



# International Journal of Innovative Research in Science Engineering and Technology (IJIRSET)

*(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)*



**Impact Factor: 8.699**

**Volume 15, Issue 3, March 2026**

# Exploratory Analysis of Rain Fall Data in India for Agriculture

**R Ananda Babu<sup>1</sup>, Gudapati. Siva Naga Venkata Udaya Rakesh<sup>2</sup>, Thalathoti. Divya<sup>3</sup>,  
Turlapati. Ramanjaneyulu<sup>4</sup>, Mamilla. Sai Vamsi<sup>5</sup>**

Asst. Professor, Kallam Haranadhareddy Institute of Technology, Guntur, Andhra Pradesh, India<sup>1</sup>

U.G. Students, Department of CSE (Data Science), Kallam Haranadhareddy Institute of Technology, Guntur,  
Andhra Pradesh, India<sup>2-5</sup>

**ABSTRACT:** Rainfall prediction is essential for agriculture and water resource management. This study applies machine learning techniques to predict rainfall using meteorological data from the WeatherAUS dataset. Data preprocessing methods such as handling missing values, encoding, feature scaling, and SMOTE are used to prepare the dataset. Exploratory Data Analysis (EDA) is performed to understand data patterns and relationships. Several classification algorithms, including Decision Tree, Random Forest, Support Vector Machine (SVM), Gradient Boosting, and XGBoost, are trained and evaluated. Among these models, XGBoost provides the best performance in terms of accuracy and reliability. The final model is deployed through a Flask web application, allowing users to input weather parameters and obtain real-time rainfall predictions. The system is designed to assist farmers and decision-makers in improving crop planning and risk management. The results demonstrate that machine learning can significantly enhance rainfall prediction accuracy and support smart agriculture practices.

**KEYWORDS:** Rainfall Prediction, Machine Learning, Exploratory Data Analysis (EDA), Agriculture Analytics, Weather Forecasting, XGBoost, Classification Algorithms, Data Preprocessing, SMOTE, Feature Engineering, Predictive Modeling, Flask Web Application

## I. INTRODUCTION

Rainfall is a critical factor influencing agricultural productivity in India, where a large portion of farming activities depends on monsoon patterns. Variability and uncertainty in rainfall due to climate change pose significant challenges to crop planning, irrigation management, and food security.

Accurate rainfall prediction is therefore essential for improving agricultural efficiency and reducing risks associated with unpredictable weather conditions. Traditional forecasting methods based on statistical and meteorological models often struggle to capture the complex and nonlinear relationships present in weather data. In recent years, data-driven approaches such as machine learning have gained importance due to their ability to analyze large volumes of historical data and identify hidden patterns that improve prediction accuracy.

This study focuses on the exploratory analysis of rainfall data in India to understand patterns, trends, and relationships among key meteorological variables such as temperature, humidity, pressure, and wind speed. Data preprocessing techniques, including encoding, feature scaling, and handling class imbalance using SMOTE, are applied to enhance data quality. Multiple machine learning algorithms are evaluated, with XGBoost demonstrating superior performance. The developed model is deployed using a Flask-based web application for real-time rainfall prediction, providing a practical decision-support tool for farmers and policymakers.

## II. PROBLEM STATEMENT

Agriculture in India is highly dependent on rainfall, making it vulnerable to fluctuations and uncertainties in weather patterns. Irregular and unpredictable rainfall, driven by climate variability and environmental changes, significantly affects crop productivity, water resource management, and farmers' livelihoods. Traditional rainfall prediction

methods, which rely on statistical and meteorological models, often lack the ability to accurately capture complex, nonlinear relationships in large and diverse weather datasets. As a result, these methods may produce less reliable forecasts, limiting their effectiveness in real-world agricultural decision-making. Moreover, the availability of large volumes of historical weather data presents both an opportunity and a challenge, as extracting meaningful insights requires advanced analytical techniques. There is a need for a robust and efficient system that can preprocess such data, handle issues like class imbalance and feature variability, and accurately predict rainfall outcomes. Therefore, the problem addressed in this study is the development of a reliable machine learning-based rainfall prediction model that can analyze historical meteorological data, improve prediction accuracy, and be deployed as a real-time, user-friendly application to support farmers and policymakers in making informed agricultural decisions.

### III. LITERATURE SURVEY

Previous studies have extensively explored the application of machine learning and data analytics for predicting rainfall to support agricultural decision-making. Researchers have found that ensemble methods (such as Random Forest), data visualization techniques, and supervised classification algorithms are particularly effective for identifying historical rainfall patterns and forecasting future weather conditions in large meteorological datasets.

