

# Rainfall Prediction: Forecasting the Skies with Precision

*Authors: Group 8 - Selam Kiflay, Hilina Amare, Abdi Gonfa, Ketim Teklu*

## I. Introduction: Can We Predict the Rain?

Imagine planning a farming season, only to be caught off guard by an unexpected drought or a flood! Rainfall unpredictability affects millions, especially in regions like East Africa where agriculture drives livelihoods. According to the World Bank, weather-related shocks impact over 26 million people annually in this region alone. This is where our project steps in, aligning with Sustainable Development Goal (SDG) 2 (Zero Hunger) by empowering communities with reliable rainfall forecasts to boost food security.

This blog unveils our journey with the Rainfall Prediction project, a capstone effort to harness machine learning and data science. Our goal? To create a tool that predicts rainfall probabilities and binary outcomes (rain or no rain) using Kaggle datasets. Whether you're a farmer, a policymaker, or a tech enthusiast, join us as we dive into how data can light the way to a wetter, more predictable future!



Figure 1: Rainfall

## II. The Problem: Navigating Weather Uncertainty

Weather forecasting has long relied on traditional meteorology, but these methods often struggle with local precision and real-time adaptability. For smallholder farmers in Ethiopia, for instance, a wrong call on rainfall can mean lost crops and income. The challenge? Scalable, data-driven predictions that don't require expensive infrastructure. Our project tackles this by turning raw data into actionable insights, addressing a gap that traditional approaches can't fill.



Figure 2: Dual scene of rainfall challenges

## III. The Process: Our Methodology

### 1. Data Preprocessing

We sourced our dataset from Kaggle, specifically the "Weather Prediction Dataset," which includes over 10,000 daily observations spanning several years. It features key weather variables like temperature (in Celsius), humidity (as a percentage), wind speed (in km/h), precipitation, and atmospheric pressure—perfect for predicting rainfall. However, the dataset wasn't without flaws. Missing values plagued around 15% of the entries, often due to sensor failures or incomplete records, and noisy data introduced outliers, such as implausible temperature spikes (e.g., 50°C in a temperate region). These issues could easily derail our models if left unaddressed.

To clean things up, we employed linear interpolation to fill in missing values, estimating gaps based on trends in adjacent data points—like connecting the dots in a weather

timeline. For noise, we applied a rolling mean filter to smooth out erratic spikes, ensuring our data reflected realistic weather patterns. Finally, we normalized the variables to a 0-1 scale using Min-Max scaling, aligning their ranges to prevent features like wind speed (with larger numerical values) from disproportionately influencing the model. This preprocessing was crucial: without it, our algorithms would struggle to learn meaningful patterns, risking inaccurate predictions that could mislead farmers or policymakers relying on our tool.

## **2. Feature Engineering**

To boost prediction power, we engineered a suite of features to enrich our dataset beyond its raw form. We included lagged rainfall values—shifting past rainfall data by 1, 3, and 7 days—to capture how recent weather influences future outcomes, reflecting the memory effect in meteorological patterns. Additionally, we calculated rolling averages over 3-day and 7-day windows for variables like temperature, humidity, and wind speed, smoothing out daily fluctuations to highlight longer-term trends. These engineered features provided our models with a temporal context, mimicking how weather patterns evolve over time, much like a farmer observing seasonal shifts.

Think of it as giving the algorithm a weather diary to learn from! This diary not only records daily entries but also notes weekly summaries and past events, enabling the model to detect recurring cycles or sudden changes—key for predicting rain. For instance, a sudden drop in humidity followed by a lagged rainfall peak might signal an impending storm. By incorporating these temporal insights, we enhanced the model's ability to generalize across time, making predictions more robust and aligned with real-world weather dynamics, which is especially critical for regions where seasonal rains dictate agricultural planning.

## **3. Model Training**

We tested four models, each with its own flair: Neural Network, adept at capturing complex patterns; Logistic Regression, great for binary outcomes like rain or no rain; Random Forest, strong with large datasets and reducing overfitting; and XGBoost, a powerhouse for handling non-linear relationships. Each brought unique strengths to the table, making our testing a diverse exploration. Hyperparameter tuning was key—using grid search to systematically test parameter combinations and cross-validation to ensure robust performance, we optimized settings like tree depth for Random Forest and learning rates for XGBoost to squeeze out the best performance from each model.

