

Research Article

Bayesian Regularized Neural Network Model Development for Predicting Daily Rainfall from Sea Level Pressure Data: Investigation on Solving Complex Hydrology Problem

Lu Ye ¹, Saadya Fahad Jabbar,² Musaddak M. Abdul Zahra,^{3,4} and Mou Leong Tan ⁵

¹School of Computer Science, Baoji University of Arts and Sciences, Baoji 721007, China

²College of Education for Human Science-ibn Rushed, University of Baghdad, Baghdad, Iraq

³Computer Techniques Engineering Department, Al-Mustaqbal University College, 51001 Hillah, Babil, Iraq

⁴Electrical Engineering Department, College of Engineering, University of Babylon, Hilla, Babil, Iraq

⁵Geoinformatic Unit, Geography Section, School of Humanities, Universiti Sains Malaysia, 11800 USM, Pulau Pinang, Malaysia

Correspondence should be addressed to Mou Leong Tan; mouleong@usm.my

Received 9 November 2020; Revised 18 December 2020; Accepted 18 March 2021; Published 1 April 2021

Academic Editor: Zaher Mundher Yaseen

Copyright © 2021 Lu Ye et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Prediction of daily rainfall is important for flood forecasting, reservoir operation, and many other hydrological applications. The artificial intelligence (AI) algorithm is generally used for stochastic forecasting rainfall which is not capable to simulate unseen extreme rainfall events which become common due to climate change. A new model is developed in this study for prediction of daily rainfall for different lead times based on sea level pressure (SLP) which is physically related to rainfall on land and thus able to predict unseen rainfall events. Daily rainfall of east coast of Peninsular Malaysia (PM) was predicted using SLP data over the climate domain. Five advanced AI algorithms such as extreme learning machine (ELM), Bayesian regularized neural networks (BRNNs), Bayesian additive regression trees (BART), extreme gradient boosting (xgBoost), and hybrid neural fuzzy inference system (HNFIS) were used considering the complex relationship of rainfall with sea level pressure. Principle components of SLP domain correlated with daily rainfall were used as predictors. The results revealed that the efficacy of AI models is predicting daily rainfall one day before. The relative performance of the models revealed the higher performance of BRNN with normalized root mean square error (NRMSE) of 0.678 compared with HNFIS (NRMSE = 0.708), BART (NRMSE = 0.784), xgBoost (NRMSE = 0.803), and ELM (NRMSE = 0.915). Visual inspection of predicted rainfall during model validation using density-scatter plot and other novel ways of visual comparison revealed the ability of BRNN to predict daily rainfall one day before reliably.

1. Introduction

Rainfall decides agricultural activities, ecology, and environment of a region. Therefore, it is considered as the key factor of social and economic development of any region [1]. At the same time, rainfall is the most influencing factor for different kinds of natural hazards [2]. Excess rainfall and floods are the most common natural hazards all over the earth [3]. Rainfall deficit and droughts are the most devastating disasters in terms of economic damages [4]. Rainfall is also the defining factor of many other kinds of natural hazards like landslides, soil erosion, and river bank subsidence [5, 6]. Therefore, the forecasting of rainfall is the major

topic of interest to hydrologists and disaster management scientists for many decades [7, 8].

Numerous attempts have been made for forecasting rainfall using various physical, empirical, and physio-empirical models. Such models have been used for forecasting rainfall at annual, seasonal, monthly, and daily scale in different regions of the world [9, 10]. Physical models are developed considering the synoptic climate and physical mechanism of interactions of different climatic variables causing rainfall in a region [11]. Mathematical models are developed to represent those physical processes to forecast rainfall in physical models. Such models are generally complex and need data and information about many land-

ocean-atmospheric variables. Still, the physical models are often not very successful in reliable forecasting of rainfall as it often makes the simple approximation of complex physical phenomena for model development [12]. Empirical models based on statistical methods have been employed as an alternative to physical models [13]. Such models are generally developed focusing a certain climate variable only like forecasting rainfall based on the statistical relationship of rainfall with its antecedent values or with other atmospheric variables related to it. The statistical model is always region-specific which means the model developed based on the statistical relationship of the regional rainfall with other atmospheric variables is only applicable for the region [14]. Therefore, such models are also known as empirical models. The advantages of empirical models are easy development and implementation and better accuracy compared with physical models in most of the cases. However, the major drawback of the empirical model is their complete dependency on historical data used for their development. Therefore, they cannot simulate rainfall for an unknown situation [15, 16]. For example, the rapid changes in rainfall that earth experienced in recent years and projected for the future cannot be forecasted by empirical models accurately as the models were not developed with such large variability in data [17]. Physical models are the option for forecasting in such situation. To tradeoff the disadvantages of both types of the modeling approach, recent focus is to develop physical-empirical models where empirical models are developed based on the variables physically responsible for the climate of the region. Therefore, more emphasis has been given on the development of the physical-empirical model for forecasting rainfall [18, 19]. Yim et al. [20] developed a physical-empirical model for the prediction of summer rainfall in southern China using several ocean-atmospheric variables. Yim et al. [9] used sea surface temperature in a regression model for the development of a physical-empirical model for the prediction of rainfall in Taiwan. Chen and Sun [21] used the physical-empirical model for performance improvement of a dynamical model in seasonal rainfall forecasting using sea level pressure (SLP) data. Pour et al. [22] employed SLP in development of physical-empirical models for the development of seasonal rainfall in Peninsular Malaysia.

Generally, statistical models are developed using different forms of regressions ranging from linear to higher-order polynomial (nonlinear) regressions. However, the relationship of rainfall with atmospheric variables responsible for rainfall is often highly intriguing and cannot be represented using generally used statistical methods. Machine learning (ML) algorithms are often used for developing complex regression analysis [23–25]. Therefore, the development of ML-based physical-empirical forecasting models has grown very fast in recent years.

Forecasting daily rainfall in tropical region is much different compared with other regions due to the presence of highly extreme values and year-round rainfall. Use of more efficient ML algorithms is required for forecasting daily rainfall in such region. Some of the ML models have shown their superiority over others in solving environmental prediction problems. Though there is no single method that

shows consistent excellence in solving all kinds of environmental problems in many regions, overall it has been noticed that some models often perform better in solving complex prediction problems which include extreme learning machine (ELM), Bayesian regularized neural networks (BRNNs), Bayesian additive regression trees (BART), extreme gradient boosting (xgBoost), and hybrid neural fuzzy inference system (HNFIS) [26–31]. Therefore, the comparative performance of these efficient algorithms in forecasting daily rainfall can provide the best forecasting model.

The main research objective is to develop a robust machine learning model for forecasting daily rainfall of East Coast of Peninsular Malaysia (ECPM). The ECPM is considered as the most vulnerable region of the peninsula to different kinds of rainfall-driven hazards [32]. Hence, a reliable rainfall pattern forecasting is highly essential for climate sustainable development of the region. Across large number of studies have been conducted for rainfall forecasting in Peninsular Malaysia, only one attempt was made for forecasting rainfall in Peninsular Malaysia (PM) using physical-empirical models [22]. However, the previous study of Pour et al. [1] was limited to the seasonal forecasting only. To the best knowledge of the current study, this is the first attempt to use an array of sophisticated ML algorithms for the development of physical-empirical models for forecasting daily rainfall of PM.

2. East Coast of Peninsular Malaysia (ECPM)

The ECPM (Figure 1) is the most interesting region of PM in terms of rainfall and rainfall-driven hydrological hazards [33]. The area has an undulating topography which varies from 4.0 m in the coastline to 2270.0 m in the interior mountainous region. A considerable part of the interior of the ECPM is mountainous, which influences the climate of the region. The region is characterized by high uniform temperature, high humidity, and copious rainfall [34]. The area receives a significant amount of rainfall even in the driest month. The study area experiences two monsoon seasons, namely, the northeast (NE) monsoon from the mid of October to end of February and the southwest (SW) monsoon from April to September [35]. The months of March and October are intermonsoonal transitional periods. The eastern coastal region is wetter compared with other parts of PM. It receives an annual average rainfall of 2800 mm. Monthly variations in rainfall in ECPM are shown in Figure 2. Heavy rainfall in the region occurs during NE monsoon. Maximum rainfall in a year is usually recorded in November or December. Some parts of the study area suffer flooding at this time of year [36]. On the other hand, the longest dry spell in a year is observed during SW monsoon, particularly in June and July.

