



**ISSN: 2321-2152**  
**IJMECE**  
*International Journal of modern  
electronics and communication engineering*

E-Mail  
[editor.ijmece@gmail.com](mailto:editor.ijmece@gmail.com)  
[editor@ijmece.com](mailto:editor@ijmece.com)

[www.ijmece.com](http://www.ijmece.com)

# DETECTING RAINFALL USING MACHINE LEARNING TECHNIQUES

<sup>1</sup>Y.Neha, <sup>2</sup>B.Tripura, <sup>3</sup>G.Karthik Reddy, <sup>4</sup>Challa Shiva, <sup>5</sup>J.S.Radhika,  
<sup>1,2,3,4</sup>U.G.Scholor, Department of IT, Sri Indu College Of Engineering & Technology, Ibrahimpatnam, Hyderabad.  
<sup>5</sup>Assistant Professor, Department of IT, Sri Indu College Of Engineering & Technology, Ibrahimpatnam, Hyderabad.

## ABSTRACT:

rainfall prediction is an important part of weather forecasting, especially when it comes to agricultural use, catastrophe control, and/or water resource explanation. In order to improve the efficacy of rainfall prevention, this research examines the use of ultrasound technology. Historical meteorological data is presented. Using non-local factors like temperature, humidity, speed of light, pressure in the atmosphere. The methods that are being suggested act as models, including spatial networks and regression analysis. Use these meteorological variables and rainfall patterns on large datasets to capture complicated relationships. By use of robust assessment and validation procedures. The models show that they can reliably predict predatory falls over short and long time periods. By using machine learning, this study increases weather forecasting. Machine learning has the ability to effectively address climate challenges through the use of predictive models.

Keywords: machine learning, Logistic Regression, Decision Tree, Random Forest, Rainfall prediction

## INTRODUCTION:

Recent arid accessible precipitation estimation is projected to achieve another time of mediation for businesses adversely affected by extreme precipitation events. These sectors include, but are not limited to, energy and agriculture, which are heavily influenced by precipitation patterns. Various scholarly investigation has revealed that both the length and force of precipitation contribute to massive environmental disasters. The effects of these precipitation-related factors include, among other things, droughts and floods. For example, in 2009, heavy rains impacted about 600,000 people in Senegal, Niger, Burkina Faso, and Ghana. Furthermore, floods in 2007 claimed the lives of about 1,000,000 people in Ethiopia, Uganda, Togo, Niger, Sudan, Mali, and Burkina Faso.

Furthermore, research suggests that the cost of hunger-related child death, which increased from 30,000 to 50,000 in Sub-Saharan Africa in 2009, may have been compounded by changes in precipitation patterns and extreme climatic events affecting the agricultural region. Aside from the impact of flooding due to floods, several investigations have documented the extensive effects of precipitation on basic sections of the Ghanaian economy. It is highlighted that a significant portion of the hydroelectric power plants, which account for more than 70% of Ghana's electricity generation, are heavily reliant on precipitation. This underscores the severe consequences for the country's power generation.

Farming, which employs around 44.7% of Ghana's labour population and plays a major role in the country's economy, is crucial for financial growth. Despite its current drop in execution, the horticulture sector remains an important component for poverty reduction and food security in Ghana. Regardless, Ghana's farming area is primarily rain-fed with roughly 3% of cultivable land supported by a water system. In the near future, precipitation is an important water source for farming, power generation, and other uses in non-industrial countries such as Ghana.

For precipitation forecasting, several operations such as Random Forest (RF), Decision Tree (DT), Neural Network (NN), K-Nearest Neighbor (KNN), and others have been investigated. The exhibition of these computations normally alters, allowing an opportunity for improvement by varying framing and testing ratios or combatting different time periods. Nonetheless, forecasting precipitation continues to be a difficult task. As a result, the careful selection of appropriate methodologies for describing precipitation patterns is crucial. A certain location is critical. AI computations have emerged as a viable alternative to improve the accuracy of precipitation forecasts. As a result, there has been a proliferation of precipitation forecast research using various approaches in various countries, including Malaysia, India, and Egypt, among others.