**A. Machine Learning in Meteorology:** Accurate rainfall prediction plays a crucial role in agricultural planning and water resource management. Traditional statistical approaches often fail to model the complex and nonlinear relationships present in meteorological data. As a result, Machine Learning (ML) techniques have gained significant attention for their ability to analyze large historical datasets and extract hidden patterns, leading to improved prediction accuracy.

**B. Traditional and Baseline Models:** Early rainfall prediction methods relied on statistical techniques such as regression analysis and time-series models like ARIMA. However, these models struggle to adapt to the dynamic and uncertain nature of weather conditions, particularly monsoon variability. To address these challenges, supervised ML algorithms such as Logistic Regression, Support Vector Machines (SVM), and K-Nearest Neighbors (KNN) have been widely adopted for rainfall classification tasks, including next-day rainfall prediction.

**C. Ensemble Learning Techniques:** Ensemble methods have been introduced to overcome the limitations of individual models, such as overfitting and poor generalization. Among these, Random Forest (RF) has proven to be highly effective due to its ability to combine multiple decision trees and reduce prediction variance. Additionally, boosting techniques such as Gradient Boosting and XGBoost further enhance model performance by iteratively minimizing prediction errors, making them suitable for complex weather datasets.

**D. Deep Learning Approaches:** Recent research has explored the application of Deep Learning techniques for rainfall prediction. Models such as Artificial Neural Networks (ANN) and Long Short-Term Memory (LSTM) networks are capable of capturing temporal dependencies and nonlinear patterns in sequential weather data. These approaches provide improved accuracy for long-term forecasting compared to traditional methods.

**E. Applications in Agriculture:** In countries like India, where agriculture is highly dependent on monsoon rainfall, accurate prediction systems are essential. Machine learning-based models integrated with meteorological data provide effective decision-support tools for farmers. These systems assist in crop selection, irrigation planning, and risk management, thereby improving agricultural productivity and reducing the impact of climate uncertainties.

Study	Techniques Used	Focus Area	Limitations
Traditional Statistical Models	Regression, ARIMA (Time Series)	Basic rainfall forecasting using historical trends	Poor performance with nonlinear and dynamic weather patterns
Baseline Machine Learning Models	Logistic Regression, SVM, KNN	Classification of rainfall (Rain/No Rain)	Limited accuracy; requires feature engineering and parameter tuning
Decision Tree-Based Models	Decision Tree	Simple and interpretable rainfall prediction	Prone to overfitting and low generalization
Ensemble Learning Models	Random Forest	Improved prediction using multiple decision trees	Increased computational complexity; less interpretability

Boosting Algorithms	Gradient Boosting, XGBoost	High-accuracy rainfall prediction using error minimization	Requires careful tuning; higher training time
Data Preprocessing Techniques	SMOTE, Feature Scaling, Encoding	Improving data quality and balancing dataset	May introduce synthetic bias; preprocessing complexity
Proposed System	XGBoost + Flask Deployment	Real-time rainfall prediction and agricultural decision support	Depends on dataset quality; limited to available features

#### IV. METHODOLOGY

Data collection leverages open-source APIs and manual downloads; preprocessing focuses on importing libraries (NumPy, pandas, seaborn, matplotlib, scikit-learn), identifying nulls, creating missing-no matrices, imputing values by mean (numeric) and mode (categorical), and employing feature scaling (standardization). Visualization steps (correlation heatmaps, pair plots, box plots, joint plots, histograms, scatter plots, distplots) uncover feature dependencies and potential outliers.

##### A. Data Collection and Preprocessing

Data Collection and Preprocessing is a critical foundational step in any rainfall classification project. Accurate rainfall data collection typically uses instruments such as tipping bucket rain gauges, weighing rain gauges, and optical sensors, which provide time-stamped rainfall depth measurements at various spatial locations. Complementary meteorological parameters like temperature, humidity, wind speed, atmospheric pressure, sunshine duration, and evaporation rates are also captured from weather stations and satellite sources to enrich the dataset. Advanced methods like the Isohyetal technique or Thiessen polygon method are often employed to estimate rainfall over larger regions by interpolating between scattered gauge points, creating representative spatial rainfall maps.