The mean temperature in ECPM is found more or less uniform throughout the year. Seasonal variation of mean temperature (27.0°C) is always less than 2.0°C. The daily maximum temperature in the study area varies from 29.0°C in January to 33.0°C in April. The minimum temperature varies from 23.0°C in January-February to 25°C in May [37].

Several rivers originate in the mountainous interior region of ECPM. Due to high variation in topography and a high amount of rainfall, the region has a dense river network.

3. Description of the Data Used

The APHRDITE gridded dataset was used as a proxy to observed data. APHRDITE rainfall was developed using gauge precipitation data obtained from the Global Telecommunication System (GTS) network and other hydro-meteorological in situ records [38]. It has been developed based on a new interpolation technique with accurate long-term gridded orographic precipitation for Asia [38]. The data are available at http://aphrodite.st.hirosakiu.ac.jp/product/APHRO_V1101EX_R1/APHRO_MA/025deg_nc/.

The recently released NCEP/NCAR ERA2 reanalysis data of SLP with a resolution of $0.25^\circ \times 0.25^\circ$ were used as predictors for the prediction of rainfall. The NCEP/NCAR reanalysis provides datasets of various atmospheric variables at 17 different pressure levels [39]. The NCEP dataset has been utilized for the development of prediction models [40]. Rainfall in ECPM mostly occurs due to moist air from the ocean. The wind circulation in the region depends on SLP. Therefore, it was considered in the present study for the development of the rainfall prediction model. The climate domain consisted of 2730 grid points for the extraction of probable predictors. The climate domain and the grid points from where SLP data were collected are shown in Figure 3. The NCEP SLP data are available at <https://www.ncdc.noaa.gov/data-access/model-data/model-datasets/reanalysis-1-reanalysis-2>.

APHRDITE rainfall data are available for the period 1951–2015 while the NCEP ERA2 reanalysis SLP data are available for the period 1948 to present. Therefore, a common period of 1951–2015 was considered in the present study for model development and validation. APHRDITE daily rainfall and NCEP ERA daily SLP data for the period 1951–2015 were collected from the corresponding websites.

4. Description of Employed Machine Learning Algorithms

Five advanced AI algorithms such as ELM, BRNN, BART, EGB, and HNFIS were used considering the complex relationship of rainfall with SLP. The SLPs over the climate domain having a significant correlation with daily rainfall of ECPM were used to compute their principal components and selection of inputs. The APHRDITE rainfall of all the grid points over the ECPM was average to prepare the daily rainfall of the area. The ML models were used for the prediction of areal averaged daily rainfall of the region. The description of the ML models used in this study is given in the following subsections.

4.1. Extreme Learning Machine (ELM). The ELM model was developed following the structure of a feedforward neural network (FNN) [41]. Basic difference of ELM and FFN is that ELM uses the Moore–Penrose generalized inverse algorithm

to calculate the weights of hidden neurons in contrast to gradient algorithms of FFN [42]. A general structure of ELM is shown in Figure 4.

The algorithm of ELM is very simple. It takes the predictors as inputs (x) and multiplies them with the estimated weights (w), adds bias (b), applies activation function (g), and repeats the steps for a time equal to the number of layers used to get the output as follows:

$$f_L(x) = \sum_{i=1}^L \beta_i g(w_i * x_i + b_i), \quad j = 1, 2, \dots, N, \quad (1)$$

in which L represents the number of hidden layers, N is training data size, and β is the weight multiplied with the output of hidden layer to get the output, which is calculated as the ratio of hidden layer output matrix H and target matrix T . The H and T can be expressed as follows:

$$T = \begin{bmatrix} t_1^T \\ \vdots \\ t_N^T \end{bmatrix}_{N \times 1}, \quad (2)$$

$$H = \begin{bmatrix} g(w_1 * x_1 + b_1) & \dots & g(w_L * x_i + b_L) \\ \vdots & \dots & \vdots \\ g(w_1 * x_N + b_i) & \dots & g(w_L * x_N + b_L) \end{bmatrix}_{N \times L}.$$

The obtained output using equation (1) is back-propagated to the network to repeat the process until the desired level of error in prediction achieved. ELM has a good nonlinear adaptive capability to overcome the limitations of many other ML algorithms [43].

4.2. Bayesian Regularized Neural Network (BRNN). BRNN is a version of an artificial neural network (ANN), which is much robust compared with conventional ANN [44]. The robustness in BRNN is achieved due to Bayesian regularization of ANN parameters. A common error function (E_D) of ANN using early stopping can be expressed as follows:

$$E_D(D | w, M) = \sum_{i=1}^n (t_i - \hat{t}_i)^2, \quad (3)$$

where w is the weight, M is the ANN structure, n is the training data size, and t_i is i -th target while \hat{t}_i is the output.

Immature convergence of ANN causes overfitting of the model. A regularization of ANN using the Bayesian method helps optimization of ANN parameters using prior values of ANN parameters. For this purpose, an additional term (E_w) is included in objective function in BRNN as follows:

$$E_D(D | w, M) = \sum_{i=1}^n (t_i - \hat{t}_i)^2 + E_w, \quad (4)$$

where E_w is used to penalize the unrealistic weights to have a better generalization and gradual conversion. A gradient-based optimization method is used to minimize the function:

$$F = \beta E_D(D | w, M) + \alpha E_w(w | M), \quad (5)$$

in which $E_w(w | M)$ is the sum of the square of ANN architecture and α and β represent the hyperparameters to be optimized.

BRNN can reveal theoretically complex input-output relationship and thus considered as an efficient predictive model [28, 29].

4.3. Bayesian Additive Regression Trees (BART). BART is an ensemble of easily adjustable regression models [26]. It uses the basic structure of random forest (RF) [45], but the outcomes of the trees are estimated using a Bayesian inferential system. This allows estimation of uncertainty and regularization of model parameters to enhance the capability of the conventional tree-based regression method.

The procedure used by BART to improve the performance of RF is like the gradient boosting (GB) where many weak models are used to learn the problem [46]. In RF, the outcomes of all trees are averaged whereas in GB, the outcomes are adjusted using constants [47]. The novelty of BART is the use of Bayesian framing. In BART, priors are used to achieve the posterior of the outcomes of the trees. This provides a higher performance of the model. Besides, a sophisticated penalty system is used which helps in parameter regularization efficiently [26].

4.4. Extreme Gradient Boosting (xgBoost). Gradient boosting (GB) integrates the outcomes of many weak models for prediction [27]. This improves or boosts the prediction capability of weak models, and thus it is called GB. The xgBoost is an advanced version of GB where model parameters are regularized properly to enhance its performance [48]. The major advantage of xgBoost is the combination of both high speed and good efficiency. It is generally found to outperform many conventional models when the prediction is done using large datasets [49].

The basic structure of xgBoost is like a random forest where multiple decision trees are generated using data samples following either bagging or boosting method [50]. However, the main difference is in the learning process. In xgBoost, the model tries to learn from the mistake in the previous step. It means, in each step, it focuses on the part which is hard to learn or the model failed to predict properly. It tries to correct the mistake in the previous step by an attempt to solve the hard part of problem. A simple diagram of xgBoost is shown in Figure 5. An ensemble of models is used in this process where different models are used to solve different parts. The more error by a model indicates that the model is weak or the problem is hard enough to solve by this model. The xgBoost tried to improve this weak model to improve the total prediction capability of the model. In each step, it estimates the error and focuses on the part which generates a higher error and try to solve the problem until the desired level of error is achieved [51].

4.5. Hybrid Neuro Fuzzy Inference System (HNFIS). The HNFIS is developed by hybridization of the fuzzy inference system (FIS) with artificial neural network (ANN) [30], where the parameters of FIS were optimized using ANN which enables the system to learn the problem systematically and helps to enhance the predictive performance of the model. ANN has a good capacity to learn automatically. However, it cannot acquire outcomes efficiently. The FIS cannot learn automatically like ANN, but it can acquire the outcomes efficiently through fuzzy logic. Therefore, the hybridization of ANN with FIS can improve the learning capability of the nonlinear problem [52].

