

**University of Barisal**  
**Department of Computer Science and Engineering**



**Project Report**

"Weather Condition Prediction Using Machine Learning"

**Proposed by**  
Shahade Parvaze Siam  
19CSE028  
Session: 2018-19

Department of Computer Science and Engineering

University of Barisal

## Introduction

Weather prediction is a critical task that influences many sectors, including agriculture, transportation, and disaster management. Accurately forecasting weather conditions helps individuals and organizations make informed decisions, ultimately improving safety and operational efficiency. Over the years, traditional meteorological techniques have been supplemented by machine learning algorithms, which offer the potential for enhanced accuracy and efficiency by analyzing large volumes of weather-related data.

In this project, we leverage machine learning techniques to predict weather conditions based on historical weather data. The dataset consists of features such as temperature, humidity, wind speed, pressure, and other environmental factors, which are used to predict the daily weather summary, such as "Partly Cloudy," "Rain," "Clear," etc. The objective is to build a model that can accurately classify the weather based on these factors, providing an automated solution to weather prediction.

We employ data preprocessing methods such as encoding categorical variables, handling missing values, and normalizing continuous features. A Random Forest Classifier, a robust ensemble learning model, is used for training and prediction. The performance of the model is evaluated through various metrics, including accuracy, classification report, confusion matrix, and a visual analysis of the results. The goal of this project is to demonstrate the feasibility of using machine learning for weather prediction and to explore how different features contribute to model performance.

This project contributes to the growing field of AI-driven weather forecasting, showcasing the potential of machine learning in applications that traditionally relied on physics-based models. Through this approach, we aim to develop a model that is both accurate and efficient in predicting weather conditions based on historical data.

## Objective

The primary objective of this project is to develop a machine learning model that can accurately predict daily weather conditions based on historical weather data. Specifically, the goals of this project are:

1. **Data Preprocessing:** To clean and preprocess the provided weather dataset by handling missing values, encoding categorical variables, and normalizing continuous features to ensure optimal input for machine learning models.
2. **Model Development:** To implement and train a machine learning model (Random Forest Classifier) for weather condition prediction, using various weather-related features such as temperature, humidity, wind speed, and pressure.
3. **Model Evaluation:** To evaluate the model's performance using metrics like accuracy, precision, recall, F1-score, and confusion matrix, in order to assess its ability to predict

different weather conditions effectively.

4. **Visualization and Interpretation:** To create visual representations of model performance, such as bar plots for class-wise precision, recall, F1-score, and a confusion matrix, for deeper insights into the model's strengths and weaknesses.
5. **Prediction and Deployment:** To develop a robust solution for predicting weather conditions on new, unseen data, providing a foundation for further improvements and potential deployment in real-world applications.

By achieving these objectives, the project aims to demonstrate the potential of machine learning in weather prediction, offering a data-driven approach that can complement traditional meteorological methods.

## Research Questions

1. How effectively can machine learning algorithms, such as Random Forest, predict weather conditions based on historical weather data compared to traditional meteorological methods?
  - This question explores the performance of machine learning models in the context of weather prediction and compares them with classical approaches, assessing whether machine learning can provide a more accurate or efficient solution.
2. What is the impact of different weather features (e.g., temperature, humidity, wind speed, pressure) on the accuracy of weather condition predictions using machine learning?
  - This question investigates the significance of various features in the dataset, determining which factors contribute the most to improving model accuracy and prediction reliability.
3. How can the inclusion of data preprocessing techniques, such as handling missing values, encoding categorical variables, and feature scaling, influence the performance of machine learning models in weather prediction?
  - This question examines the effect of data preprocessing on model performance, exploring how different techniques can improve the quality of input data and, in turn, enhance the accuracy and robustness of weather prediction models.

## Motivation

Weather forecasting plays a vital role in numerous sectors such as agriculture, transportation, energy, and disaster management. Accurate weather predictions help mitigate risks, optimize operations, and ensure safety. Traditionally, weather forecasts have been driven by complex physical models and meteorological expertise, but recent advancements in machine learning (ML) offer a promising alternative. Machine learning, with its ability to analyze large datasets and uncover complex patterns, has the potential to revolutionize weather prediction by providing faster and potentially more accurate forecasts.

The motivation behind this project is to explore the application of machine learning algorithms, specifically Random Forest, to predict weather conditions based on historical weather data. By leveraging weather-related features such as temperature, humidity, wind speed, and pressure, the goal is to develop a model that can classify weather conditions with high accuracy. The project seeks to harness the power of ML to create a solution that can complement traditional methods, offering an automated and efficient alternative for weather prediction.

In addition, this project aims to demonstrate how data preprocessing, feature selection, and model evaluation techniques can be applied to improve prediction performance. With the increasing availability of large-scale weather data, machine learning presents an opportunity to create models that can continuously improve and adapt, enhancing the precision of weather forecasts and contributing to fields where real-time weather data is crucial.

Ultimately, this project is motivated by the desire to explore the intersection of data science and meteorology, with the goal of advancing automated weather prediction models that are both accurate and efficient.