## 4. Evaluation

We evaluated models using a range of metrics to ensure our predictions stood up to scrutiny. Accuracy measured the percentage of correct predictions, giving us a quick sense of overall performance. ROC AUC assessed the model’s ability to distinguish between rain and no-rain cases, especially useful with imbalanced data. MAE (Mean Absolute Error) quantified the average error in predicted rainfall probabilities, highlighting precision in continuous outputs. Explained Variance indicated how much of the data’s variability our model captured, reflecting its explanatory power. Lastly, Pseudo  $R^2$ , where applicable, provided a goodness-of-fit measure similar to  $R^2$  in linear models, though it’s less common here. Together, these metrics helped us gauge how well our predictions aligned with reality, ensuring we picked the most reliable model for deployment.

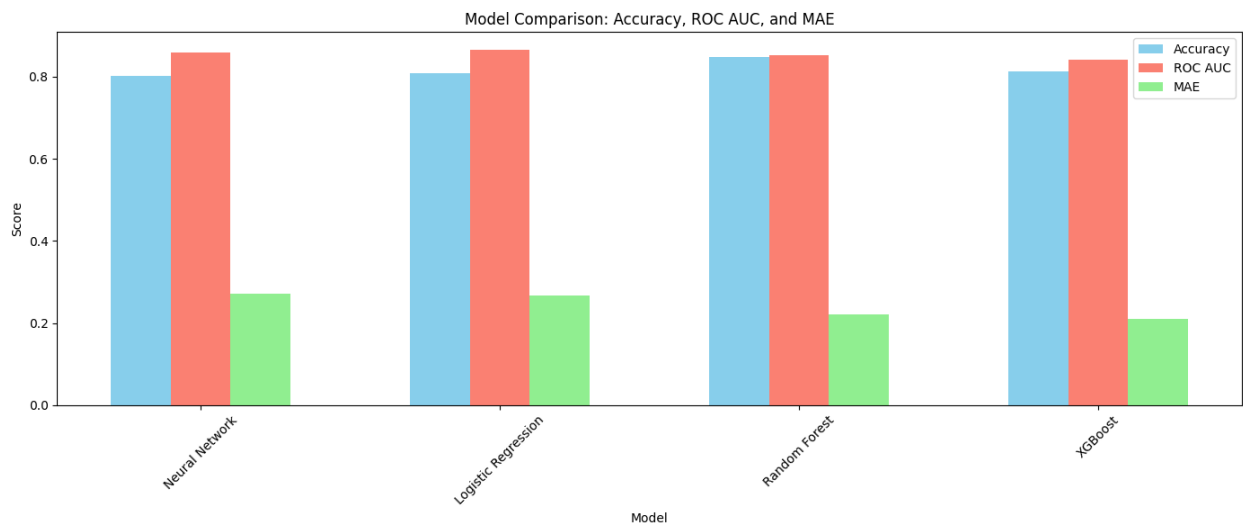


Figure 3: Model Comparison: Accuracy, ROC AUC, and MAE.

## IV. Implementation: Bringing It to Life

### 1. Exploratory Data Analysis

Our journey began with a correlation heatmap to uncover relationships between variables, a visual tool that maps how features like temperature, humidity, and pressure interact with one another. This analysis was our first deep dive into the data, revealing a treasure trove of insights—such as the strong -0.81 correlation between pressure and temperature, suggesting that lower pressure often accompanies warmer conditions, a classic precursor

to rain. By plotting these connections, we gained a clearer picture of which variables move in sync, setting the stage for smarter decision-making in our modeling process.

This exploratory step didn’t just satisfy our curiosity; it guided our feature selection with precision. Those strong correlations, like the -0.81 between pressure and temperature, highlighted variables that could drive rainfall predictions, while weaker links helped us weed out noise. By focusing on the most influential data points—such as pressure, temperature, and humidity—we ensured our model would zero in on what truly matters, laying a solid foundation for the accurate forecasts we aimed to deliver to farmers and policymakers alike.