A typical HNFIS system is presented using the diagram in Figure 6. It is very similar to the adaptive neuro fuzzy inference system (ANFIS) where a multilayer feedforward ANN is the core [53]. The AND operation is employed first, followed by product, normalization, and sum operations, and finally centroid to get the output. However, the weights connecting the nodes are defined using fuzzy sets. The least-squares approach is used to decide the set of weights. The weights are adjusted using fuzzy rules to reduce model error [54].

4.6. Generation of Model Inputs. The SLP of all the NCEP grid points over the selected climate domain was correlated with the areal averaged daily APHRODITE rainfall over the ECPM. Spearman rank correlation (Spearman, 1904) was used considering high skewness in daily rainfall data [55]. The significance of the correlation was estimated at a significance level of 0.5 ($p < 0.05$). The SLP data which found significantly correlated with the areal average rainfall of ECPM were used to estimate their principal components (PCs). The principal component analysis (PCA) allows reduction of dimensionality in data [56]. The PCs which are found to represent the maximum variability in SLP data were used as input. In this study, the above procedure was conducted for SLPs for 1- to 3-day lags and the PCs of each lag were used as model input. The PCA has been found as an efficient method for the selection of inputs from a large dataset in previous studies [57, 58].

4.7. Development of Models and Assessment of Model Performance. APHRODITE daily rainfall and PCs of daily SLP data for the period 1951–2015 were used for the development of the model. There is no universal rule of using the proportion of data for training and validation. However, conventionally 70% of available data are used for model training. Following the generally adopted division rule, 70% of time series data (01 January 1951 to 30 June 1995) were used to train the models while the data for the period July 1995 to December 2015 were used for model validation.

The performance of the ML model depends on optimum values of model parameters. In this study, a k -fold validation was used for the optimization of ML model parameters, where k was considered as 10. Data were shuffled in each

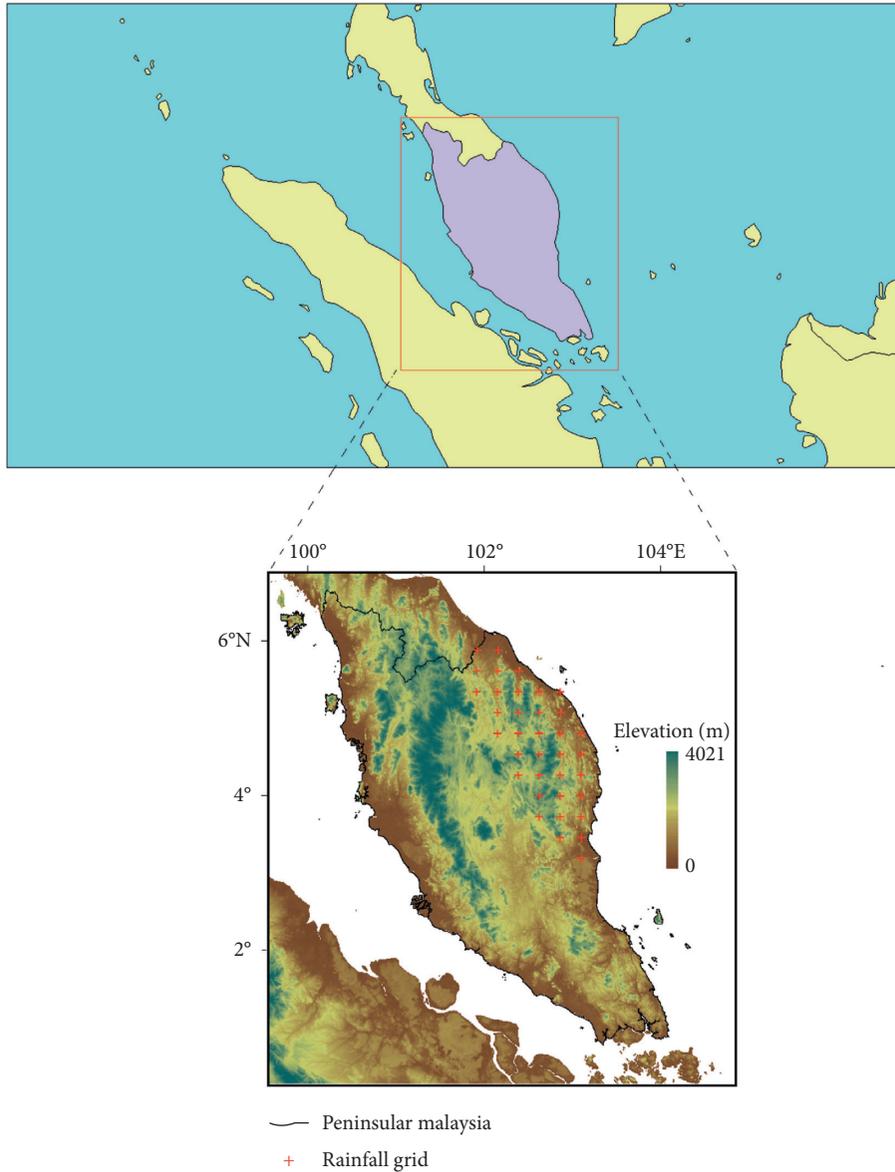


FIGURE 1: Location of Peninsular Malaysia in Southeast Asia. The elevation and location of grid points from where rainfall data were extracted.

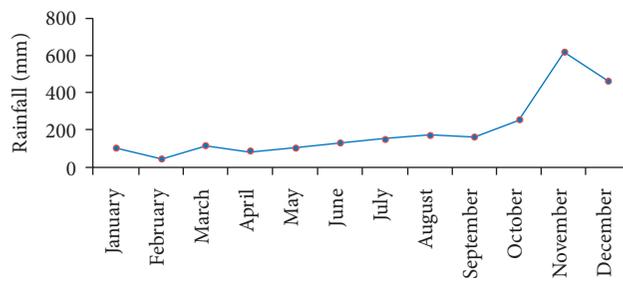


FIGURE 2: Mean monthly rainfall in the study area (1961–2000).

iteration so that all the data are used for model validation. The parameter values estimated in each iteration were averaged to get the optimum parameters. In this study, two parameters for ELM (number of hidden units and activation

function), three parameters for BRNN (α , β , and the number of neurons), five parameters for BART (number of trees, prior boundary, base terminal not, power terminal node, and degrees of freedom), four parameters for xgBoost (number

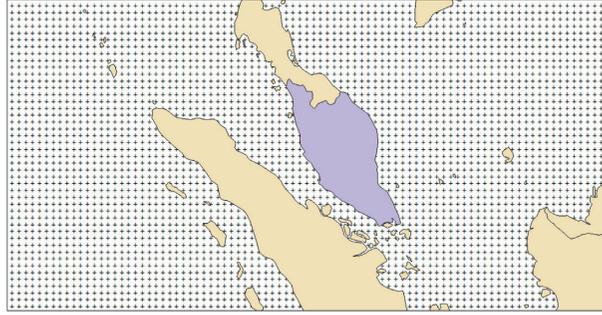


FIGURE 3: Climate domain considered for extraction of sea level pressure data for the prediction of rainfall in Peninsular Malaysia.

of boosting iterations, two regularization parameters, and learning rate), and two HNFIS parameters (number of fuzzy terms and maximum iterations) were optimized.

The *caret* package of statistical software *R* was used for this purpose. All the models were also developed using different packages available in *R*.

Two error metrics, namely, mean absolute error (MAE) and normalized root mean square error (NRMSE) and two association metrics such as Wilmott's modified index of agreement (WI) and Kling-Gupta efficiency (KGE) were used to assess the performance of the models statistically [59, 60]. The model performance was evaluated during the validation period to show the relative performance of the models in predicting an areal average of daily rainfall in ECPM with a lag time of one day. The mathematical formulas are as follows:

$$\begin{aligned} \text{MAE} &= 1/N \sum_{i=1}^N |y_f - y_o|, \\ \text{NRMSE} &= \frac{\sqrt{1/N \sum_{i=1}^N (y_f - y_o)^2}}{\bar{y}_o}, \\ \text{WI} &= 1 - \frac{\sum_{i=1}^N |y_f - y_o|}{\sum_{i=1}^N |y_f - \bar{y}_o| + |y_o - \bar{y}_o|}, \\ \text{KGE} &= 1 - \sqrt{(r-1)^2 + (\gamma-1)^2 + (\beta-1)^2}, \end{aligned} \quad (6)$$

where N is the number of observations, y_f is the forecasted values, y_o is the observed values, \bar{y}_f and \bar{y}_o are the mean values of the forecasted and observed, β is the bias, γ is the fraction of the coefficient of variation, and r is the linear correlation.