In one example, AI methodologies are used to create precipitation forecast models for major Australian cities. A number of calculations were performed to assess the performance of various approaches in various countries, including Decision Tree, Random Forest, Logistic Regression, AdaBoost, Gradient Boosting, and K-Nearest Neighbor.

## LATEDVOR] A :

The Papa [1] demonstrates that flood-related fatalities and economic losses have increased in Africa over the past 50 years, prompting the need to identify the causes for this increase. The study finds that extensive and unplanned human settlements in flood-prone areas are a major factor or increasing flood risk, and urgent actions such as discouraging settlements in these areas and implementing early warning systems are needed. Papa [2] focuses on the use of machine learning classification for rainfall prediction in Malaysia, comparing the performance of different techniques. The study concludes that the Neural Network (NN) classifier is the most effective for rainfall prediction in Malaysia. In paper [3], it finds distinct differences in rainfall characteristics and length of the rainy season across different climatic zones in Ghana. The forest and coastal zones have their rainfall onset in March, while the transition zone has its onset from March to April, and the savannah zone has its onset from April to May. The length of the rainy season varies across zones, with the forest zone having the longest rainy season. The paper [4] compares the onset, cessation, and duration of the rainy season in Ghana using simulated rainfall data from the Regional Climate Model (RegCM4) and range gauge measurements from the Ghana Meteorological Agency (GMA). The paper [5] highlights the negative consequences of decreasing rainfall on agricultural practices, water resource management, and food security, which is a well-documented consequence in literature on climate change and its impacts on agriculture and food systems. The paper [6] proposes a deep-learning-based classification method for data pages used in hydrographic memory. The paper focuses on the classification of data pages used in hydrographic memory using deep learning techniques. The paper [7] presents a conjunction model for drought forecasting that combines dyadic wavelet analysis and neural networks. The paper [8] proposes a k-Tree method for kNN classification that assigns different optimal k values to different test samples, resulting in higher classification accuracy compared to traditional kNN methods. The paper [9] aims to derive optimal data-driven machine learning methods for forecasting rainfall in Odisha, by comparing three techniques: linear regression analysis, random forest method, and Artificial Neural Network (ANN) method. The study found that the maximum rainfall occurred in 1961 (385.3 mm) and the minimum in 1974 (197.2 mm). The paper [10] describes understanding and quantifying long-term rainfall variability at a regional scale is important for a country like India where economic growth is very much dependent on agricultural production which is closely linked to rainfall distribution.

## 2. Methodology

### 2.1 Data Collection :

The initial stage of forecasting rainfall is to collect a dataset comprising numerous weather-related elements which include air pressure, wind speed, atmospheric humidity, temperature, and other important qualities.

### 2.2 Data Preprocessing:

Once the data has been acquired, it must be prepared. To maintain data quality and consistency, operations such as resolving missing values, correcting outliers, and normalizing feature values are performed. There are some stages in data preprocessing:

#### 2.2.1 Formatting:

It's unlikely that the data we collected is accessible in a way that makes it usable. The information we obtained is presented as raw data. These data were saved in a JSON file and a CSS file was created. If you intend to use it for a social database or content record, you might want it to be in a sample document or a specific record configuration.

#### 2.2.2 Cleaning :

Information cleaning also includes the elimination of missing or incomplete data. This may be instances where the data is lacking or problematic, making it impossible for you to fix the issue. It could be necessary to eliminate these cases. Additionally, some of the characters could contain invalid data; as a result, it might be necessary to reveal the characteristics.

#### 2.2.3 Sampling :

It can be challenging to work with information that is more scattered than what is readily accessible. Calculations that require manual intervention may take longer to perform and consume more memory and computing resources. You may conduct a small test on the selected data while analyzing the entire dataset, which might help in figuring out and testing the settings more easily.

