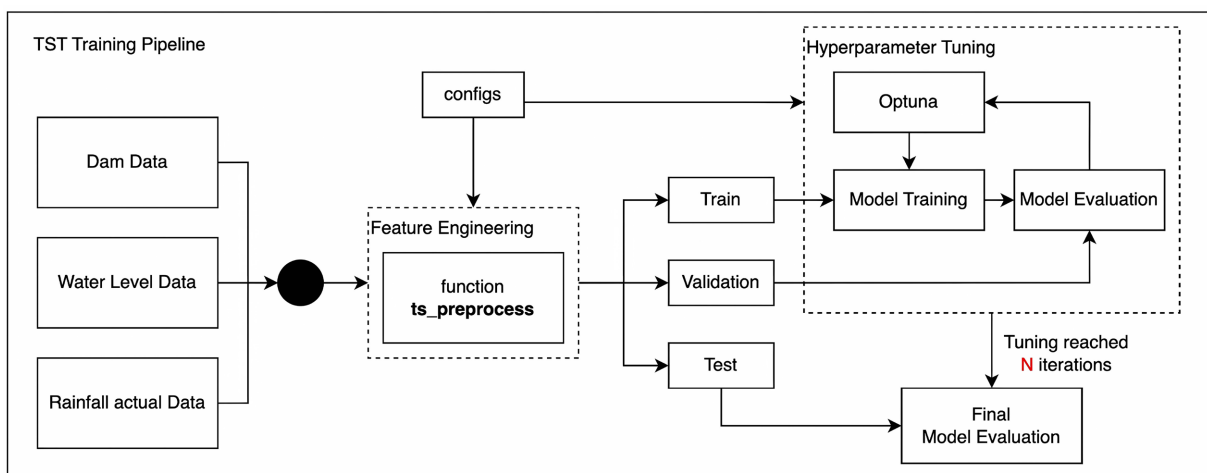


**Figure 7.** The overview of the TST model development process.

Step 2: Input the data as mentioned above into the data extracting process. The variables were selected based on the "configs.json" file which stored the variables used for training model. In this step, the function was "ts_preprocess".

Step 3: Divide the data into three groups that were Train, Validation and Test datasets based on the "configs.json" file.

Step 4: Train and select the best model by finding the optimal parameters using the "run_optuna_study" function of "OptunaTS" class. The details of finding parameters are as follows.

Step 4.1: Define the variables, the input data lengths, and the output data lengths which are the predicted data with the reference in the "configs.json" file.

Step 4.2: Run the model from the random parameters for the first time, and then the parameters were generated based on the parameters of the previous running time or the parameters of the model that produced the best results.

Step 4.3: Train and measure the model results. If the results were better than the previous running time, the model parameters were used in the next running time.

Step 4.4: Repeat steps 4.2 and 4.3 until the specified cycle was completed.

Step 4.5: Use the parameters of the best results to train the final model.