5. Results and Discussion

5.1. Analysis of SLP Data for Generation of Inputs. The areal average rainfall of ECPM was first correlated with different lags of the SLP data at different NCEP grid points over the climate domain. Obtained results are presented in Figure 7. The correlation of rainfall with 1- to 3-lag days is presented in three maps in the figure. The colour ramp in the maps is used to show the positive (negative) correlation. The significant correlation ($p < 0.05$) is presented with dots. The figure shows that rainfall in ECPM is positively correlated

with SLP in the north and negatively correlated with SLP in the south. The differences in SLP cause movement of air over the region. The moist air from the sea when enters the land causes rainfall. Therefore, SLP data can be used for prediction of rainfall in ECPM. The SLP data of nearly 2321 to 2486 are found to correlate with rainfall for different lags. Such a big number of variables cannot be used for the development of the model. Therefore, PCA was used to reduce the dimensionality of data.

The significantly correlated SLP data were used to estimate their PCs using PCA. Obtained PCs for one-day lag SLP are shown in Figure 8. The scree plot of the first 10 PCs for one-day lag SLP is shown in Figure 8(a). The figure shows that the first PC explains most of the variance compared with other PCs. The eigenvalue at 1 is shown using a red dashed horizontal line in the figure. It can be seen from the figure that the first 5 PCs are above the red line. The cumulative variance plot of the first 10 PCs is shown in Figure 8(b). The vertical blue dashed line represents the cutoff line for PC-5. The results showed that the first 5 PCs represent 99% of the variance in the SLP data over the climate domain. The first PC covered 91.7% of variability, the first 2 PCs covered 95.4% of variability, the first 3 PCs covered 98.3% of the variability, and so on. The results indicated increasing of PCs after the first PC does not increase the explained variability by the PCs significantly. Similar results were obtained for SLP with 2- and 3-day lags. The first PC along was able to explain more than 91% of the variability in data for all the three lags. The higher number of PCs as inputs increases the complexity of the model. Therefore, to make a parsimonious model, only the first PC of each lag was selected. These three PCs were finally used as input to train the models.

5.2. Comparison of Model Performance Using Statistical Metrics. The performance of the models was first evaluated using statistical metrics. Obtained results are given in Table 1. Only two error and association metrics were used as they can assess the capability of the models in replicating all the aspects including bias in mean and variability and association in rainfall amount and volume. For a better comparison of the relative performance of the model based on all the metrics together, the results are also presented using a radar chart in Figure 9.

The results showed that none of the models was able to attain very high metric values such as KGE or WI above 0.8.

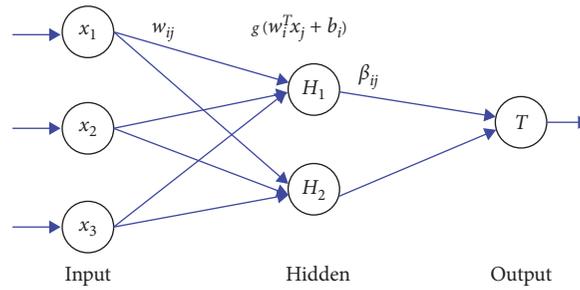


FIGURE 4: Schematic diagram showing the basic structure of an extreme learning machine.

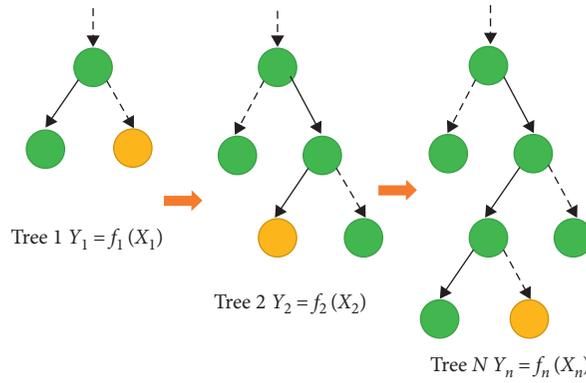


FIGURE 5: Schematic diagram explaining the gradient boosting algorithm.

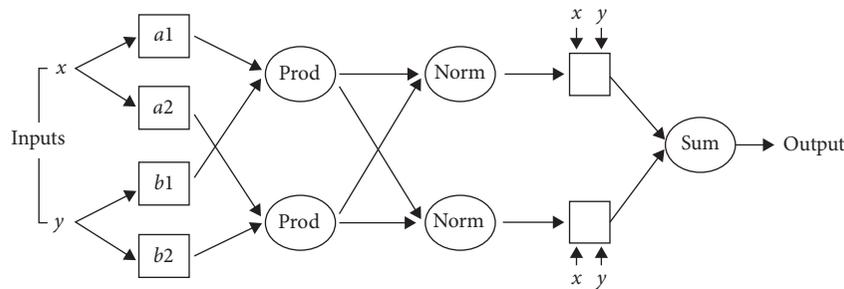


FIGURE 6: The basic structure of a hybrid neuro fuzzy inference system.

However, metric values obtained by the models should be judged considering the amount of sample size of data. The obtained statistics were estimated for the validation period July 1995 to December 2015 which means the sample size of 7120. Computation of any statistical metrics depends on sample size (n), where metric value decreases with the increase in n , even when the model behave the same. Therefore, the KGE or WI values more than 0.5 should be considered very good for this high amount of data size. This can be justified with a small amount of error in terms of MAE and NRMSE.

Overall, all the models were found to perform similarly with a slightly better or less performance in terms of different indices. However, BRNN was found to perform best in term of all metrics. Figure 9 shows lower values of BRNN in terms of error metrics and high values in terms of association-

based metrics. The HNFIS can be considered as the second-best model followed by BART and xgBoost, while ELM showed the least performance.

It is a well-established fact that the performance of ML models depends on the data used. There is no universal ML model which shows good performance for all kinds of data. It is not possible to justify the better performance of BRNN and the least performance of ELM. Only it can be remarked that the present study revealed that BRNN is more suitable for prediction of rainfall in ECPM using SLP data over the climate domain of PM.

5.3. Comparison of Model Performance through Graphical Presentations. The performance of the five ML models used in this study was evaluated using graphical presentation of