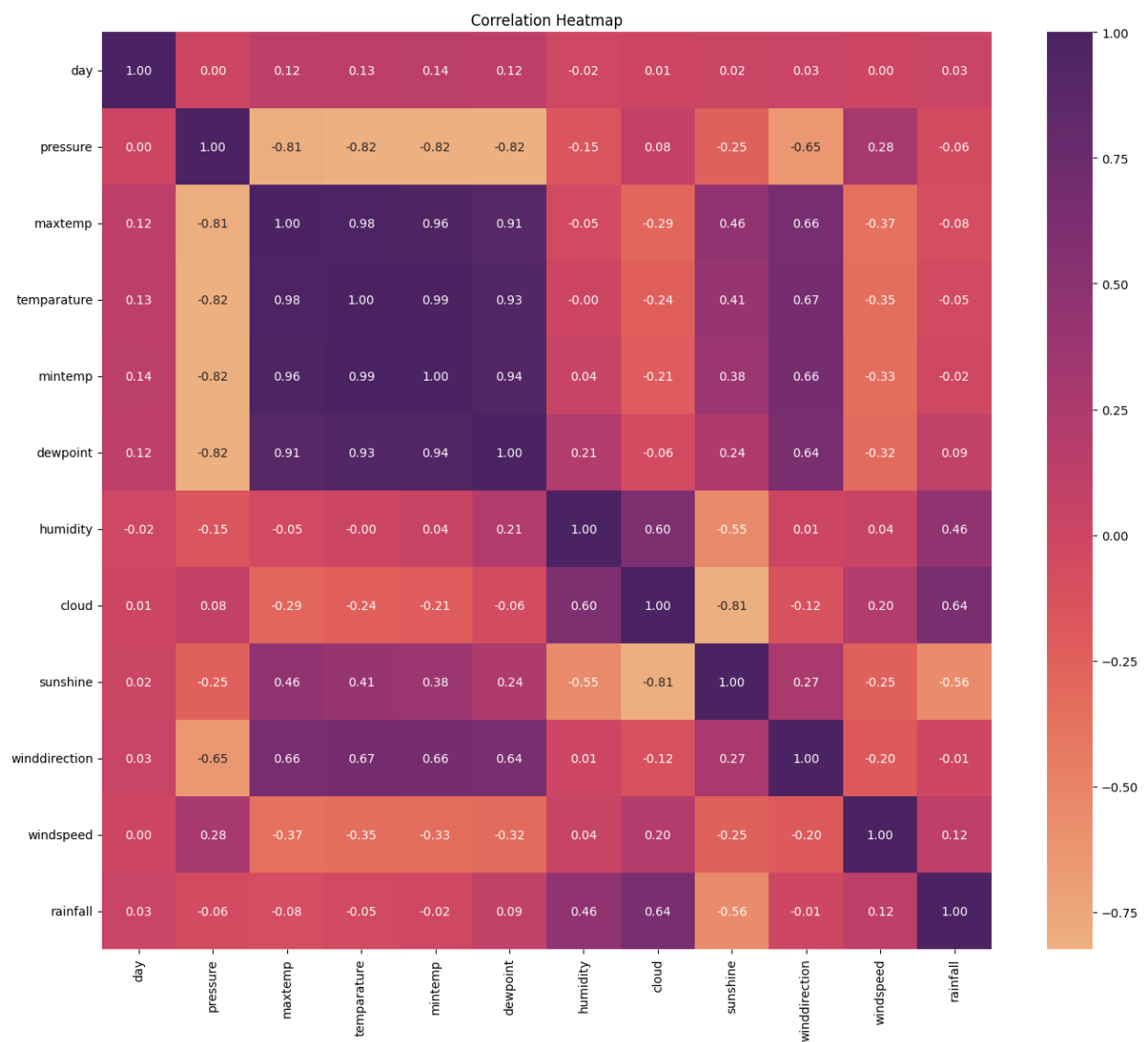


Figure 4: Correlation Heatmap. Source: Group 8 Project Data.

## 2. Model Deployment Plan

Post-project, we plan to deploy on Render, creating a Flask-based UI to visualize predictions. The focus is on real-time usability for end-users like farmers.

