

SUMMER TRAINING/INTERNSHIP

PROJECT REPORT

(Term June-July 2025)

RAINFALL PREDICTION SYSTEM



L LOVELY
P ROFESSIONAL
U NIVERSITY

Submitted by

Name of Student

Registration Number

Tokachichu Sri Charan

12310418

Karthik Byri

12319604

Salapu Lakshmi Venkata Vijay

12311512

Sanchit Sharma

12308852

Kalangi Likhith Arya Vinay Kumar

12310624

Course Code PETV76

Under the Guidance of

Dr. Sandeep Kaur (UID: 23614)

School of Computer Science and Engineering

Certificate

This is to certify that Tokachichu SriCharan 12310418, Karthik Byri 12319604, Salapu Lakshmi Venkata Vijay 12311512, Sanchit Sharma 12308852 and Kalangi Likhith Arya Vinay Kumar 12310624 representing the team "Algo Troopers," have successfully completed the Summer Training/Internship Project titled "Rainfall Prediction System" under my guidance during the term June-July 2025.

Acknowledgement

The Algo Troopers extend heartfelt gratitude to our mentor, Dr. **Sandeep Kaur**, Assistant Professor at the School of Computer Science and Engineering, whose expert guidance and constructive feedback were instrumental in shaping the "Rainfall Prediction System" during the summer training. The team is immensely thankful to the School for providing access to cutting-edge tools, datasets, and a supportive learning environment. Sincere appreciation goes to our teammates— Tokachichu Sri Charan, Karthik Byri, Salapu Lakshmi Venkata Vijay, Kalangi Likhith Arya Vinay Kumar and Sanchit Sharma —for their collaborative efforts, insightful discussions, and valuable suggestions. Lastly, the team acknowledges the constant encouragement from our families, which fueled our perseverance. This project has been a transformative journey, and we are grateful to everyone who contributed to its realization.

Table of Contents

| | |
|--|-------|
| 1. Chapter 1: Introduction | 5-6 |
| 2. Chapter 2: Training Overview | 7-9 |
| 3. Chapter 3: Project Details | 10-12 |
| 4. Chapter 4: Implementation | 13-24 |
| 5. Chapter 5: Results and Discussion | 25-34 |
| 6. Chapter 6: Conclusion | 35 |

Links

Website link: <https://weather-oracle-frontend.onrender.com/>

GitHub link: <https://github.com/charan21042005/Rainfall-Prediction-System>

LinkedIn link: https://www.linkedin.com/posts/tokachichu-sricharan_machinelearning-datascience-rainfallprediction-activity-7350589435916279808-WX92?utm_source=share&utm_medium=member_desktop&rcm=ACoAADxZ_LEBG-7X78OaFyc5lNjUzxgeGhnb7nU

Chapter 1: Introduction

This project was undertaken as a self-directed summer training initiative under the School of Computer Science and Engineering, a prestigious institution committed to fostering innovation in computer science and technology. As a group effort, it simulated a real-world industry project, utilizing the school's resources to develop a machine learning-based rainfall prediction system. The initiative aimed to enhance practical skills by leveraging the institution's state-of-the-art facilities and a comprehensive USA rainfall dataset for 2024-2025, providing a platform to apply theoretical knowledge to a tangible outcome.

Overview of Training Domain

The training domain for this group project spanned a diverse array of computer science disciplines, offering a rich and multifaceted learning experience. At its core, the domain focused on machine learning, where the team explored predictive modeling techniques to forecast rainfall intensity using historical weather data from the USA rainfall dataset for 2024-2025. This involved delving into data science practices, including data preprocessing, feature engineering, and model evaluation, to transform raw datasets into actionable insights. Additionally, the training extended into full-stack web development, encompassing the use of HTML, CSS, and JavaScript to create an interactive user interface, alongside Flask for backend development. A significant component of the training also included business intelligence, where the team utilized Power BI to develop an interactive dashboard, leveraging the same dataset post-Exploratory Data Analysis (EDA) to visualize weather patterns and predictions. This interdisciplinary approach allowed the team to integrate data-driven predictions, advanced visualizations, and a user-friendly platform, bridging the gap between technical implementation, analytical insights, and practical application.