#### 2.2.4 Label Encoding :

We will need to convert the target variable and the category characteristics to binary representation in this. Label encoding forms the basis of the label into binary form that is readable by machines. The functionality of these tags can be understood using ML. It is critical for supervised learning phase or the structured dataset preparation process.

### 2.2.5 Feature Scaling :

The process of "feature scaling" dish4 but the independent features in the data or API will be scaled way throughout a certain range. It happens while the data is being pre-processed. Additionally, the learning and generalization phases of the ML process are sped up by machine learning and data reduction to create variable combinations (features). Finally, we use the collected named dataset and the classifier approach to obtain our models using the classify module of the Python Natural Language Toolkit package. To assess the models, the remaining categorized data in our sample will be used. The already processed data was categorized using a few machine learning techniques. The classifier chosen was Random Forests. These methods are frequently employed in text classification tasks.

### 2.3 Model Training :

After that, the dataset is divided into training and testing sets. Using the training data, a Random Forest model is trained. To produce predictions, this ensemble learning approach mixes many decision trees.

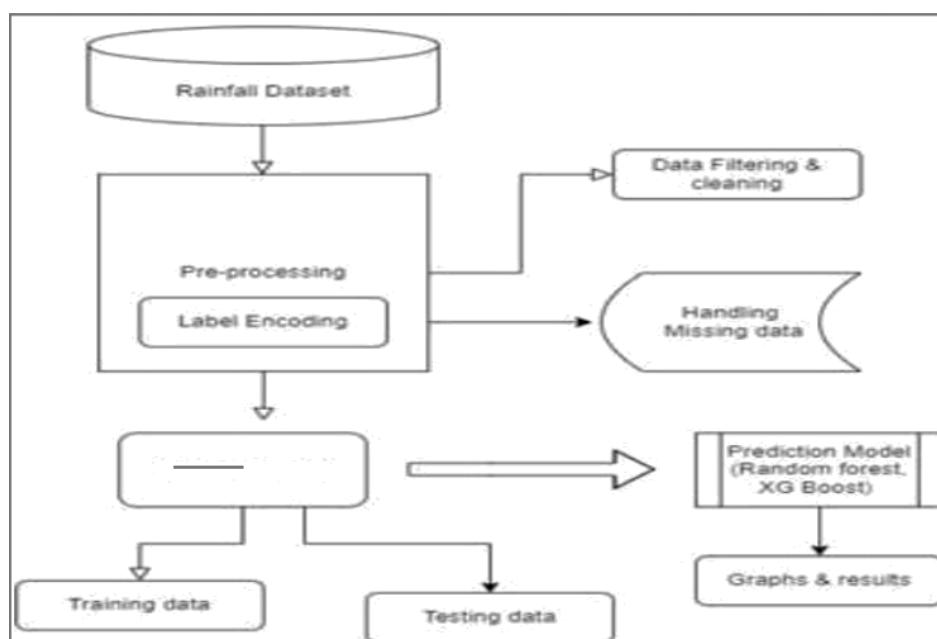
### 2.4 Model Evaluation :

The testing set is used to evaluate the model's performance. To quantify how successfully the model forecasts rainfall based on the testing data, several evaluation criteria, such as F1 score, Precision, Recall, and Accuracy, may be utilized.

### 2.5 Result :

When desired accuracy is not obtained, it is important to go back and carefully go over a specific method.

### 2.6 Architecture Model



### 2.7 Applying algorithms

#### 2.7.1 Logistic Regression :

In order to represent the connection between a dependent variable and one or more independent variables, a statistical method known as logistic regression is used. It is used to determine the likelihood that an event will occur based on predictor factors. Logistic regression may be used to forecast rainfall using a variety of weather-related factors. These predictors are chosen based on their correlation with rainfall, excluding air pressure, wind speed, and humidity. The logistic regression model is then trained using the obtained data. By examining the values of these factors, a model that can forecast the possibility of rainfall must be developed. The associations between the predictor variables and the incidence of rainfall are assessed using statistical techniques.