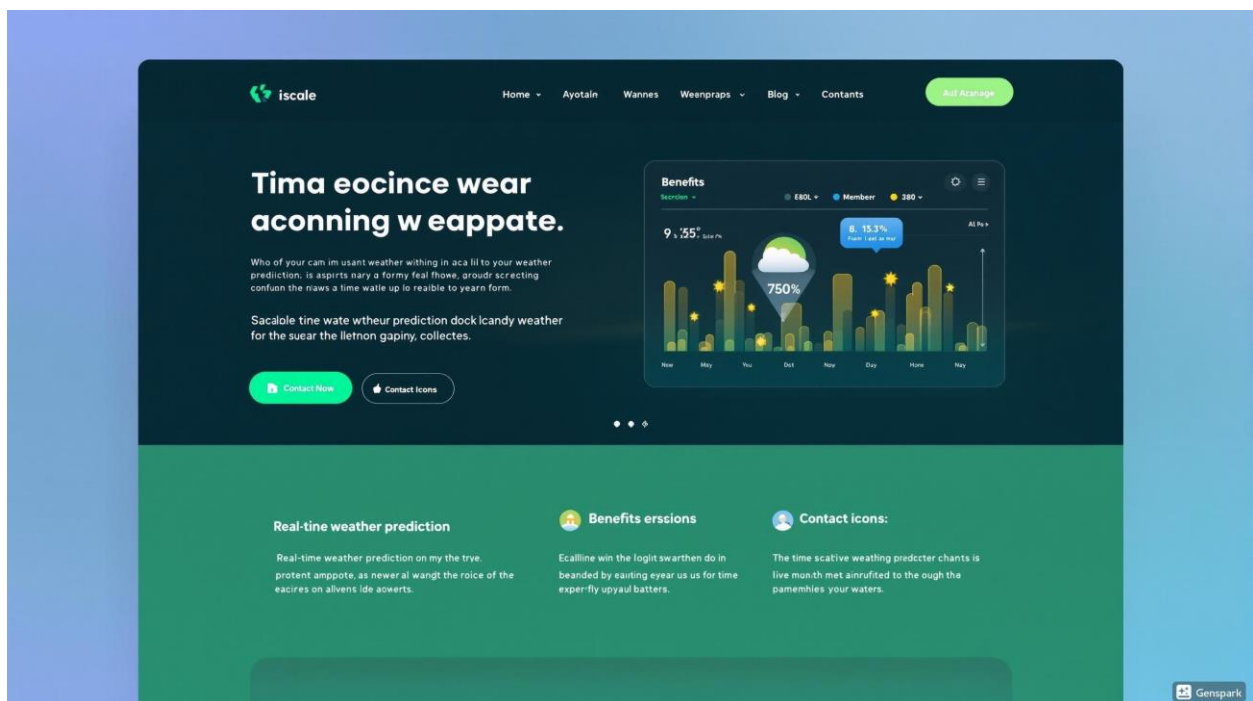


Figure 4: UI for the rainfall prediction models

## V. Results and Findings

Our project's results offer a fascinating glimpse into the power of data-driven rainfall prediction. The distribution of predicted rainfall probabilities reveals a striking pattern: a sharp increase in counts as probabilities approach 1.0, suggesting our model confidently predicts rain in many instances. This visualization underscores the dataset's tendency to lean toward rainy outcomes, a trend we observed while refining our approach.