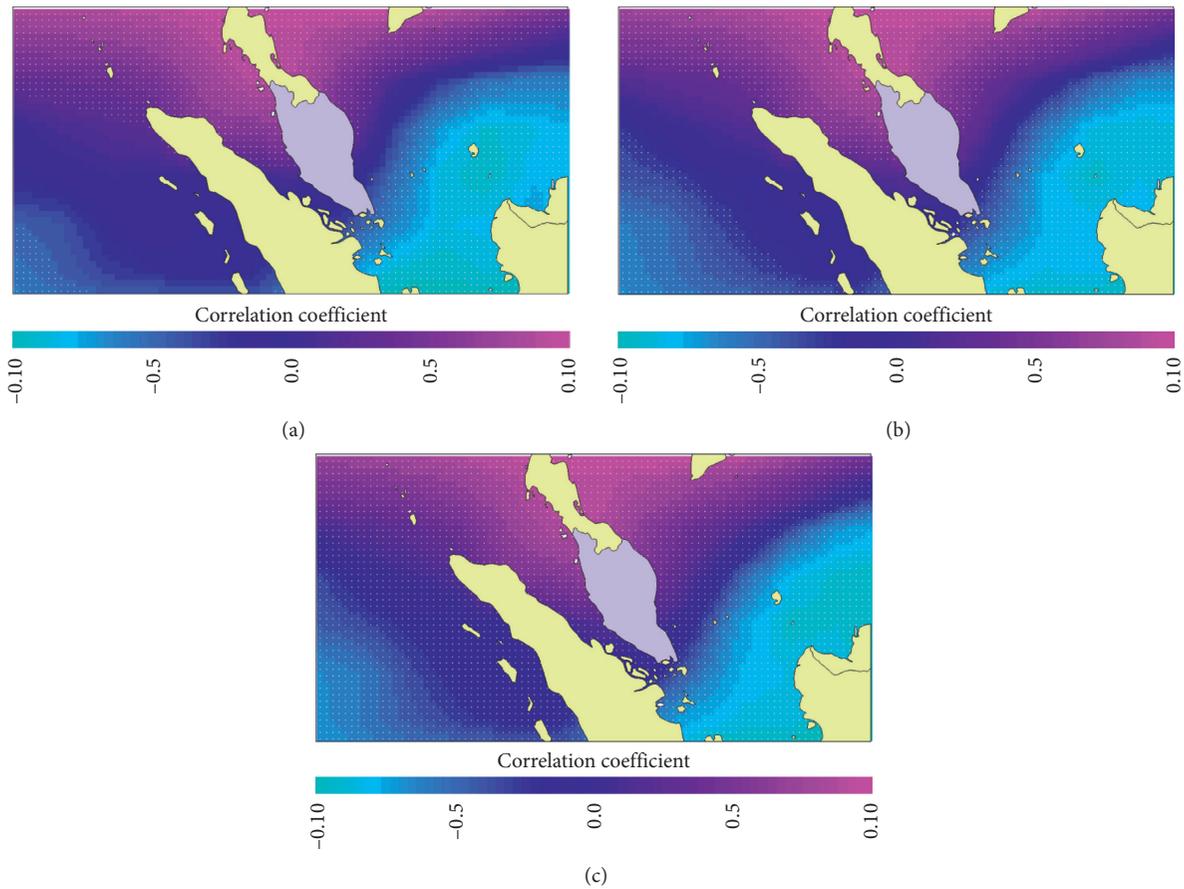


FIGURE 7: Correlation of areal average rainfall in the East Coast of Peninsular Malaysia with sea level pressure at different locations for (a) 1-day; (b) 2-day, and (c) 3-day lags. The colour ramp is used to show the positive (negative) correlation. The significant correlation ($p < 0.05$) is presented using white dots.

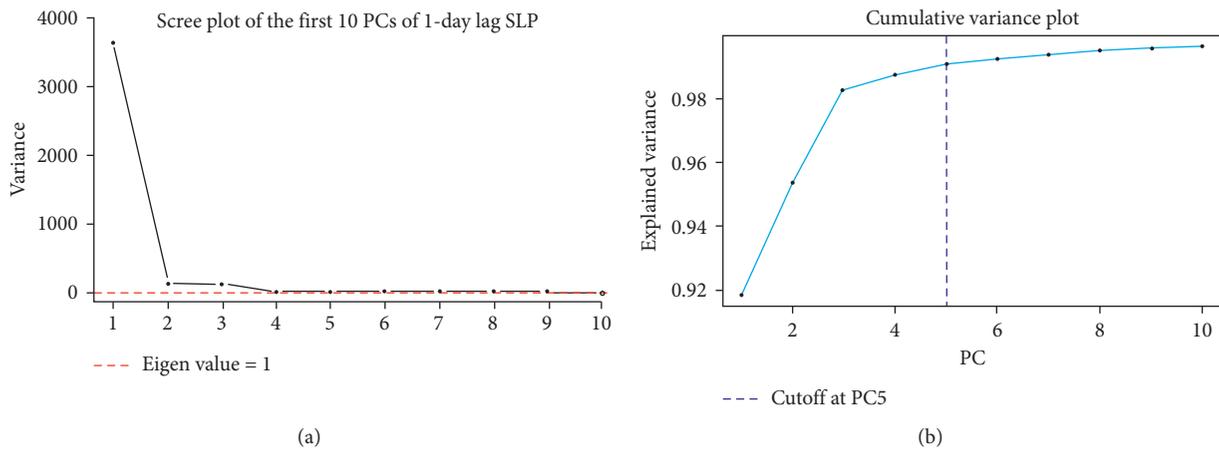


FIGURE 8: (a) The scree plot and (b) cumulative variance plot of the first ten principal components of one-day lag sea level pressure data over the climate domain of Peninsular Malaysia.

observed and model data in three different ways. Considering the large volume of data, different innovative approaches were employed for the presentation of data.

First, the density-scatter plots were prepared. This type of plot is suitable for comparison of a large amount of data.

As the data size was 7120, it was not possible to compare effectively using a simple scatter plot. Therefore, density-scatter plots are used in this study. Obtained results are presented in Figure 10. The red colour in the plot represents a higher density of data, while the blue represents a lower

TABLE 1: Performance of the ML models in terms of two error and two association-based statistical metrics. The bold number indicates the best performance.

Model	MAE	NRMSE	WI	KGE
ELM	0.045	0.915	0.53	0.53
BRNN	0.035	0.678	0.62	0.69
BART	0.039	0.784	0.58	0.64
xgBoost	0.039	0.803	0.58	0.63
HNFIS	0.036	0.708	0.61	0.68

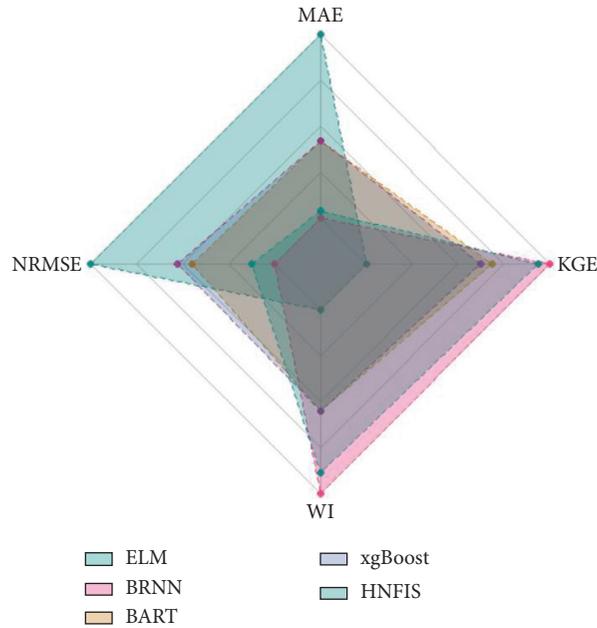


FIGURE 9: Radar chart showing the relative performance of five ML models considered in the present study based on four statistical metrics.

density of data. As the days with zero or low rainfall is very high compared with days with higher rainfall, the red zones (high density of data) were noticed only for the lower rainfall. If the diagonal line of the plot bisects the red and other higher rainfall zones symmetrically, the model is considered better. The results presented in Figure 10 reveal that high-density data were more above the bisecting line for all the models which indicates over prediction of low rainfall values by all the models. However, the BRNN was most perfect followed by HNFIS and BART in symmetrically bisecting the high-density rainfall data. The results also showed that none of the models was capable to simulate the high rainfall of values. All the model underpredicted the high rainfall amount. Overall, the underprediction was less for BRNN followed by BART. Daily rainfall data are highly skewed. It follows a gamma distribution. The number of no rainfall days is very high compared with days having different amounts of rainfall. Again, the days with low rainfall values are always very high compared with the days having high rainfall values. Another major problem is the outliers in rainfall data. Among thousands of rainfall data, only a few are very high and exceed the range of three standard deviations from the mean. Prediction of such outliers by training a model with all rainfall data is still a challenge. Overprediction of low rainfall values and underprediction of

high rainfall values in modeling daily or hourly rainfall are still a subject of investigation for both physical and empirical modelers.

The capability of the models in estimating rainfall of all the months and the seasonal variability was evaluated using side-by-side bar plots. Obtained results for the five models are presented in Figure 11. The results revealed that all the models were able to replicate the seasonal variability in rainfall (high rainfall during November-December and low rainfall during June-August). However, it was also noticed that all the models overpredicted rainfall during low rainfall months and underpredicted rainfall during high rainfall months. Almost no difference among the models was noticed based on the monthly presentation of data. However, close observation of results revealed a bit less under-prediction and overprediction by BRNN.