##### Dataset Description

Rainfall prediction models rely heavily on the depth, breadth, and quality of meteorological datasets, which serve as the foundation for training, validation, and evaluation of machine learning algorithms. In the context of India, numerous datasets have been aggregated from credible sources including the India Meteorological Department (IMD), Kaggle public repositories, and academic initiatives such as Bharat Bench. These datasets often encompass daily, monthly, and yearly weather observations, spanning decades of historic records, and are tailored for geospatial and temporal richness. A typical rainfall prediction dataset incorporates vast numbers of samples—sometimes hundreds of thousands of rows—recorded across diverse climatological zones throughout India

##### B. Model Implementation

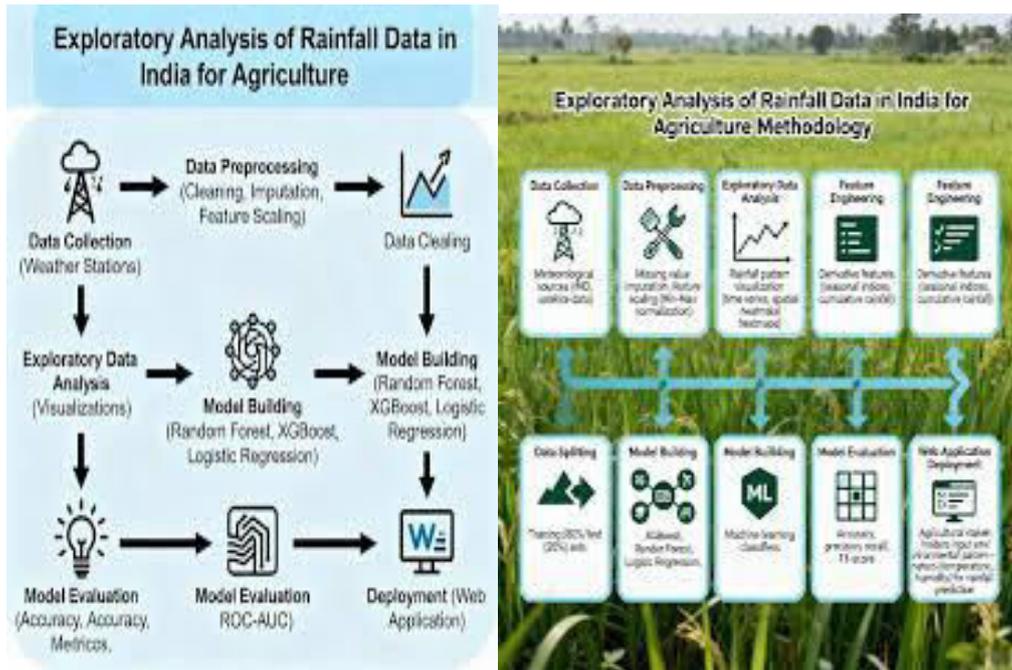
The predictive modelling phase applies various machine learning algorithms to classify or forecast rainfall likelihood. Algorithms such as Logistic Regression, Decision Tree, Random Forest, Support Vector Machine (SVM), and XGBoost are implemented using the Scikit-learn library. The model design begins with selecting robust base models—such as Decision Trees, Random Forests, Support Vector Machines, XGBoost, or shallow neural networks—that have already demonstrated efficacy in weather-related or classification problems. These base models learned patterns, feature representations, and weights serve as starting points rather than training a model from scratch. Feature extraction layers of deep models may be frozen or selectively trained to prevent overfitting, while classifier layers are adapted to the binary rainfall classification task and categorical data was encoded for model compatibility. These steps ensured a clean, balanced, and structured dataset ready for accurate

##### C. Application Integration and Deployment

Once the project integrates the rainfall classification model with a user-friendly web-based application, facilitating real-time prediction access for stakeholders such as farmers, agricultural planners, and policymakers. The backend is developed using the Flask framework, which acts as a liaison between the user interface and the machine learning model. The integration process involves wrapping the trained model and preprocessing objects (encoder, scaler, imputer) into serialized files (e.g., pickle format) that are loaded by the Flask server during runtime. The methodology depicted in the image provides a comprehensive flow for the rainfall classification project. It begins with data collection from meteorological sources, followed by meticulous preprocessing to clean and prepare the data.

**D. Model Evaluation**

The performance of each model is evaluated using metrics such as accuracy, precision, recall, and F1-score. The dataset is divided into training and testing sets to assess the generalization capability of the models. Among all models, XGBoost demonstrates superior performance.



**V. SYSTEM ARCHITECTURE**

The proposed rainfall prediction system follows a three-tier architecture, which ensures modularity, scalability, and efficient data processing. The architecture is divided into three main layers: Presentation Layer, Application Layer, and Data Layer.