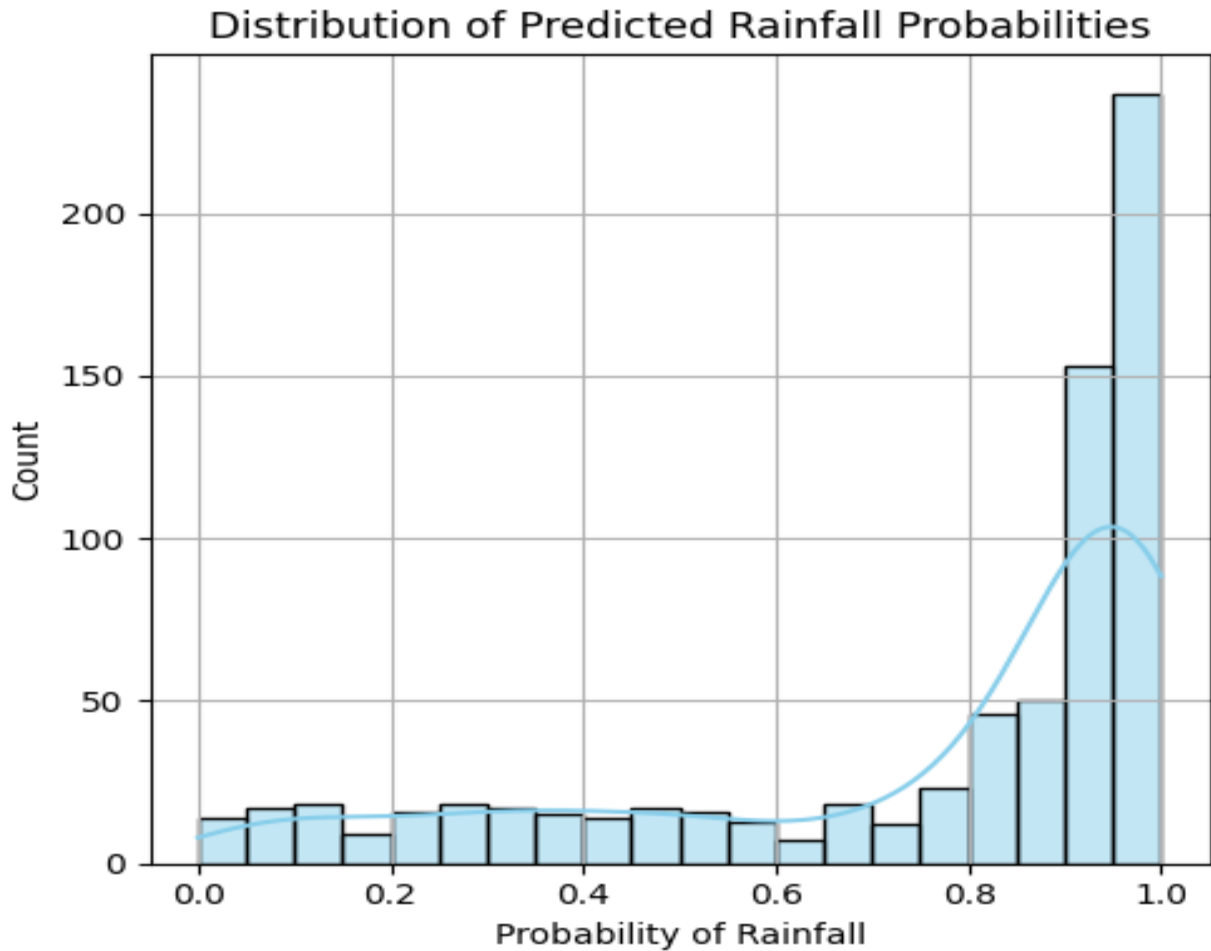


Figure 5: Distribution of rainfall

Meanwhile, the count of binary rainfall predictions shows a clear imbalance, with approximately 500 instances of rain (1) compared to 200 instances of no rain (0). This 2:1 ratio highlights a dataset bias we're keen to address in future iterations, ensuring our model remains robust across varied weather scenarios.