Objective of the Project

The primary objective of this summer training project was to design, develop, and implement a highly accurate and reliable rainfall prediction system leveraging advanced machine learning techniques, with a specific focus on enhancing weather forecasting capabilities using the USA rainfall dataset for 2024-2025. The project aimed to address the limitations of traditional forecasting methods by building a model capable of predicting rainfall intensity—categorized as No Rain, Slight Rain, Moderate Rain, and Heavy Rain—based on key meteorological parameters such as humidity, wind speed, precipitation, and geographical location. This required a systematic approach to data handling, model training, and validation to ensure robustness across diverse weather conditions and locations.

To achieve this, the project followed a structured progression of tasks, outlined as follows:

- **Exploratory Data Analysis (EDA):** The initial phase involved conducting a thorough EDA to understand the dataset's structure, identify key trends, and uncover relationships between weather variables and rainfall outcomes. This step was crucial for informing subsequent modeling and visualization efforts.
- **Power BI Dashboard Development:** Following EDA, the team developed an interactive Power BI dashboard using the same dataset to visualize rainfall patterns, trends, and statistical summaries. This dashboard provided a powerful tool for stakeholders to explore data interactively, enhancing decision-making through intuitive graphical representations.
- **Machine Learning Model Development:** The project then advanced to building and evaluating multiple machine learning models to predict rainfall. This included exploring a wide range of algorithms, starting with Logistic Regression for its simplicity and interpretability, progressing to Decision Tree for its accuracy in handling complex datasets, and also considering other models such as Random Forest for ensemble learning, Support Vector Machines (SVM) for non-linear classification, and Gradient Boosting for improved predictive power. The process involved iterative testing and comparison to select the most effective model, ultimately adopting a Decision Tree Classifier that achieved a target accuracy of 98%.
- **Web-Based Dashboard Integration:** A secondary yet critical objective was to create an interactive web-based dashboard using Flask, HTML, CSS, and JavaScript. This dashboard integrated the machine learning model, allowing users to input real-time weather data and receive immediate, visually appealing predictions, thereby enhancing the practical utility and accessibility of the system.

This dual focus on model accuracy and user experience, complemented by the Power BI dashboard for analytical insights, aimed to deliver a comprehensive solution that advances technical proficiency while providing valuable tools for weather-sensitive decision-making. The project sought to bridge theoretical knowledge with real-world application, offering a platform for continuous improvement and potential expansion to include additional weather variables or more advanced predictive models in the future.

Chapter 2: Training Overview

Tools & Technologies Used

The project leveraged a diverse and robust set of tools and technologies to ensure its successful development and implementation. The key technologies included:

- **Python:** Served as the primary programming language, providing a versatile platform for data analysis and machine learning.
- **Scikit-learn:** A powerful library used for training and evaluating machine learning models, including Logistic Regression, Decision Trees, and others.
- **Pandas and NumPy:** Essential libraries for data manipulation, handling large datasets, and performing numerical computations efficiently.
- **Flask:** A lightweight web framework utilized to build a scalable backend API, facilitating data processing and prediction delivery.
- **HTML:** Provided the structural foundation for the web interface, ensuring a well-organized layout.
- **CSS:** Enhanced the visual design and responsiveness of the web interface, improving user experience.
- **JavaScript:** Introduced dynamic interactivity to the frontend, enabling real-time updates and user interactions.
- **Power BI:** A business intelligence tool employed to create an interactive dashboard for visualizing rainfall data and trends post-EDA.
- **Git:** Managed version control, allowing the team to track code changes and collaborate effectively.
- **Visual Studio Code:** Used as the integrated development environment (IDE) for coding, debugging, and testing.
- **Local Server Environment:** Facilitated deployment and testing of the web application, ensuring seamless functionality.