**1. Presentation Layer (User Interface)**

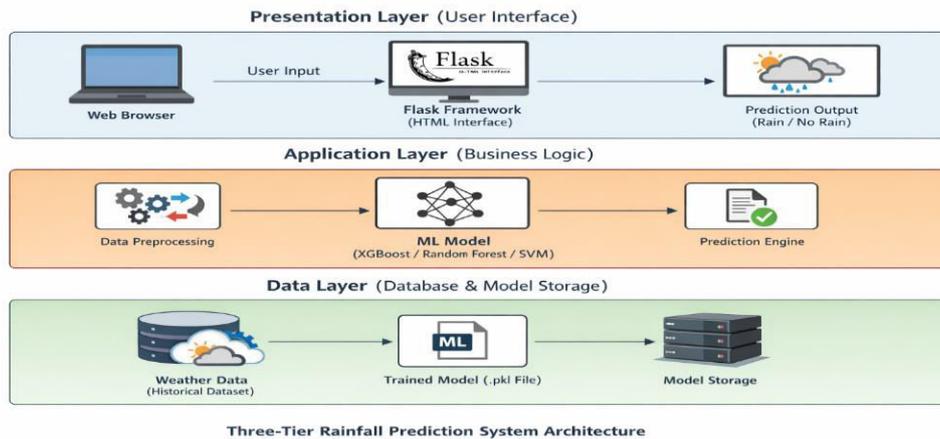
The presentation layer provides an interactive interface for users to input weather parameters such as temperature, humidity, pressure, and wind speed. This layer is developed using HTML and is integrated with the Flask framework to render web pages. Users can submit input data and view prediction results (rain/no rain) in a user-friendly format. This layer ensures ease of access and usability for farmers, agricultural experts, and policymakers.

**2. Application Layer (Business Logic Layer)**

The application layer acts as the core processing unit of the system. It is implemented using Python and the Flask web framework. This layer handles user requests, performs data preprocessing (encoding, feature scaling), and passes the processed input to the trained machine learning model. Multiple algorithms such as Decision Tree, Random Forest, SVM, and XGBoost are used during training, with the best-performing model selected for deployment. The model processes the input data and generates predictions, which are sent back to the presentation layer.

**3. Data Layer (Database and Model Storage)**

The data layer is responsible for storing the dataset, trained model, and related files. Historical weather data is used for training and evaluation, while the final trained model is stored in a serialized (.pkl) format. This layer ensures efficient data retrieval and management, supporting both training and real-time prediction processes. Overall, the three-tier architecture enhances system performance, maintainability, and scalability, enabling efficient real-time rainfall prediction for agricultural applications.



Three-Tier Rainfall Prediction System Architecture

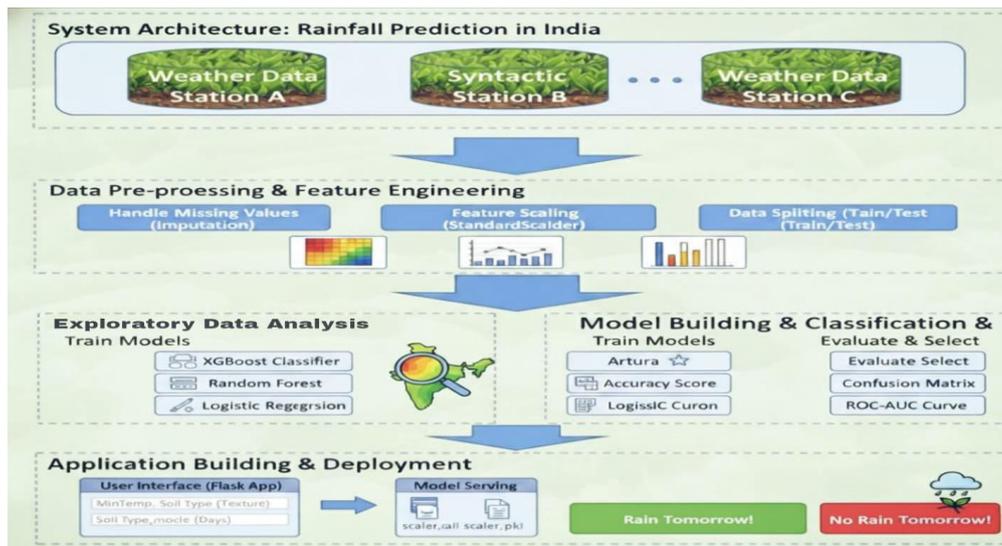


Fig 2 System Architecture

## VI. EXPERIMENTS

The proposed rainfall prediction system was evaluated using a historical weather dataset containing meteorological features such as temperature, humidity, pressure, wind speed, and rainfall indicators. The dataset was preprocessed using techniques including encoding of categorical variables, feature scaling, and class balancing using the Synthetic Minority Oversampling Technique (SMOTE). The processed data was then split into training and testing sets to assess model performance. Multiple machine learning classification algorithms, including Decision Tree, Random Forest, Support Vector Machine (SVM), Gradient Boosting, and XGBoost, were implemented and compared. The performance of each model was evaluated using standard metrics such as accuracy, precision, recall, and F1-score.