### 2.7.2 Decision Tree :

An effective supervised machine learning technique used for both classification and regression applications is the decision tree. It operates by repeatedly dividing the dataset into subgroups according to the most important attribute, producing a decision tree-like structure. The interior nodes of the tree stand for features, the branches for decision-making processes, and the leaf nodes for results or class labels.

Based on decision trees, ensemble approaches like Random Forest and Gradient Boosting improve their accuracy and resilience by mixing numerous trees. Decision trees are frequently employed in real-world applications for generating informed and data-driven decisions and serve a crucial role in the basis of more complex machine learning algorithms.

By training a model to understand the complex correlations between different input factors, such as temperature, humidity, wind speed, and others, and the output variable, precipitation, a machine learning model can be used to forecast rainfall. By taking into account variables like temperature, humidity, wind speed, cloud cover, and geographical characteristics like elevation and proximity to water bodies, it collects information on rainfall patterns and other relevant factors. By addressing missing values, outliers, and anomalies, preparing the data. To verify that the input variables' values are comparable, normalize them. Following that, the dataset is divided into training and testing sets. The training set will be used to build the machine learning model, and the testing set will be used to assess it.

### 2.7.4 Random Forest :

A Random Forest constructs a large number of decision trees, each using a random part of the training data and a random collection of features at each split. The final forecast is calculated by averaging the predictions of all the distinct trees. The use of randomization in data and feature selection reduces overfitting and improves the model's resilience. Unlike a Random Forest, decision trees are conducted independently of one another. Because of this capability, Random Forests are very scalable for huge datasets. Every tree is built with a distinct subset of the data and features, creating diversity and decreasing correlations between trees. This variety helps to increase model generalization.

Random Forest is a versatile and efficient algorithm that can be used for both regression and classification applications. It handles categorical and numerical data equally well, giving it a wide range of prediction applications. Estimation of Feature Significance: Random Forest offers estimations of feature significance. This useful data may be used for feature selection and feature engineering, assisting in the identification of the most significant factors during the prediction process. Random Forest is a powerful and versatile machine learning algorithm that can be used for a wide range of applications. Its capacity to reduce overfitting, handle varied data formats, and provide insights on feature relevance makes it a dependable alternative for real-world problem solving across all disciplines.

### 2.7.5 LightGBM :

LightGBM is a gradient boosting framework developed by Microsoft that uses tree-based learning algorithms. Large datasets and high-dimensional feature spaces make good use of its speed and efficiency design. LightGBM is built on the gradient boosting framework and develops a powerful predictive model by assembling a group of weak learners (often decision trees). Sequential tree construction is used, where each tree fixes the flaws of the preceding one. LightGBM has been enhanced for speed and effectiveness. It is faster than many other gradient boosting implementations, especially when working with large datasets. Histogram-based learning, which uses less memory and is faster than traditional tree methods, is one of the key strengths of LightGBM. LightGBM is highly scalable and capable of handling enormous datasets, with millions of instances per feature. Big data applications can benefit from its effective algorithms and parallel processing capability. There are several machine learning applications that make use of LightGBM, including click-Enough rate prediction, picture categorization, recommendation systems, and more. Large-scale and real-time applications can benefit from its speed and efficiency. A significant distinction between LightGBM and other gradient boosting implementations is its unique tree-building approach. Instead of constructing trees level by level, LightGBM builds them leaf-by-leaf, selecting splits that maximize loss reduction, resulting in an efficient and fast model training.

### 2.7.6 XG Boosting :

XGBoost is a well-known and efficient algorithm for rainfall prediction. It is a variety of fields, including agriculture and forestry, management, prediction, rainfall, and weather. Accurate rainfall forecasting can help in planning and decision-making processes like irrigation systems. The first step in utilizing XGBoost to predict rainfall is collecting historical rainfall data from diverse sources, such as meteorological stations or remote sensing data. The report should contain the amount of rainfall as well as other important elements like temperature, humidity, and wind speed. After it has been collected, the data can be cleaned up and preprocessed, including feature engineering, dealing with missing data, and removing anomalies. By modifying the data and creating new variables, feature engineering may include capturing the fundamental patterns in the data.