- Linear Regression Model is classified as supervised learning, that is, it has to learn the data pattern from training datasets to calculate in statistical regression and then returns the result which is a linear correlation. The equation used for prediction in this case is in the form of Multiple Linear Regression which means more than an independent variable can be added to predict the water levels in the future. All operations of the system are shown in Figure 8 and these functions are the main parts of model development. There are the steps of reading the data file, extracting the data from variables, and finding the optimal parameters for the model. The details of the procedure are as follows.
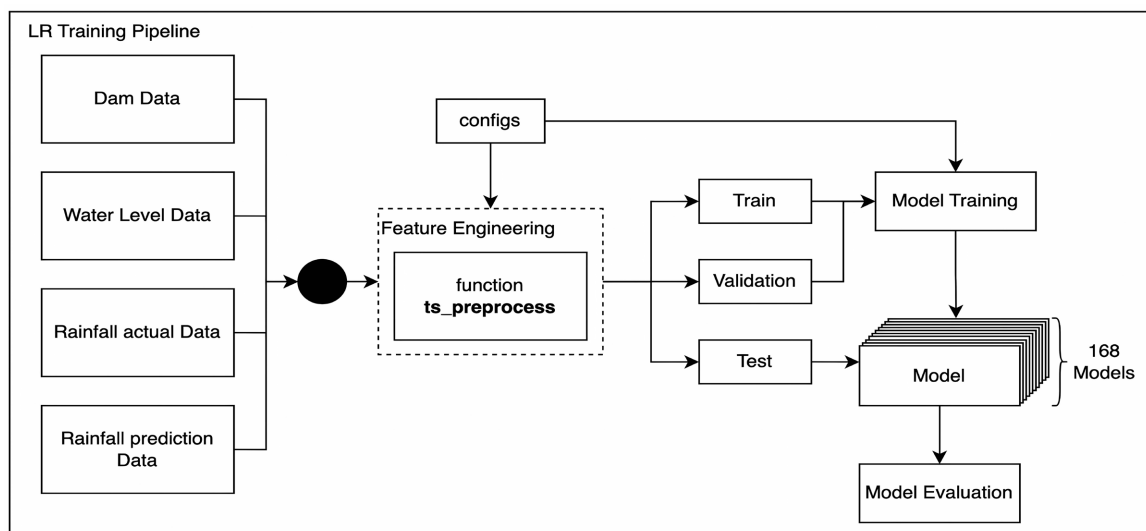


**Figure 8.** The overview of the linear regression model development process.

Step 1: Input the water levels data, the rainfall data and the dam discharge volume data obtained from the data preparation process together.

Step 2: Input the data as mentioned above into the data extracting process. The variables were selected based on the "configs.json" file which stored the variables used for training model. In this step, the function was "ts_preprocess".

Step 3: Divide the data into three groups that were Train, Validation and Test datasets based on the "configs.json" file.

Step 4: Define each variable type according to the following topics.

○  The valued variables in history: water levels, rainfall, and dam discharge volume data.

○  Forecasting variables in the future: forecast rainfall data. In the model training process, the actual forecasting rainfall variables were used for learning the best variable of the model.

Step 5: Develop 168 Linear Regression models, each of which was only responsible for one-hour prediction, e.g. the first model produces the prediction of the water level in the first hour. Therefore, the second model produces the prediction of the water level in the second hour.

Step 6: Train all 168 models using the variables defined in the "configs.json" file.

## 4. Water Level Forecasting Results and Dashboard

### 4.1. Evaluation

In evaluation, the data used to test the model's performance were divided for training the model and testing the model's performance. The period of data could be separated as follows:

1) Training Data: January 1, 2021 to October 31, 2022

2) Validation Data: November 1, 2022 to November 30, 2022

3) Testing Data: December 1, 2022 to December 31, 2022

The experimental results of TST model and Linear Regression model were measured by Coefficient of Determination or $R^2$ and Root Mean Square Errors or RMSE methods that were computed between the actual data and the predicted data. TST performance and Linear Regression performance were compared. Table 1 summarizes the results of $R^2$ and RMSE for TST model in 1 day forecast, 3 days forecast, and 7 days forecast, respectively. Table 2 summarizes the results of $R^2$ and RMSE for the Linear Regression model with future known values which are rainfall forecast data in 1 day forecast, 3 days forecast, and 7 days forecast, respectively.

### 4.2. Water Levels Forecasting Dashboard

In this section, the dashboard display is shown. It can be divided into 8 parts as follows: As shown in Figure 9, part 1 is the display of the Pattani River area map. Users can select symbols of interesting positions which have 3 symbols, 5 water level stations, 18 rainfall stations, and 1 reservoir. If users select any stations on

**Table 1.** The results of $R^2$ and RMSE for TST model.

| Water Level Stations | 1 day forecast | | 3 days forecast | | 7 days forecast | |
|---|---|---|---|---|---|---|
| | $R^2$ | RMSE | $R^2$ | RMSE | $R^2$ | RMSE |
| X40A | 0.63 | 0.20 | 0.64 | 0.31 | 0.70 | 0.42 |
| X10A | 0.73 | 0.11 | 0.61 | 0.13 | 0.54 | 0.16 |
| FOP045 | 0.68 | 0.27 | 0.49 | 0.36 | 0.38 | 0.45 |
| BLGTD03 | 0.63 | 0.25 | 0.43 | 0.34 | 0.46 | 0.46 |
| BLGTD05 | 0.81 | 0.11 | 0.65 | 0.13 | 0.54 | 0.15 |