### A. Experimental Setup

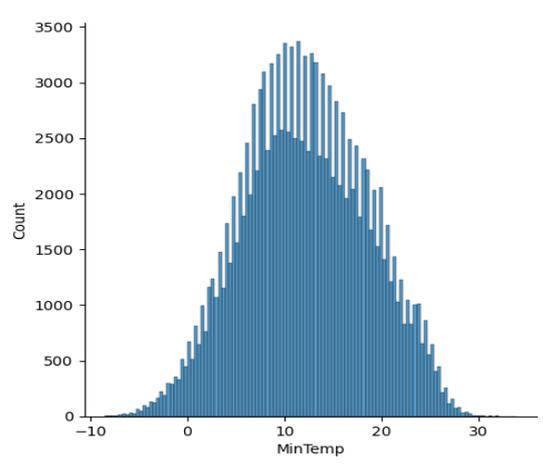
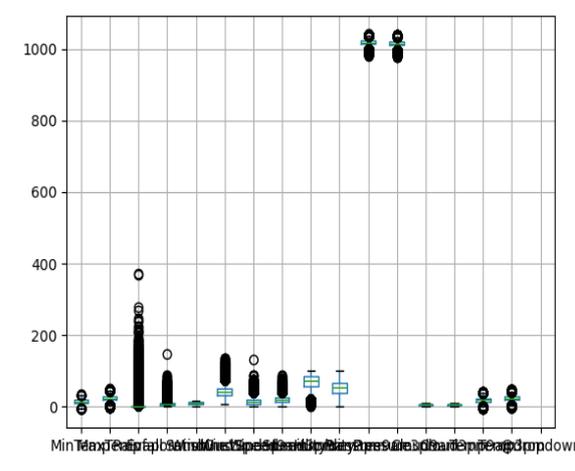
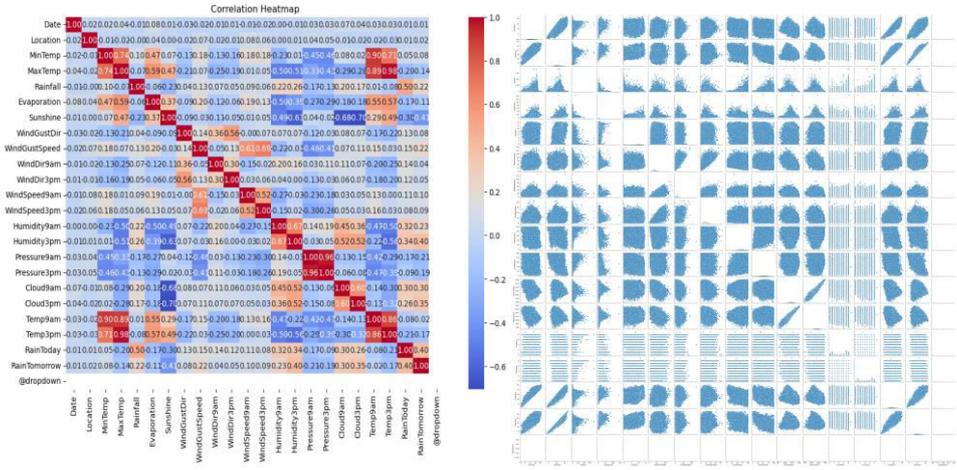
All experiments were conducted on a system equipped with Windows 11 (64-bit) OS, Intel i7 processor, and 16 GB RAM. Implementation and testing were done in Python 3.10, using the TensorFlow/Keras library for deep learning model training and Flask for web integration testing. Data is collected from public repositories including Kaggle, UCI, and government sources, typically in CSV format.