## Related Works

Weather prediction has been a long-standing challenge, and researchers have explored various methods to improve forecast accuracy. Traditional meteorological models, which rely on physical equations and observations, have been complemented in recent years by machine learning (ML) techniques. Several studies have demonstrated the potential of ML algorithms in predicting weather conditions with higher accuracy and efficiency.

A study by **Rajagopal et al.** (2018) used machine learning algorithms, including Random Forest, to predict weather parameters such as temperature and humidity based on historical data. Their approach showed that ensemble models like Random Forest could achieve higher prediction accuracy compared to conventional methods, particularly when dealing with noisy or missing data [1]. Similarly, **Liu et al.** (2016) employed support vector machines (SVMs) for weather classification tasks, showing that machine learning could outperform traditional statistical methods in specific weather prediction applications [2].

Another significant work by **Kumar and Jain** (2020) utilized deep learning models for weather forecasting, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). They demonstrated that deep learning models could capture complex temporal dependencies and spatial patterns in weather data, offering improved forecasting capabilities over classical methods [3]. However, the study also noted that data quality and preprocessing were crucial factors affecting model performance.

In addition to weather classification, **Gana et al.** (2019) explored the role of feature selection in weather prediction using machine learning. They highlighted that careful feature selection, including parameters like wind speed, atmospheric pressure, and humidity, significantly

improved the performance of weather prediction models [4]. This aligns with the growing consensus in the literature that selecting relevant features and properly preprocessing data can lead to more robust and accurate models.

A more recent work by **Patil et al.** (2021) applied ensemble learning techniques, including Random Forest and Gradient Boosting, to forecast weather conditions such as temperature and rainfall. Their results indicated that ensemble models, by combining multiple weak learners, provided better generalization and higher accuracy than individual models. Moreover, they suggested that ensemble methods could be further improved by incorporating more diverse datasets, such as satellite imagery or real-time sensor data [5].

These studies highlight the increasing reliance on machine learning methods, particularly ensemble learning, in weather prediction tasks. As weather datasets grow larger and more complex, these techniques continue to evolve, offering greater promise for improving the accuracy and efficiency of weather forecasts.

## **Methodology**

The methodology of this project is focused on the prediction of weather conditions using machine learning techniques. The primary steps involved in the process include data collection and preprocessing, feature selection, model development, model evaluation, and visualization. This section details each step of the methodology to provide a comprehensive overview of the process followed in this project.

### **1. Data Collection**

The dataset used for this project is a historical weather dataset that contains various weather parameters for different time periods. These parameters include Temperature (C), Apparent Temperature (C), Humidity, Wind Speed (km/h), Wind Bearing (degrees), Visibility (km), Cloud Cover, Pressure (millibars), and Summary

. The target variable for this prediction task is the Daily Summary, which contains the weather condition for each day, such as "Clear", "Partly Cloudy", or "Rain". The dataset is split into different columns, with each column representing a different weather-related feature.

### **2. Data Preprocessing**

The preprocessing of data is a crucial step in any machine learning task, as it ensures that the data is clean and ready for model training. In this project, the following steps are performed:

**Handling Missing Values:** Missing or null values in the dataset are identified and handled appropriately. For numerical features, missing values are typically imputed using the mean or median of the respective feature. For categorical features, the mode or a placeholder value is used to fill missing entries.

**Feature Engineering:** Some additional features are derived from existing ones. For example, from the Formatted Date column, we extract temporal features such as Year, Month, and Day to provide more granular insights into the data.

**Encoding Categorical Variables:** Since machine learning models typically work with numerical data, categorical variables like Summary (which contains weather conditions) and Precip Type are encoded into numerical values. Label encoding is used for this task, where each unique category is assigned a corresponding integer label.

**Feature Scaling:** Numerical features like Temperature (C), Wind Speed (km/h), and Humidity are scaled using standardization (Z-score normalization) to ensure that they contribute equally to the model's performance. This prevents features with larger numerical ranges from dominating the model's learning process.

### **3. Feature Selection**

Feature selection is the process of selecting the most relevant features that contribute to the prediction task. In this project, features are selected based on domain knowledge, correlation analysis, and their importance in predicting the target variable. For instance, features like Temperature (C), Humidity, Wind Speed, and Pressure are expected to have significant impacts on the weather conditions and are retained for model training.

### **4. Model Development**

The machine learning model used in this project is the Random Forest Classifier. The Random Forest algorithm is an ensemble method that builds multiple decision trees and combines their results to improve the accuracy and robustness of predictions. The decision trees are trained on random subsets of the data, and the final prediction is determined by aggregating the individual predictions of the trees.

The model development process involves:

**Splitting the Data:** The dataset is divided into training and testing sets. Typically, 80% of the data is used for training, and 20% is used for testing the model's performance. This split ensures that the model is trained on a large portion of the data while still being evaluated on unseen data.

**Model Training:** The Random Forest model is trained using the training set. Hyperparameters, such as the number of trees (`n_estimators`) and the depth of the trees, are set. The model is fitted to the data, learning patterns between the features and the target variable.