**Table 2.** The results of $R^2$ and RMSE for Linear Regression model.

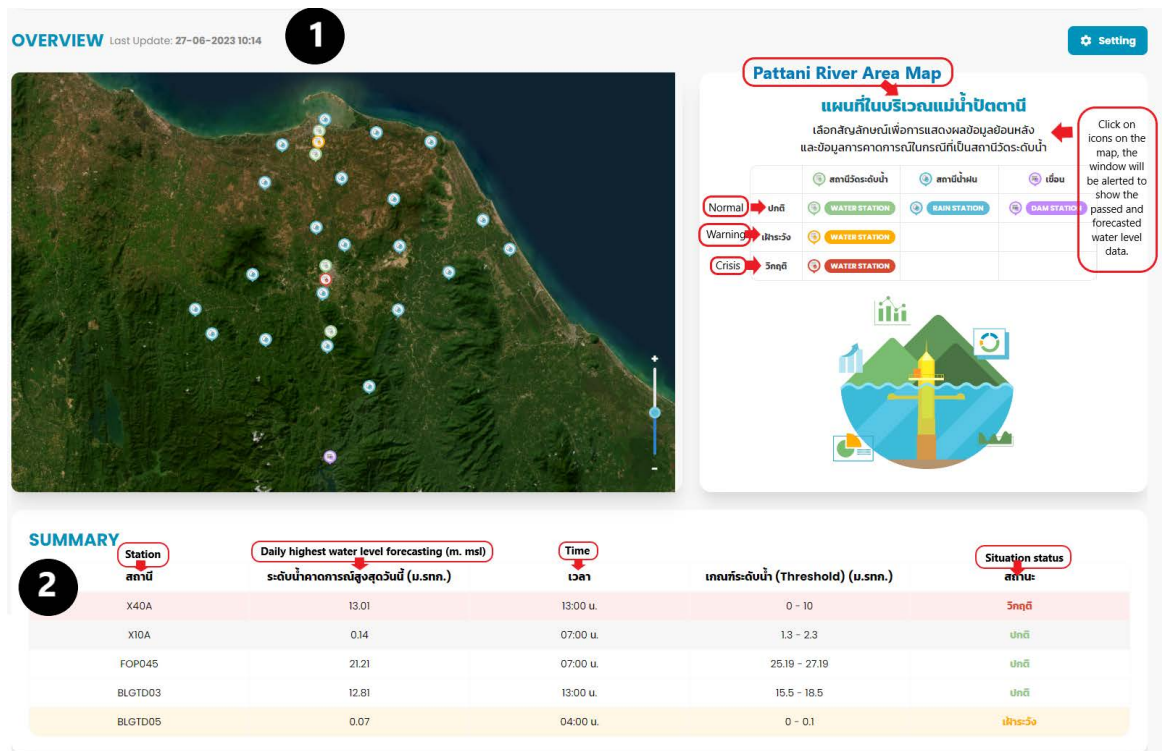| Water Level Stations | 1 day forecast | | 3 days forecast | | 7 days forecast | |
|---|---|---|---|---|---|---|
| | $R^2$ | RMSE | $R^2$ | RMSE | $R^2$ | RMSE |
| X40A | 0.64 | 0.29 | 0.64 | 0.53 | 0.70 | 0.58 |
| X10A | 0.59 | 0.12 | 0.48 | 0.19 | 0.54 | 0.23 |
| FOP045 | 0.14 | 8.93 | 0.07 | 10.06 | 0.02 | 10.16 |
| BLGTD03 | 0.41 | 0.51 | 0.43 | 0.86 | 0.46 | 0.93 |
| BLGTD05 | 0.44 | 0.20 | 0.40 | 0.31 | 0.28 | 0.32 |