Finally, the performance of the models was visually compared by plotting the first two PCs of the data. For this purpose, the PCs of observed and model outputs were estimated. The PCs were then plotted (first PC in the x -axis and the second PC in the y -axis). This provides a visualization of how the data are related to each other in PCs. It gives an assessment of which samples are similar and which are not. Obtained results are presented in Figure 12. The results showed similarity in results of BRNN, xgBoost, and

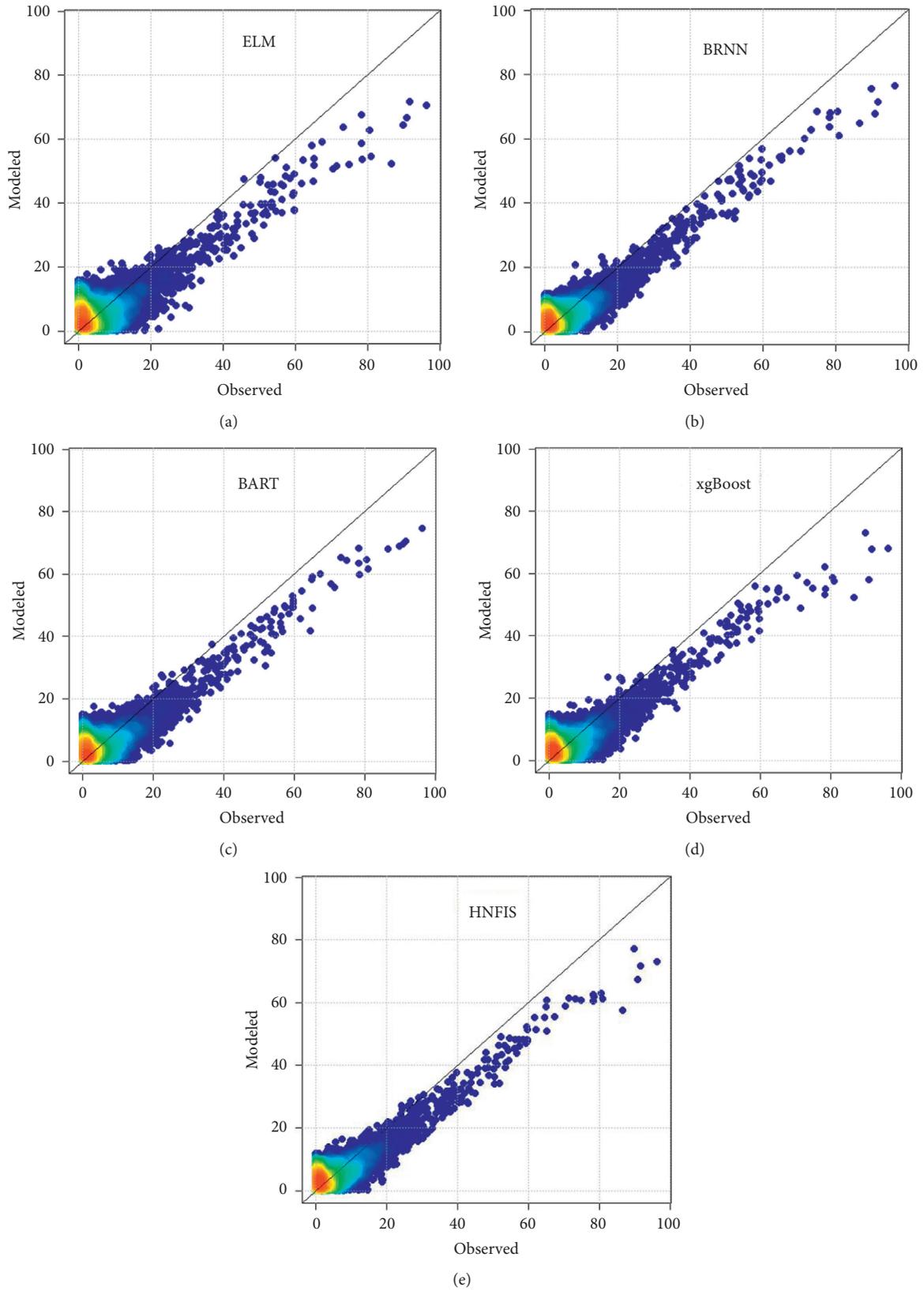


FIGURE 10: Comparison of observed and model simulated rainfall using density-scatter plot.

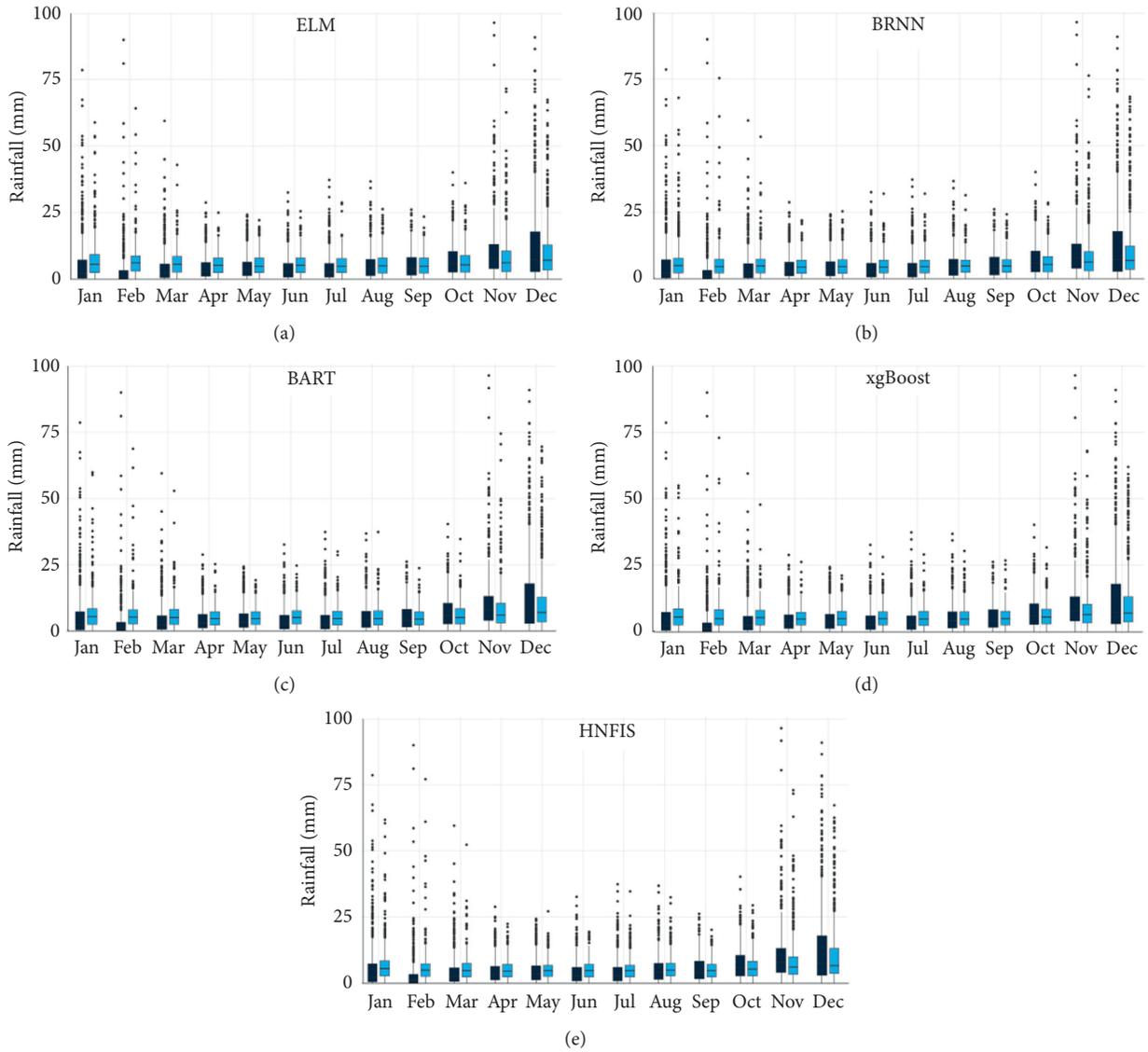


FIGURE 11: Side-by-side presentation of whisker-box plot of observed and simulated rainfall by the model showing the capability of the models in simulating rainfall for each month and replicate the seasonal variability of rainfall in the study area.