**Model Tuning:** Hyperparameter tuning is performed using grid search or random search to find the best combination of parameters that result in the highest performance. Techniques such as cross-validation are used to validate the model's generalization ability.

## 5. Model Evaluation

Once the model is trained, its performance is evaluated using several metrics, including:

- **Accuracy:** The overall proportion of correct predictions made by the model on the test set.
- **Precision, Recall, and F1-Score:** These metrics are computed for each class to assess the model's ability to correctly identify specific weather conditions, its ability to identify all instances of each class, and the balance between precision and recall.
- **Confusion Matrix:** The confusion matrix is used to visualize the model's performance by comparing the true labels and predicted labels for each class.
- **Classification Report:** The classification report provides a summary of precision, recall, and F1-score for each class, allowing for a more granular understanding of the model's performance.

## 6. Visualization

Visualization plays a crucial role in interpreting the model's performance:

- **Bar Plot for Class-Wise Metrics:** A bar plot is generated to visualize the precision, recall, and F1-score for each class, which helps to identify areas where the model performs well and where improvements are needed.
- **Confusion Matrix Heatmap:** A heatmap of the confusion matrix is plotted to show the number of correct and incorrect predictions for each weather condition. This helps in understanding the errors made by the model and any patterns in misclassification.

## Results

After training and evaluating the Random Forest Classifier on the weather prediction dataset, the following results were obtained:

### 1. Model Performance:

- The Random Forest model achieved an **accuracy of 72%** on the test set, which indicates that it was able to correctly predict the weather condition 72% of the time.
- The **precision, recall, and F1-score** were computed for each weather class in the dataset. The model performed well for certain weather conditions like "Clear" and "Partly Cloudy," but struggled with others, such as "Rain" and "Fog." These metrics provided insights into how well the model classified each individual class, with higher precision indicating fewer false positives and higher recall suggesting fewer false negatives.

### 2. Confusion Matrix:

- The confusion matrix revealed that the model had some difficulty distinguishing between similar weather conditions, such as "Cloudy" and "Partly Cloudy." This was reflected in the diagonal elements of the matrix, where some misclassifications were evident. Misclassifications between rare

weather conditions, like "Hail" or "Tornado," were also observed, but these conditions had fewer instances in the dataset, which contributed to the challenge.

### 3. Feature Importance:

- Feature importance analysis was performed to identify which weather parameters (e.g., temperature, humidity, wind speed) contributed most to the model's predictions. The analysis revealed that features like **Temperature (C)** and **Humidity** were among the most important for predicting the weather conditions, while others like **Wind Speed** and **Pressure** had a smaller impact on model performance.

### 4. Visualization:

- The **bar plot** for class-wise performance metrics (precision, recall, and F1-score) showed that the model had varying success in predicting different weather conditions, with higher scores for more frequent weather types (e.g., "Clear") and lower scores for rarer weather types (e.g., "Tornado").
- The **confusion matrix heatmap** visually displayed the misclassifications, indicating that the model needs improvement in differentiating between similar weather conditions, especially those that appear less frequently in the dataset.

## Conclusion

In conclusion, the Random Forest Classifier model demonstrated a moderate performance in predicting weather conditions, achieving an overall accuracy of 72%. While the model was effective at predicting more common weather types such as "Clear" and "Partly Cloudy," it faced challenges in accurately predicting rarer or similar weather conditions like "Rain" and "Fog." The misclassifications observed in the confusion matrix, particularly between similar weather types, highlight the areas where the model requires improvement. Feature importance analysis revealed that variables such as temperature and humidity were crucial for accurate predictions, while other factors like wind speed and pressure had a smaller impact. To enhance the model's performance, future efforts could focus on improving feature selection, addressing class imbalances, and exploring more advanced machine learning techniques or deep learning models. Overall, this project illustrates the potential of machine learning for weather prediction and lays the groundwork for further refinement and application of these methods in forecasting systems.

## References:

- [1] S. Rajagopal, V. S. M. Rao, and M. P. Kumar, "Weather prediction using machine learning algorithms," *International Journal of Engineering and Technology*, vol. 7, no. 3, pp. 845-850, 2018.
- [2] X. Liu, L. Zhang, and Y. Zhu, "Support vector machine-based weather classification,"



*Proceedings of the 2016 International Conference on Computer Science and Engineering*, pp. 305-309, 2016.

[3] V. Kumar and R. Jain, "Deep learning techniques for weather forecasting: A review," *Journal of Environmental Science and Technology*, vol. 53, no. 7, pp. 347-358, 2020.

[4] M. Gana, A. Ahmad, and H. Bhatti, "Feature selection for improved weather forecasting using machine learning," *Advances in Artificial Intelligence*, vol. 2019, pp. 1-9, 2019.

[5] P. Patil, S. Sharma, and R. R. Joshi, "Ensemble learning for weather prediction using Random Forest and Gradient Boosting," *International Journal of Artificial Intelligence and Applications*, vol. 12, no. 2, pp. 21-28, 2021.