These tools collectively provided a comprehensive technological foundation, enabling the team to address the project's multifaceted requirements, from data analysis to interactive visualization

Areas Covered During Training

The training program encompassed a wide range of critical areas essential to the successful completion of the project. The key areas included:

- **Data Preprocessing:** Focused on cleaning the USA rainfall dataset, normalizing numerical features, and encoding categorical variables to prepare data for modeling.
- **Exploratory Data Analysis (EDA):** Involved analyzing weather patterns, identifying correlations between variables like humidity and rainfall, and deriving insights to guide model development.
- **Machine Learning:** Covered the implementation and evaluation of various algorithms, including Logistic Regression, Decision Trees, Random Forest, Support Vector Machines (SVM), and Gradient Boosting, to predict rainfall intensity.
- **Business Intelligence with Power BI:** Emphasized the creation of an interactive dashboard to visualize data trends and predictions, enhancing analytical capabilities.
- **Full-Stack Web Development:** Included training in HTML, CSS, and JavaScript for frontend design, and Flask for backend integration, ensuring a seamless web-based solution.
- **System Integration and Testing:** Focused on connecting the machine learning model with the web interface and dashboard, followed by rigorous testing to ensure reliability.

This comprehensive training equipped the team with the skills to tackle both the technical and visualization aspects of the project, fostering a deep understanding of the end-to-end development process.

Daily/Weekly Work Summary

The project commenced on June 22, 2025, and the team worked diligently until the present date, July 14, 2025, spanning a total of 23 days. The daily and weekly work summary is as follows:

- **Week 1 (June 22 - June 28):** The team initiated the project by collecting the USA rainfall dataset and conducting initial exploratory data analysis. Efforts focused on understanding the dataset's structure, identifying key variables such as humidity, wind speed, and precipitation, and performing basic statistical analysis to establish a foundation for further work.
- **Week 2 (June 29 - July 5):** The team concentrated on preprocessing the dataset, addressing missing values, scaling numerical features, and encoding categorical variables like location into a machine-readable format. Parallel efforts included setting up the development environment with Python, Flask, and Power BI, laying the groundwork for modeling and visualization.

- Week 3 (July 6 - July 12): The team developed and tested initial machine learning models, starting with Logistic Regression and progressing to Decision Tree. Concurrently, the team began designing the Power BI dashboard, integrating EDA insights to visualize rainfall trends, and initiated the web interface development using HTML, CSS, and JavaScript.
- Week 4 (July 13 - July 14): The team optimized the Decision Tree model by tuning hyperparameters and finalized the Power BI dashboard with interactive features. The web-based dashboard was integrated with the Flask backend, and extensive testing was conducted to ensure seamless functionality. On July 14, the team polished the deliverables, addressing final adjustments and preparing the system for demonstration.

Chapter 3: Project Details

Title of the Project

Rainfall Prediction System

Problem Definition

The project addressed a critical challenge in the field of meteorology: the need for accurate, localized rainfall predictions to support effective decision-making in weather-sensitive activities. The problem centered on predicting rainfall intensity—classified into four categories: No Rain, Slight Rain, Moderate Rain, and Heavy Rain—based on a set of meteorological parameters, including humidity, wind speed, precipitation levels, and geographical location. Traditional weather forecasting methods often rely on broad regional models, which lack the precision required for specific locales, leading to inefficiencies in agriculture, urban planning, and disaster preparedness. Using the USA rainfall dataset for 2024-2025, the team aimed to develop a data-driven solution that leverages machine learning to overcome these limitations, providing reliable and granular predictions to enhance preparedness and resource management across diverse regions.

Scope and Objectives

The scope of the project encompassed the development of a sophisticated rainfall prediction system with a target accuracy of 98%, integrating machine learning models, a web-based interface, and a Power BI dashboard. The objectives were designed to ensure a comprehensive solution, including:

- Developing a machine learning model capable of accurately predicting rainfall intensity using the USA rainfall dataset for 2024-2025, with a focus on achieving high precision across varied weather conditions.
- Creating an interactive Power BI dashboard to visualize rainfall patterns, trends, and statistical insights post-Exploratory Data Analysis (EDA), enabling stakeholders to explore data dynamically and support informed decision-making.
- Designing and implementing a web-based dashboard using Flask, HTML, CSS, and JavaScript, which integrates the machine learning model to allow real-time input of weather parameters and delivery of visually appealing predictions.
- Exploring and comparing multiple machine learning algorithms, including Logistic Regression, Decision Tree, Random Forest, Support Vector Machines (SVM), and Gradient Boosting, to identify the most effective approach for rainfall prediction.