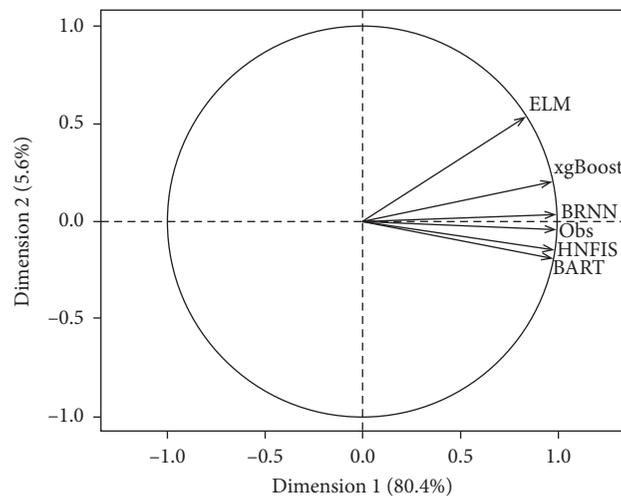


FIGURE 12: Assessment of similarity of model outputs with the observed rainfall-based variability plots of principal components of observed and simulated rainfall.

ELM while the results of HNFIS and BART were different from the other three model outputs. However, when compared with the observed data, the BRNN was found closer to the observed line in the plot followed by HNFIS and BART. The ELM is found to locate most apart from the observed line and also from the other models. The results indicate the best performance of BRNN in replicating observed rainfall followed by HNFIS and BART.

6. Conclusion

Five advanced AI algorithms, namely, ELM, BRNN, BART, xgBoost, and HNFIS were employed for the developed of the physical-empirical models for the prediction of rainfall in the eastern coastal region of PM from SLP data collected from the climate domain of PM. The parameters of the models were optimized using a k -fold validation approach. The performance of the models was evaluated using various innovative visuals presented along with statistical metrics. The results showed better performance of BRNN compared with other models in predicting daily rainfall one-day ahead. Prediction of daily rainfall is an extremely challenging task for the tropical region where daily rainfall data can be compared with a high fluctuating extremely noisy data without any exportable temporal pattern. The capability of the physical-empirical model of incorporation of physical mechanism and use of advanced ML algorithms has made the models developed in this study efficient in rainfall prediction. The Bayesian regularization of neural network parameters makes it more capable to reveal theoretically complex input-output relationship. It might be concluded that the complexity due to a very chaotic relationship of SLP and rainfall in the ECPM was able to capture through Bayesian regularization of the ANN model. The results revealed the potential of the model to be employed for the development of early warning of rainfall. However, it should be noted that the model was not capable to replicate very high rainfall which means it cannot be employed for prediction of extreme rainfall and probable flash flood. Besides, the model overpredicted low rainfall and thus more rainfall days which restricted this applicability in predicting dry spells. In the future, more attention should be given to improve the performance of the models in simulating extreme rainfall for expanding their applicability in prediction of extreme rainfall.

Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

Conflicts of Interest

The authors declare that they have no conflicts of interest.

Acknowledgments

This research was funded by the Ministry of Higher Education Malaysia under the Fundamental Research Grant Scheme, (203.PHUMANITI.6711695) and the USM

Publication Fund. The first author would like to acknowledge the support from the Science and Technology Plan Project of Shaanxi Province grant no. 2020GY-041 and Doctoral Research Initiation Project grant no. 209040080.

References

- [1] S. H. Pour, A. K. A. Wahab, and S. Shahid, "Spatiotemporal changes in aridity and the shift of drylands in Iran," *Atmospheric Research*, vol. 233, 2020.
- [2] D. B. Wright, R. Mantilla, and C. D. Peters-Lidard, "A remote sensing-based tool for assessing rainfall-driven hazards," *Environmental Modelling & Software*, vol. 90, 2017.
- [3] I. Belachsen, F. Marra, N. Peleg, and E. Morin, "Convective rainfall in a dry climate: relations with synoptic systems and flash-flood generation in the Dead Sea region," *Hydrology and Earth System Sciences*, vol. 21, 2017.
- [4] S. S. B. Brito, A. P. M. A. Cunha, C. C. Cunningham, R. C. Alvalá, J. A. Marengo, and M. A. Carvalho, "Frequency, duration and severity of drought in the Semiarid Northeast Brazil region," *International Journal of Climatology*, vol. 38, 2018.
- [5] V. Joshi and K. Kumar, "Extreme rainfall events and associated natural hazards in Alaknanda valley, Indian Himalayan region," *Journal of Mountain Science*, vol. 3, 2006.
- [6] H. R. Pourghasemi, A. Gayen, M. Edalat, M. Zarafshar, and J. P. Tiefenbacher, "Is multi-hazard mapping effective in assessing natural hazards and integrated watershed management?" *Geoscience Frontiers*, vol. 11, 2020.
- [7] B. Praveen, S. Talukdar, Shahfahad et al., "Analyzing trend and forecasting of rainfall changes in India using non-parametrical and machine learning approaches," *Scientific Reports*, vol. 10, 2020.
- [8] M. Zeynoddin, H. Bonakdari, A. Azari, I. Ebtehaj, B. Gharabaghi, and H. Riahi Madavar, "Novel hybrid linear stochastic with non-linear extreme learning machine methods for forecasting monthly rainfall a tropical climate," *Journal of Environmental Management*, vol. 222, 2018.
- [9] S. Y. Yim, B. Wang, W. Xing, and M. M. Lu, "Prediction of Meiyu rainfall in Taiwan by multi-lead physical-empirical models," *Climate Dynamics*, vol. 44, 2015.
- [10] Y. Luo, J. Sun, Y. Li et al., "Science and prediction of heavy rainfall over China: research progress since the reform and opening-up of new China," *Journal of Meteorological Research*, vol. 34, 2020.
- [11] X. Yang, R. He, J. Ye et al., "Integrating an hourly weather generator with an hourly rainfall SWAT model for climate change impact assessment in the Ru River Basin, China," *Atmospheric Research*, vol. 244, p. 105062, 2020.
- [12] E. Toth, A. Brath, and A. Montanari, "Comparison of short-term rainfall prediction models for real-time flood forecasting," *Journal of Hydrology*, vol. 239, no. 1-4, pp. 132-147, 2000.
- [13] C. L. Wu and K. W. Chau, "Prediction of rainfall time series using modular soft computing methods," *Engineering Applications of Artificial Intelligence*, vol. 26, no. 3, pp. 997-1007, 2013.
- [14] P. Satyanarayana and V. V. Srinivas, "Regional frequency analysis of precipitation using large-scale atmospheric variables," *Journal of Geophysical Research*, vol. 113, 2008.
- [15] S. Q. Salih, A. Sharafati, I. Ebtehaj et al., "Integrative stochastic model standardization with genetic algorithm for rainfall pattern forecasting in tropical and semi-arid environments," *Hydrological Sciences Journal*, vol. 65, no. 1-13, 2020.