Crash data has undergone a process of cleaning, it may be separated into framing, validation, and testing datasets. It is a well-known option for many machine learning tasks, including classification, regression, ranking, and recommendation systems, because it is adaptable and performs well. XGBoost, or Extreme Gradient Boosting, is a versatile machine learning algorithm used for both regression and classification tasks. It is known for its exceptional

performance in various applications. XGBoost is an ensemble learning method that combines the predictions of multiple decision trees to create a more predictive model. Each tree corrects the errors made by the previous ones, making the model more accurate. The algorithm optimizes an objective function, which combines a loss function (measuring prediction accuracy) and a regularization term (to prevent overfitting). XGBoost can calculate feature importance scores, helping with feature selection and understanding which features drive predictions. The algorithm is designed for efficiency, allowing it to run parallel and distributed computing, making it suitable for large datasets or distributed systems.

### 3. RESULTS AND DISCUSSIONS:

The dataset for rainfall in Iihima, which includes factors like rainfall, speed, and atmospheric humidity, among others, is used to assess the performance of various machine learning models. The primary goal of the project is to utilize current rainfall data to estimate future precipitation from 1901 to 2015, and from those predictions, select the best model. The accuracy of the quick outcomes depends on the climate, people's ability to predict rainfall, and the skill of farmers because it relies significantly on rainfall. The forecast of rainfall is more suitable for agriculture than for other purposes. The accuracy for a decision tree is 55.92%, whereas it is 79.514% for logistic regression. Neural networks' accuracy is 55.40%, Random forests' accuracy is 92.90%, LightGBM's accuracy is 57.43%, and XGBoost's accuracy is 95.79%.

Model used	Accuracy	ROC Area under curve	Cohen's Kappa	F1 score
Logistic Regression	0.79	0.75	0.78	4.051
DecisionTree	0.57	0.8*	0.74	0.60
Neural Network	0.88	0.88	0.76	422.22
RandomForest	0.92	0.93	0.81	40.44
LightGBM	0.67	0.57	0.74	3.97
XGBoost	0.9*	0.96	0.91	192.10

TABLE 1: Comparison Table of Used Algorithms

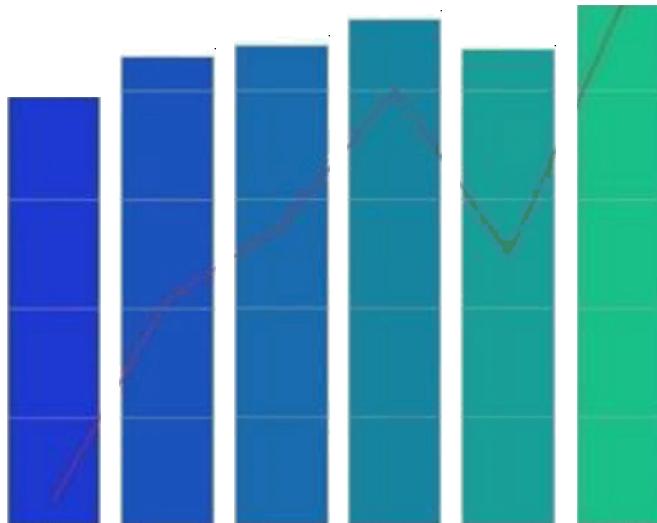


FIGURE 1: Area under ROC and Cohen's Kappa



**FIGURE2:AccuracyandTime takenforexecution**

### **CONCLUSIONSANDFUTURESCOPE:**

Chicofthemostsignificantnaturaloccurrences that hasaninfluenceonbodipeoplearidodierlivugthingsisramRainfallcyclesarealteringasaresultofchanging climatic conditions and irieasuiq global temperatures. Keeping ari eye on diis riatual event is crucial suire climate changeliuts trade,agriculture,andsporadicallycauseslandsipsaridflooding. Therefore,predictingrainfallhasbenefitsforagriculture, conservation, andpublicav'aieriesofnatural disasters like Ooods. Asystem toanticipate rainfall usuigrieual amfieial intelligence, v'lieli iscommon inmodern teeluiology,isneededtoaddress diese problems arid provide basic needs.