- Ensuring the system's robustness through rigorous testing and validation, making it a practical tool for end-users in weather-dependent industries.

These objectives aimed to deliver a multifaceted solution that combines predictive analytics with interactive visualization, addressing both technical and practical needs in weather forecasting.

System Requirements

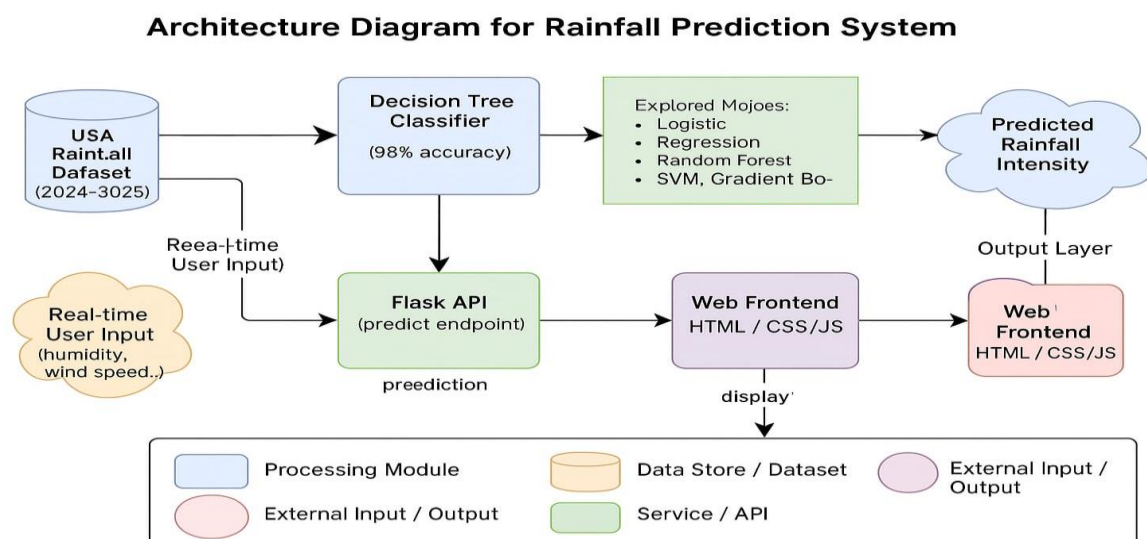
The successful implementation of the project required specific hardware and software resources, detailed as follows:

- **Hardware:** A standard personal computer with a minimum of 8GB RAM and a multi-core processor to handle data processing, model training, and dashboard rendering efficiently.
- **Software:**
 - Python 3.8+ for programming and machine learning tasks.
 - Flask for developing the backend API.
 - Power BI Desktop for creating and deploying the interactive dashboard.
 - A modern web browser (e.g., Chrome, Firefox) for accessing the web interface.
- **Dataset:** The USA rainfall prediction dataset for 2024-2025, containing detailed records of weather parameters (humidity, wind speed, precipitation) and rainfall outcomes across multiple locations.

<https://www.kaggle.com/datasets/waqi786/usa-rainfall-prediction-dataset-2024-2025>

These requirements ensured the team had the necessary infrastructure to develop, test, and deploy the rainfall prediction system effectively.