**Figure 9.** Map of the Pattani River area and summary of the daily water level situation.

the map, the graphs in part 3 as shown in **Figure 10** and 4 are shown the details. In the right part of that map, users can mouse-click each symbol to let that type is active or inactive on the map.
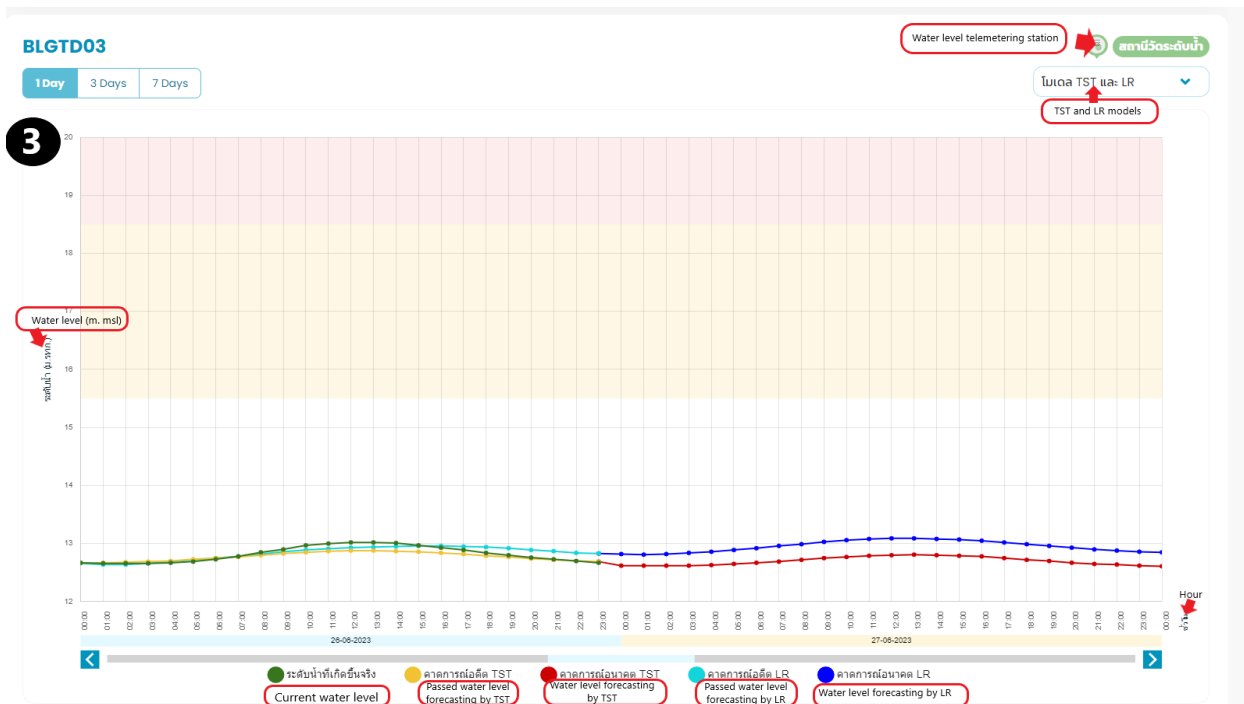
**Figure 10.** The graph of the water level, rainfall, dam discharge rate, and comparison of the actual and forecasted water levels.

Part 2 shown in **Figure 9** is the summary table of the daily water level situation of each water level station. The table shows the results of the highest water level prediction on that day, and the situation status of water level in each station in 3 criteria:

- Normal (green)—In case of the water level forecasting is lower than the water level threshold.
- Warning (yellow)—In case of the water level forecasting is nearly the water level threshold.
- Crisis (red)—In case of the water level forecasting is exceeded the water level threshold.

Part 3 is the graphs showing water level stations, rainfall stations and dam data which are changed according to the symbols that the user has chosen in Part 1 as shown in **Figure 10**. This part can choose to view data in 1, 3, and 7 days forecast for both models, TST and Linear Regression. Each data displays different information as follows.

- Water level stations display all 4 values: the past forecasting values, forecasting values, actual values, and water level threshold.
- Rainfall stations show the historical rainfall.
- Dam displays the historical water discharge rate.