Researchers should investigate v'aysto adaptanfall prediction models to account forchanging climatic factors and evaluate themodels' efficacy ui arange of climatic conditions. creating sofhvae to help decisions: Chice a reliable prediction model has been developedit may be applied to thedevelopment of decision-inakuigapplications for a number of sectors, such as agriculture and v'ater management. Researchers may look intov'aysto develop user-friendlyuiterfaeesanddecision-support tools tohelpstakeholders make gooddecisions based on die expectedrainfall.

### **REFERENCES:**

1. G.Di Baldassalve,A.hloritalali, H.Luis, D. Koutsoyaliuus,L.Blandilnarte, andG.Blöselil,"Flood fatalities iriAfiea: From diagriosis tolmtigation," Geophys. Res.Lett., vol. 3.7,rio. 2.2, pp. 529-546, Nov. 2010.
2. N.SalnsidiSam,I.Sirlaeli,hI.Hassan,A.Hali, andhI.Aliff,"Eulariririghlalaysia1ainfallpledietioriusuigelassificationitehluiques,"J.Appl. Envilm. Biol. Sci., vol.7, no. 2S, pp. 2m-29,201.7.
3. L.Ainekudzi,E. Yainba, K.Reko,E. Asale, J. Aryee,hI. Baidu, and S. Codjoe,"5'ariabilitiessir1ailfalloliset,eessationiandlerighthof lairyseasorifor dre variousagro-ecologiral zoriesofGhalia,"Chlnate,yol.3,rio.2,pp. 416M34,Jun.2015.
4. C.hleisali,L.Ainekudzi,N.Klutse,J.Aryee, andKAsare,Coinparisoroflamyseasorionset,eessationiarifdulationforGhaliafiolnRegChl4andGhletdatasets,"Atunos.ChlnateSci.,vol.6,rio.1,pp.300—309,2016,doi:10.4236/aes.2016.62025.
5. hI.Baidu,L.KAinekudzi,J. Aryee,arifT. Arrior,"Assessimeritoflong-terimspatio-temporalailfallvariabilityoverGhaliausing v'aveletahalysis,"Chlnate,yol. 5, rio. 2, p. 30, hfal.201.7.
6. T. Sirunobaba,N. Kuv'ata, hI. Hoinina,T.Takaliaslu,Y. Nagalaina, hI. Sario, S. Hasegav'a,R.Hitayalna,T. Kakue, A. Sirndaki,N. Takada, andT.Ito,"Deep-learning-basedhatapageclassificationfor holographlneinoly,"2017,arXiv:1 707.00654.
7. T.-V'.KunandJ.B.5'aldös,"Nonlirrealmodel fordought foleeastırıgbasedoriaeoriurirtiorfov'avelettrarisfolinsarifcuelalrienvoiks," J.Hydrol.Erig.,yol.S,rio.6, pp.319—32S,Nov.2003.
8. S. Zhang,X. Li, hI. Zorig,X. Zliu, andR. V'arig,"Efficient kNN classificationv'idiifferent numbers of nearest neighbors," IEEE Trans.Neural Nenv. Learn. Syst., vol. 2.9, rio. 5, pp. 1774-1755, hlay 2015.
9. hlisla,R.K,Panda,P.K.,Salu,AK.,Saloo,S.,&Beliera,D.P.(2021).Raifall predictioriusuiglnaelirrie learrning approael: A casestudyfor the state of ofisha. hdiari Jounial of Natulal Seierires.
10. hlohiapatia,G.,Rakeshi,5'.,Puru'ar,S.,&Dilmi,A.P.(2021).Spatio-temporalailfallvariabilityoverdiffereritneteoilogiralsubdivisionisirihidia:analysis usingdifferent Inaeluiulearning teeluiques.TlieoletiealarifAppliedChlnatology,145(1-2),673-656.