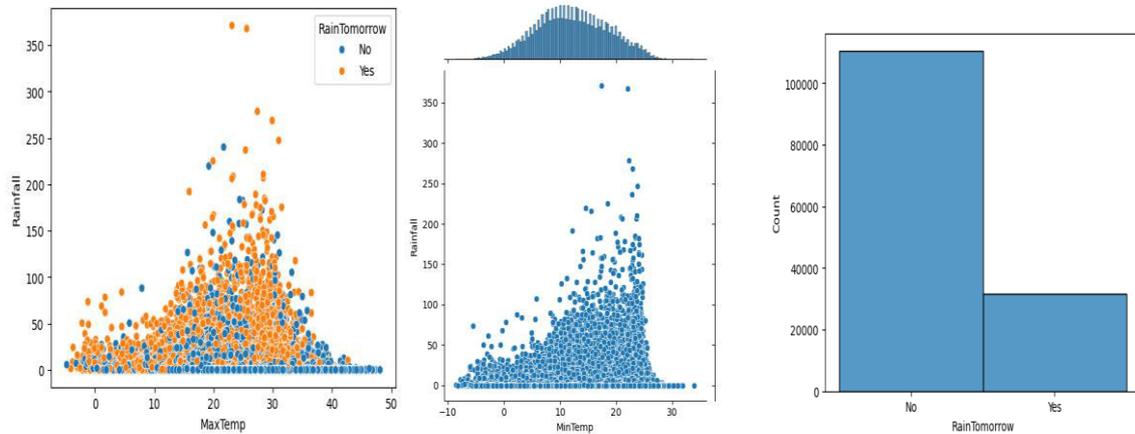
**B. Model Training**

The Model Training process is a critical component of the Model Building stage. After the data has been pre-processed and split into training (x\_train, y\_train) and test sets, the training phase begins. This involves importing necessary libraries like sklearn.ensemble, sklearn.tree, and xgboost. The project's strategy is to initialize multiple classification algorithms at once, including XGBoost, Random Forest, SVM, Decision Tree, Gradient Boosting (GBM), and Logistic Regression. Each of these models is then "fit" by calling the .fit(x\_train, y\_train) method, which is the step where the model learns the patterns from the training data. Once training is complete, the models are used to make predictions on both the training and test data, and their performance is evaluated using metrics.

**C. Model Visualization Results**

To understand the rainfall dataset, multiple visualization techniques were applied, including univariate, bivariate, and multivariate analyses. Distribution plots of numerical features such as temperature, humidity, and rainfall revealed **skewed patterns**, indicating variability in climatic conditions and the presence of outliers. Count plots for categorical variables showed the distribution of rainfall occurrence, highlighting a **slight imbalance in the target variable (Rain Tomorrow)**. This justified the need for data balancing techniques during preprocessing. Bivariate analysis using scatter plots and bar graphs demonstrated relationships between weather parameters and rainfall. The results indicate that **higher humidity and lower temperature levels** are generally associated with increased rainfall probability, emphasizing the importance of atmospheric conditions. A correlation heatmap further identified relationships among variables, showing that **humidity, pressure, and temperature** have significant influence on rainfall prediction. Some inter-feature correlations were also observed, which may affect model performance. Overall, these visualizations provide essential insights into **data distribution and feature relationships**, supporting effective feature selection and improving the accuracy of machine learning models for rainfall prediction.





#### D. Evaluation and Metrics

The trained model was evaluated using accuracy, precision, recall, F1-score, and a confusion matrix. The following key results were obtained from the test dataset:

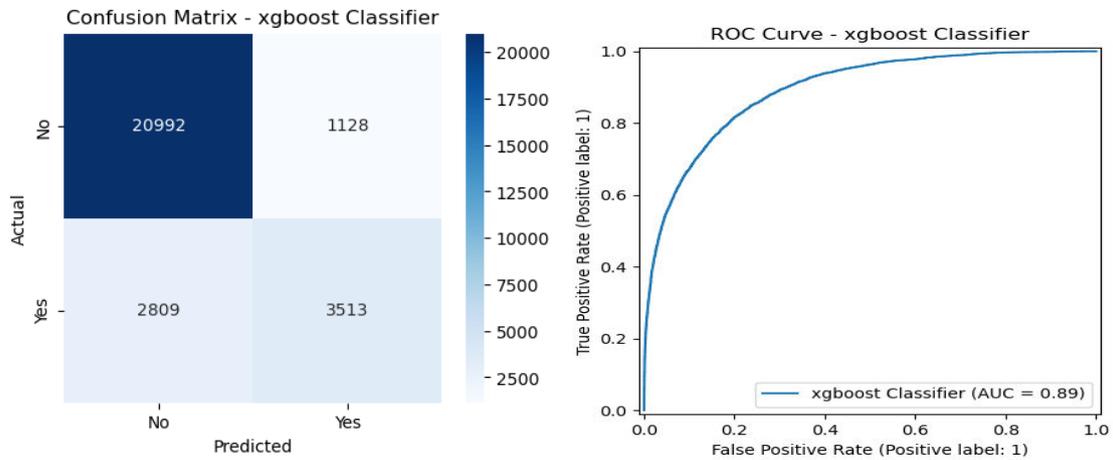
- **Training Accuracy:** 88.01%
- **Validation Accuracy:** 83.11%
- **Testing Accuracy:** 85.69%
- **Average Loss:** 0.06