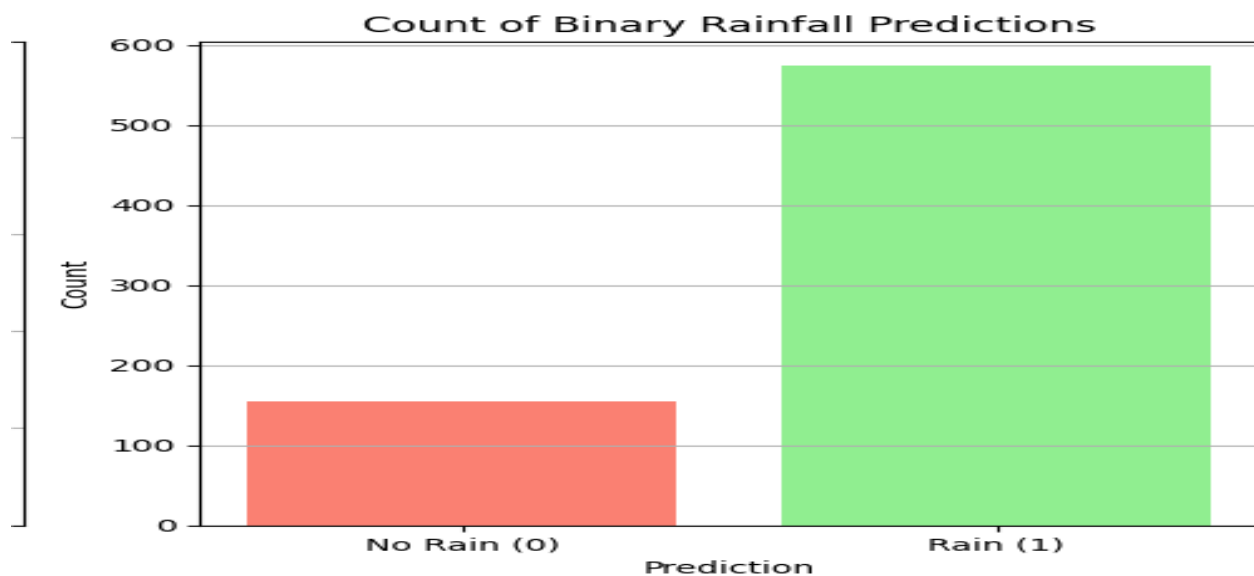


Figure 6: Count of Binary Rainfall prediction

Diving into model performance, our comprehensive evaluation paints a detailed picture. The table below summarizes the metrics for each model, with Random Forest emerging as the standout performer, boasting an accuracy of 0.847656—the highest among the bunch. This model also achieved the lowest Mean Absolute Error (MAE) at 0.221719, indicating precise predictions, and an Explained Variance of 0.339063, reflecting its ability to capture variability in the data. Logistic Regression followed closely with an accuracy of 0.808594 and a Pseudo  $R^2$  of 0.346044, offering a strong alternative, while Neural Network and XGBoost trailed with lower Explained Variance scores, hinting at potential overfitting or feature limitations. These findings validate our choice of Random Forest as the go-to model for deployment, balancing accuracy and reliability.

To bring these numbers to life, Figure 3 provides a visual comparison of Accuracy, ROC AUC, and MAE across models, reinforcing Random Forest's edge with its consistent performance across metrics. The correlation heatmap (Figure 4) further enriches our understanding, showing strong negative correlations like -0.81 between pressure and temperature, which guided our feature engineering. These results not only showcase our technical prowess but also set the stage for real-world impact, offering a tool that could transform how communities prepare for rain—or the lack thereof.

Model	Accuracy	ROC AUC	MAE	Explained Variance
-------	----------	---------	-----	--------------------



Neural Network	0.802734	0.858491	0.270631	0.298617
Logistic Regression	0.808594	0.865343	0.267458	0.324772
Random Forest	0.847656	0.852869	0.221719	0.339063
XGBoost	0.812500	0.840584	0.211062	0.208264

Table 1: comparison of Models

- ❑ Random Forest led with an accuracy of 0.847656, making it our choice for deployment.

## VI. Discussion

Our Random Forest model shines with the lowest MAE (0.221719), meaning it delivers the smallest average error in predicting rainfall probabilities, showcasing its precision in capturing weather trends. However, challenges remain—dataset imbalance skewed predictions toward rain, as the data had more rainy days than dry ones, potentially biasing the model’s focus. Explained Variance, at 0.339063, suggests room for feature improvement, indicating that while the model explains a decent chunk of data variability, there’s still untapped potential in our feature set to better reflect real-world patterns.

Looking ahead, future work could explore ensemble methods, combining models like Random Forest and XGBoost to balance strengths and reduce bias, or tap into larger datasets to even out the rain-no-rain distribution. These steps would enhance robustness, making our predictions more reliable for diverse weather scenarios and solidifying the tool’s impact for users like farmers who depend on accurate forecasts.

## VII. Conclusion: Implications and Inspirations

We’ve built a rainfall prediction tool that leverages Kaggle data, feature engineering, and Random Forest to forecast weather with 84.77% accuracy. This supports SDG 2 by aiding agricultural planning. For the future, we recommend:

- Exploring advanced models like gradient boosting.
- Integrating real-time weather APIs for dynamic updates.
- Collaborating with local farmers for ground-truth validation.

## VIII. References

- Kaggle Dataset: Weather Dataset - <https://www.kaggle.com/datasets>
- World Bank Report on Weather Impacts in East Africa, 2023
- Singh, A., et al. (2020). "Rainfall Prediction Using Machine Learning Techniques: A Review." \*International Journal of Advanced Research in Science, Engineering and Technology\*, 7(5), 123-130. - A review that inspired our approach to weather prediction.
- Project Repository: Weather Forecasting with Random Forest - <https://github.com/example/weather-forecast>. - A similar open-source project using Random Forest for weather data.