- [16] P. Deb, A. S. Kiem, and G. Willgoose, "A linked surface water-groundwater modelling approach to more realistically simulate rainfall-runoff non-stationarity in semi-arid regions," *Journal of Hydrology*, vol. 575, 2019.
- [17] A. Banerjee, R. Chen, M. E. Meadows, R. B. Singh, S. Mal, and D. Sengupta, "An analysis of long-term rainfall trends and variability in the uttarakhand himalaya using google earth engine," *Remote Sensing*, vol. 12, 2020.
- [18] S. K. Jain, P. Mani, S. K. Jain et al., "A Brief review of flood forecasting techniques and their applications," *International Journal of River Basin Management*, vol. 16, no. 3, pp. 329–344, 2018.
- [19] S. Segoni, L. Piciullo, and S. L. Gariano, "A review of the recent literature on rainfall thresholds for landslide occurrence," *Landslides*, vol. 15, 2018.
- [20] S. Y. Yim, B. Wang, and W. Xing, "Prediction of early summer rainfall over south China by a physical-empirical model," *Climate Dynamics*, vol. 47, 2014.
- [21] P. Chen and B. Sun, "Improving the dynamical seasonal prediction of western Pacific warm pool sea surface temperatures using a physical-empirical model," *International Journal of Climatology*, vol. 40, 2020.
- [22] S. H. Pour, A. K. A. Wahab, and S. Shahid, "Physical-empirical models for prediction of seasonal rainfall extremes of Peninsular Malaysia," *Atmospheric Research*, vol. 233, p. 104720, 2020.
- [23] T. Hai, A. Sharafati, A. Mohammed et al., "Global solar radiation estimation and climatic variability analysis using extreme learning machine based predictive model," *IEEE Access*, vol. 8, pp. 12026–12042, 2020.
- [24] J. Yuval and P. A. O’Gorman, "Stable machine-learning parameterization of subgrid processes for climate modeling at a range of resolutions," *Nature Communications*, vol. 11, 2020.
- [25] X. Chen, J. Parajka, B. Széles, P. Strauss, and G. Blöschl, "Spatial and temporal variability of event runoff characteristics in a small agricultural catchment," *Hydrological Sciences Journal=Journal des Sciences Hydrologiques*, vol. 65, no. 13, 2020.
- [26] H. A. Chipman, E. I. George, and R. E. McCulloch, "BART: Bayesian additive regression trees," *The Annals of Applied Statistics*, vol. 4, no. 1, pp. 266–298, 2010.
- [27] R. Bekkerman, M. Bilenko, and J. Langford, *Scaling Up Machine Learning: Introduction*, pp. 1–20, Cambridge University Press, Cambridge, UK, 2012.
- [28] H. Okut, "Bayesian regularized neural networks for small n big p data," *Artificial Neural Networks-Models and Applications*, InTech, London, UK, 2016.
- [29] M. Kayri, "Predictive abilities of bayesian regularization and levenberg-marquardt algorithms in artificial neural networks: a comparative empirical study on social data," *Mathematical and Computational Applications*, vol. 21, no. 2, p. 20, 2016.
- [30] S. Supatmi, R. Hou, and I. D. Sumitra, "Study of hybrid neurofuzzy inference system for forecasting flood event vulnerability in Indonesia," *Computational Intelligence and Neuroscience*, vol. 2019, Article ID 6203510, 13 pages, 2019.
- [31] A. Dormishi, M. Ataei, R. Khaloo Kakaie, R. Mikaeil, and S. Shaffiee Haghshenas, "Performance evaluation of gang saw using hybrid ANFIS-DE and hybrid ANFIS-PSO algorithms," *Journal of Mining and Environment*, vol. 10, 2018.
- [32] L. Juneng, F. T. Tangang, and C. J. C. Reason, "Numerical case study of an extreme rainfall event during 9–11 December 2004 over the east coast of Peninsular Malaysia," *Meteorology and Atmospheric Physics*, vol. 98, 2007.
- [33] O. O. Mayowa, S. H. Pour, S. Shahid et al., "Trends in rainfall and rainfall-related extremes in the east coast of peninsular Malaysia," *Journal of Earth System Science*, vol. 124, no. 8, pp. 1609–1622, 2015.
- [34] N. Khan, S. H. Pour, S. Shahid et al., "Spatial distribution of secular trends in rainfall indices of peninsular Malaysia in the presence of long-term persistence," *Meteorological Applications*, vol. 26, 2019.
- [35] M. F. Mohd Akhir, N. Z. Zakaria, and F. Tangang, "Intermonsoon variation of physical characteristics and current circulation along the east coast of peninsular Malaysia," *International Journal of Oceanography*, vol. 2014, Article ID 527587, 9 pages, 2014.
- [36] H. Brammer, "Floods in Bangladesh: II. Flood mitigation and environmental aspects," *The Geographical Journal*, vol. 156, 1990.
- [37] C. L. Wong, R. Venneker, S. Uhlenbrook, A. B. M. Jamil, and Y. Zhou, "Variability of rainfall in peninsular Malaysia," *Hydrology and Earth System Sciences*, vol. 183, 2009.
- [38] A. Yatagai, O. Arakawa, K. Kamiguchi, H. Kawamoto, M. I. Nodzu, and A. Hamada, "A 44-year daily gridded precipitation dataset for Asia based on a dense network of rain gauges," *SOLA*, vol. 5, pp. 137–140, 2009.
- [39] E. Kalnay, M. Kanamitsu, R. Kistler et al., "The NCEP/NCAR 40-year reanalysis Project," *Bulletin of the American Meteorological Society*, vol. 77, no. 3, pp. 437–471, 1996.
- [40] C. Folberth, A. Baklanov, J. Balković, R. Skalský, N. Khabarov, and M. Obersteiner, "Spatio-temporal downscaling of gridded crop model yield estimates based on machine learning," *Agricultural and Forest Meteorology*, vol. 264, pp. 1–15, 2019.
- [41] G.-B. Huang, Q.-Y. Zhu, and C.-K. Siew, "Extreme learning machine: theory and applications," *Neurocomputing*, vol. 70, no. 1–3, pp. 489–501, 2006.
- [42] G.-B. Huang and C.-K. Siew, "Extreme learning machine: RBF network case," in *Proceedings of the ICARCV 2004 8th Control, Automation, Robotics and Vision Conference*, vol. 2, pp. 1029–1036, Kunming, China, December 2004.
- [43] J.-M. Park and J.-H. Kim, "Online recurrent extreme learning machine and its application to time-series prediction," in *Proceedings of the 2017 International Joint Conference On Neural Networks (IJCNN)*, IEEE, Anchorage, AK, USA, May 2017.
- [44] F. Burden and D. Winkler, "Bayesian regularization of neural networks," *Methods In Molecular Biology™*, pp. 23–42, Humana Press, Totowa, NJ, USA, 2008.
- [45] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [46] S. Abu-Nimeh, D. Nappa, X. Wang, and S. Nair, "A distributed architecture for phishing detection using Bayesian additive regression trees," in *Proceedings of the 2008 eCrime Researchers Summit*, IEEE, Atlanta, GA, USA, October 2008.
- [47] S. Zhang, Y.-C. T. Shih, and P. Müller, "A spatially-adjusted Bayesian additive regression tree model to merge two datasets," *Bayesian Analysis*, vol. 2, no. 3, pp. 611–633, 2007.
- [48] H. Zhang, D. Qiu, R. Wu, Y. Deng, D. Ji, and T. Li, "Novel framework for image attribute annotation with gene selection XGBoost algorithm and relative attribute model," *Applied Soft Computing*, vol. 80, 2019.
- [49] T. Chen, H. Li, Q. Yang, and Y. Yu, "General functional matrix factorization using gradient boosting," in *Proceedings of the 30th International Conference On Machine Learning, ICML 2013*, Atlanta, GA, USA, June 2013.
- [50] L. Breiman, "Random forrests," *Machine Learning*, vol. 45, pp. 5–32, 2001.

- [51] A. Asuncion and D. Newman, *UCI Machine Learning Repository*, UCI, Irvine, CA, USA, 2007.
- [52] M. Reza kazemi, A. Dashti, M. Asghari, and S. Shirazian, "H2-selective mixed matrix membranes modeling using ANFIS, PSO-ANFIS, GA-ANFIS," *International Journal of Hydrogen Energy*, vol. 42, no. 22, pp. 15211–15225, 2017.
- [53] K. Li, H. Su, and J. Chu, "Forecasting building energy consumption using neural networks and hybrid neuro-fuzzy system: a comparative study," *Energy and Buildings*, vol. 43, 2011.
- [54] E. Kim, M. Park, S. Ji, and M. Park, "A new approach to fuzzy modeling," *IEEE Transactions on Fuzzy Systems*, vol. 5, no. 3, pp. 328–337, 1997.
- [55] J. H. Zar, "Spearman rank correlation: overview," in *Wiley StatsRef: Statistics Reference Online* John Wiley & Sons, Hoboken, NJ, USA, 2014.
- [56] I. Jolliffe, "Principal component analysis," *International Encyclopedia Of Statistical Science*, pp. 1094–1096, Springer, Berlin, Germany, 2011.
- [57] U. Demšar, P. Harris, C. Brunson, A. S. Fotheringham, and S. McLoone, "Principal component analysis on spatial data: an overview," *Annals of the Association of American Geographers*, vol. 103, 2013.
- [58] Y. Liu, A. Singleton, and D. Arribas-Bel, "A Principal Component Analysis (PCA)-based framework for automated variable selection in geodemographic classification," *Geospatial Information Science*, vol. 22, no. 4, pp. 251–264, 2019.
- [59] B. Choubin, H. Darabi, O. Rahmati, F. Sajedi-Hosseini, and B. Kløve, "River suspended sediment modelling using the CART model: a comparative study of machine learning techniques," *Science of The Total Environment*, vol. 615, pp. 272–281, 2018.
- [60] M. S. Al-Musaylh, R. C. Deo, J. F. Adamowski, and Y. Li, "Short-term electricity demand forecasting with MARS, SVR and ARIMA models using aggregated demand data in Queensland, Australia," *Advanced Engineering Informatics*, vol. 35, 2018.