Part 4 shown in **Figure 11** is the comparison graph of actual and forecasting water levels in 1, 3, and 7 days forecast for both TST and Linear Regression models by showing the average difference in each period. Users can select to compare data in the average values or percentage values formats.

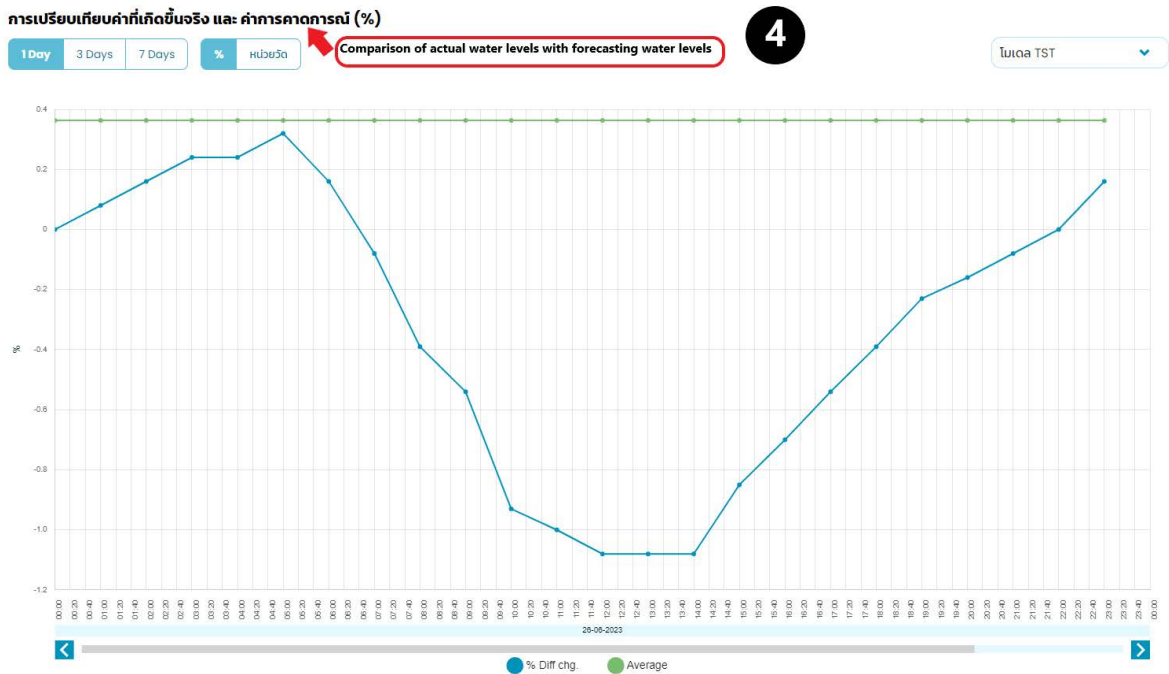Part 5 and 6 shown in **Figure 12**, display the model assessment results of the 1