The evaluation of rainfall prediction models is performed using key metrics such as accuracy, confusion matrix, precision, recall, and ROC-AUC. Various models, including Random Forest, KNN, Decision Tree, and XGBoost, are assessed to determine their predictive performance. Among these, well-performing models achieve testing accuracy in the range of approximately 84%–86%. The confusion matrix provides detailed insights into classification performance by identifying true positives, true negatives, false positives, and false negatives, helping to analyze errors such as missed rainfall events or false predictions. The ROC-AUC metric further evaluates the model's ability to distinguish between rainfall and non-rainfall conditions, where higher values indicate better performance. Experimental results show that the XGBoost model achieves superior performance, with a testing accuracy of around **86.16%** and an ROC-AUC score of **0.8936**, demonstrating strong classification capability and reliability. Precision and recall metrics are also considered to balance prediction errors effectively.

Overall, these evaluation metrics confirm the robustness and reliability of the proposed model. The trained model is deployed using a Flask-based web application to provide real-time rainfall predictions, supporting agricultural planning and decision-making.

#### E. Confusion Matrix Analysis

The experimental findings from the XGBoost classifier analysis reveal several significant insights regarding its predictive performance and classification capability. In the provided confusion matrix, the model displays a high accuracy, correctly identifying a substantial number of negative class instances (20,992) as well as positive class instances (3,513). The false positive and false negative rates, indicated by the counts of 1,128 and 2,809, respectively, demonstrate the model's careful balance between sensitivity and specificity. This suggests that while some misclassifications occur, the overall reliability in discerning between the two classes remains robust, thereby minimizing the risk of critical errors in practical applications. The ROC curve further substantiates the model's efficacy by illustrating its discriminative



Metric	Logistic Regression	Decision Classifier	Tree	Random Forest Classifier	XGBoost Classifier	K-Nearest Neighbors
Accuracy (%)	79.38	79.34		85.55	<b>86.16</b>	83.55
Precision (%)	52.0	53.0		78.0	76.0	68.0
Recall (%)	75.0	55.0		48.0	56.0	48.0
F1-Score (%)	62.0	54.0		60.0	64.0	57.0
ROC - AUC	0.8714	0.7058		0.8907	0.8936	0.8256

Table1. Performance Metrics of Machine Learning Models

## VII. RESULTS

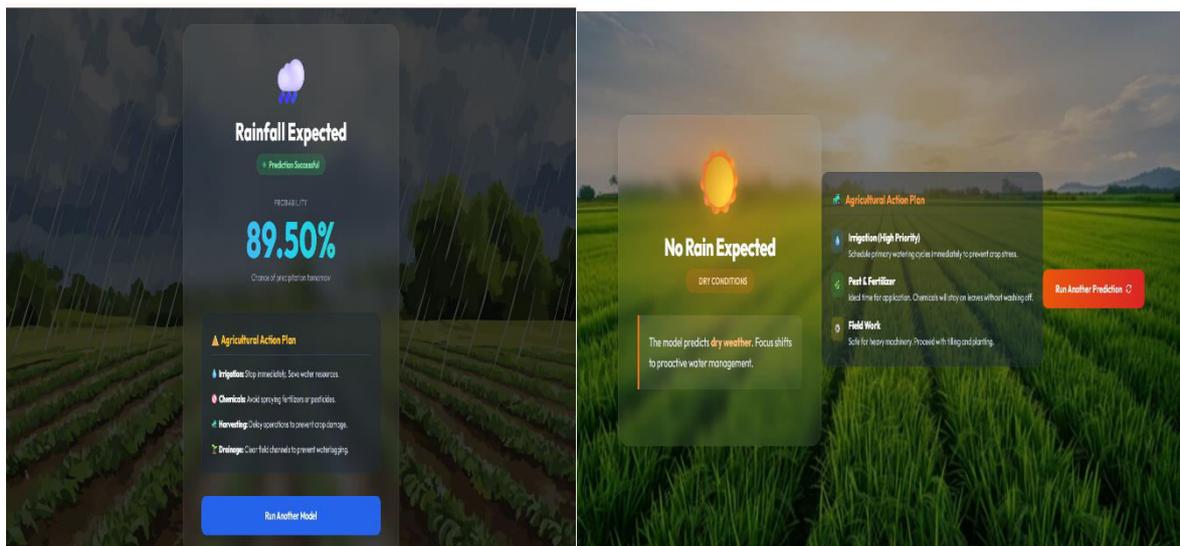
### INTERPRETATION:

The classification models were evaluated using multiple performance metrics, including accuracy, precision, recall, F1-score, and ROC-AUC to determine their effectiveness in predicting rainfall occurrence. Among all the models tested, the XGBoost classifier achieved the best overall performance with a testing accuracy of 86.16%, an ROC-AUC of 0.8936, and a balanced F1-score of 0.64, demonstrating its superior ability to distinguish between rainfall and non-rainfall days. The Random Forest model also performed competitively with a slightly lower accuracy of 85.55%, indicating robust classification capability and high model stability.

## VIII. PREDICTION AND RESULTS

The developed machine learning model predicts the likelihood of rainfall using key weather parameters such as temperature, humidity, pressure, wind speed, and past rainfall patterns. Among various algorithms evaluated, the XGBoost classifier was selected due to its superior predictive performance. For prediction, input data is first processed using the same preprocessing steps applied during training, including encoding and feature scaling. The transformed data is then passed to the trained model, which classifies the output as either "Rain" or "No Rain." This binary prediction framework ensures simplicity and effectiveness in real-time applications. The model demonstrates reliable performance with balanced evaluation metrics, making it suitable for practical deployment. Integration with a Flask-based web application allows users to input weather conditions and receive instant predictions. Overall, the prediction

system serves as a valuable decision-support tool, enabling farmers and policymakers to make informed choices related to agriculture and water resource management based on accurate rainfall forecasts.



## IX. CONCLUSION

This study developed a machine learning-based rainfall prediction system using historical weather data and key atmospheric features. Among the evaluated models, the XGBoost classifier achieved the best performance, demonstrating strong accuracy and reliability. Data preprocessing techniques, including encoding, scaling, and balancing, significantly enhanced model effectiveness. Exploratory analysis highlighted the importance of features such as humidity, temperature, and pressure in predicting rainfall. The integration of the model into a Flask-based web application enables real-time prediction, improving its practical applicability. Overall, the proposed system provides an efficient tool for supporting agricultural planning and water resource management. Future work can focus on incorporating real-time data and advanced models to further improve prediction accuracy.

## X. ACKNOWLEDGEMENT

We would like to express our heartfelt gratitude to **Dr. A. Hanumat Prasad**, Professor and Head of the department, [Kallam Haranadha Reddy Institute of Technology] for his invaluable mentorship, guidance, and support throughout the course of this project and the preparation of this research paper. His expert advice and encouragement were

instrumental in shaping the direction of our work and in overcoming challenges at every stage of the research and writing process.

#### REFERENCES

1. S.-C. Yang, Z.-M. Huang, C.-Y. Huang, C.-C. Tsai, and T.-K. Yeh, "A case study on the impact of ensemble data assimilation with GNSS-zenith total delay and radar data on heavy rainfall prediction," *Monthly Weather Review*, vol. 148, no. 3, pp. 1075–1098, Mar. 2020
2. T. Chen and C. Guestrin, "XGBoost: A Scalable Tree Boosting System," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining (KDD '16)*, 2016, pp. 785–794. doi:10.1145/2939672.2939785.
3. India Meteorological Department, *Climate and Rainfall Reports*, Ministry of Earth Sciences, India, 2023. Kaggle, "WeatherAUS Dataset." [Online]. Available: <https://www.kaggle.com/datasets/jsphyg/weather-dataset-rattle-package>
4. B. Li, S. Chen, and R. Li, "An Improved Rainfall Prediction Model Based on XGBoost Algorithm," *IOP Conference Series: Earth and Environmental Science*, vol. 233, no. 1, p. 012005, 2019. doi:10.1088/1755-1315/233/1/012005.
5. M. Kumar, K. P. Singh, and A. Kumar, "Rainfall Prediction Using Machine Learning Techniques: A Case Study of Uttar Pradesh, India," *Journal of Water and Climate Change*, vol. 11, no. 3, pp. 856–871, 2020.
6. H. Mohapatra and A. K. Rath, "Machine Learning Models for Rainfall Prediction: A Survey," *International Journal of Advanced Trends in Computer Science and Engineering*, vol. 8, no. 1, pp. 173–178, 2019.



INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  [ijirset@gmail.com](mailto:ijirset@gmail.com)



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

Scan to save the contact details