**Figure 11.** The comparison graph of actual and forecasting water levels.
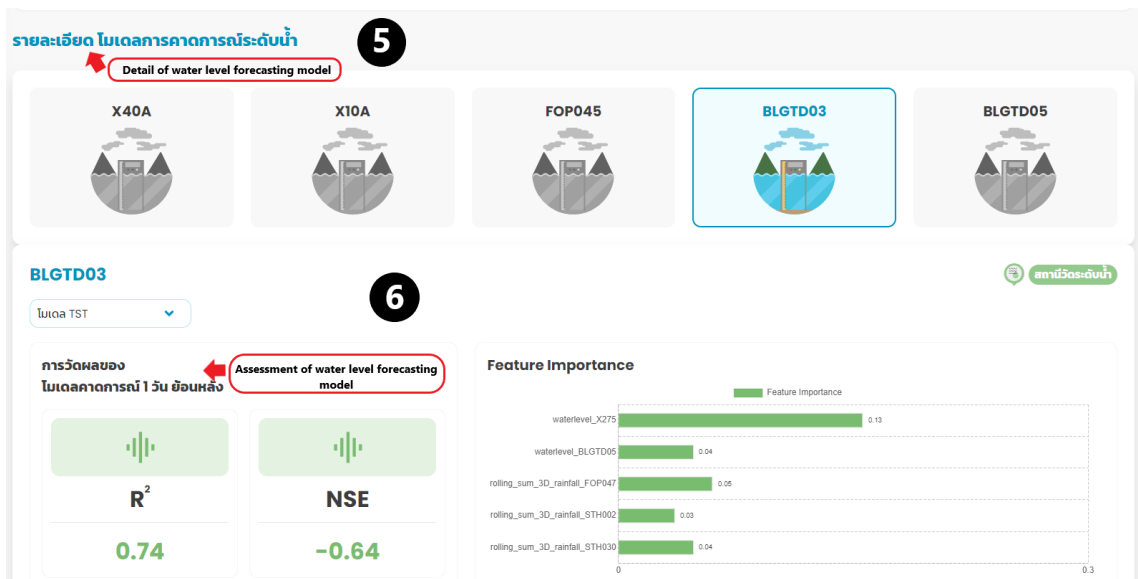


**Figure 12.** The model performance and feature importance.

day forecasting which the user can choose to view in each station. The display is shown in the R2 and RMSE values and the Feature Importance of the model.

Part 7 shown in **Figure 13** is the notifications updated every day that the model is run, and the status of notifications in each station corresponds to the summary table in Part 2.

Part 8 shown in **Figure 13** is the water level threshold setting. It can be set separately for each station criteria that the system administrator can be adjusted according to the suitability of each station.
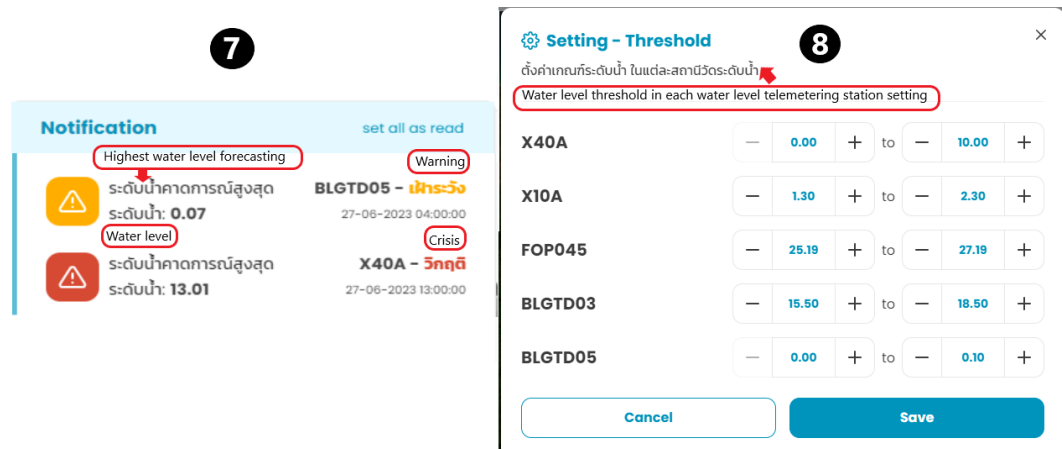
**Figure 13.** The notification of daily water level status at each station and the water levels threshold setting.

## 5. Conclusions

Because the number of datasets using in the model was quite limited, especially the data in the situations where the water level in the Pattani River was higher than normal, the results of TST model had the insufficient accuracy for 7 days forecast. Therefore, we have developed the further model that includes 7 days of rainfall forecasting obtained from HII rainfall forecasting model to be another variable for model.

In this work, we applied TST and Linear Regression models to predict the water levels every hour for 7 days forecasting in the Pattani River which produced the results with higher accuracy than other tested models.

## Acknowledgements

We would like to express our sincere thanks to the development team of Data Wow Company in Thailand for working on this project together with HII research team, especially the system.

## Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

## References

[1] Panigrahi, B.K., Das, S., Nath, T.K., *et al.* (2018) An Application of Data Mining Techniques for Flood Forecasting: Application in Rivers Daya and Bhargavi, India. *Journal of the Institution of Engineers* (*India*): *Series B*, **99**, 331-342. https://doi.org/10.1007/s40031-018-0333-9

[2] Cho, M., Kim, C., Jung, K. and Jung, H. (2022) Water Level Prediction Model Applying a Long Short-Term Memory (LSTM)-Gated Recurrent Unit (GRU) Method for Flood Prediction. *Water*, **14**, Article No. 2221. https://doi.org/10.3390/w14142221

[3] Liu, C., Guo, L., Ye, L., Zhang, S., Zhao, Y. and Song, T. (2018) A Review of Advances in China's Flash Flood Early-Warning System. *Natural Hazards*, **92**, 619-634.

https://doi.org/10.1007/s11069-018-3173-7

[4]  Maspo, N.A., Harun, A.N.B., Goto, M., Cheros, F., Haron, N.A. and Nawi, M.N.M. (2020) Evaluation of Machine Learning Approach in Flood Prediction Scenarios and Its Input Parameters: A Systematic Review. *IOP Conference Series: Earth and Environmental Science*, **2020**, Article ID: 012038. https://doi.org/10.1088/1755-1315/479/1/012038

[5]  Govindaraju, R.S. (2000) Artificial Neural Networks in Hydrology. I: Preliminary Concepts. *Journal of Hydrologic Engineering*, **5**, 115-123. https://doi.org/10.1061/(ASCE)1084-0699(2000)5:2(115)

[6]  Mijwel, M.M. (2018) Artificial Neural Networks Advantages and Disadvantages. https://www.linkedin.com/pulse/artificial-neural-networks-advantages-disadvantages-maad-m-mijwel

[7]  Le, X.-H., Ho, H.V., Lee, G. and Jung, S. (2019) Application of Long Short-Term Memory (LSTM) Neural Network for Flood Forecasting. *Water*, **11**, Article No. 387. https://doi.org/10.3390/w11071387

[8]  Zhao, Y., Li, J. and Yu, L. (2017) A Deep Learning Ensemble Approach for Crude Oil Price Forecasting. *Energy Economics*, **66**, 9-16. https://doi.org/10.1016/j.eneco.2017.05.023

[9]  Cheng, Y., Dai, Z., Ji, Y., Li, S., Jia, Z., Hirota, K. and Dai, Y. (2020) Student Action Recognition Based on Deep Convolutional Generative Adversarial Network. 2020 *Chinese Control and Decision Conference* (*CCDC*), Hefei, 22-24 August 2020, 128-133. https://doi.org/10.1109/CCDC49329.2020.9164040

[10] Rashmi, M., Ashwin, T.S. and Guddeti, R.M.R. (2021) Surveillance Video Analysis for Student Action Recognition and Localization inside Computer Laboratories of a Smart Campus. *Multimedia Tools and Applications*, **80**, 2907-2929. https://doi.org/10.1007/s11042-020-09741-5

[11] Chen, F.W. and Liu, C.W. (2012) Estimation of the Spatial Rainfall Distribution Using Inverse Distance Weighting (IDW) in the Middle of Taiwan. *Paddy and Water Environment*, **10**, 209-222. https://doi.org/10.1007/s10333-012-0319-1

[12] Childress, H.R. (2023) What Is a Spline? https://www.allthescience.org/what-is-a-spline.htm

[13] Prabhakaran, S. (2023) MICE Imputation: How to Predict Missing Values Using Machine Learning in Python. https://www.machinelearningplus.com/machine-learning/mice-imputation

[14] Zerveas, G., Jayaraman, S., Patel, D., Bhamidipaty, A. and Eickhoff, C. (2020) A Transformer-Based Framework for Multivariate Time Series Representation Learning. *KDD* '21: *Proceedings of the* 27*th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, Singapore, 14-18 August 2021, 2114-2124. https://doi.org/10.1145/3447548.3467401

[15] Kanade, V. (2023) What Is Linear Regression? Types, Equation, Examples, and Best Practices for 2022. https://www.spiceworks.com/tech/artificial-intelligence/articles/what-is-linear-regression