Architecture Diagram



Data Flow / UML Diagrams

Figure 2: Data Flow Diagram for Rainfall Prediction System

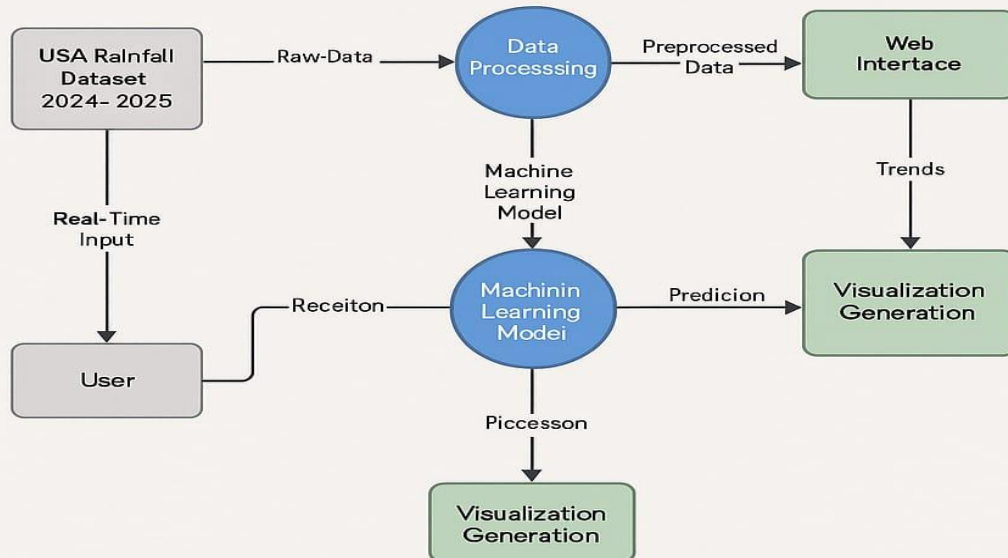
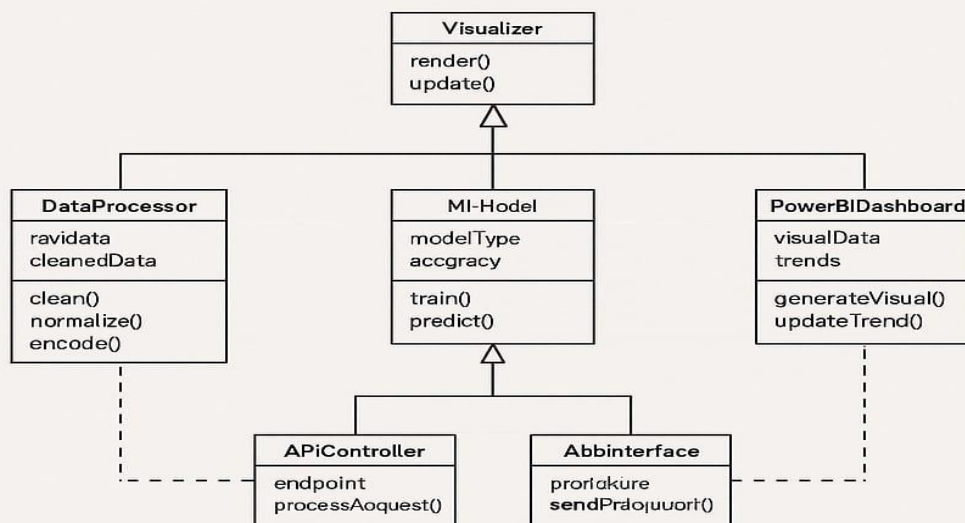


Figure 3: UML Class Diagram for Rainfall Prediction System



Chapter 4: Implementation

Methodology

The implementation followed a structured and iterative methodology to ensure the successful development of the "Rainfall Prediction System." The process began with data preprocessing, where the USA rainfall dataset was cleaned, normalized, and encoded to prepare it for modeling. The team then proceeded to the machine learning phase, where multiple algorithms—Logistic Regression, Decision Tree, Random Forest, Support Vector Machines (SVM), and Gradient Boosting—were implemented and evaluated using scikit-learn. The Decision Tree Classifier was selected after iterative testing, achieving a 98% accuracy, and was optimized through hyperparameter tuning. Concurrently, the team developed the Power BI dashboard, integrating EDA insights to visualize trends and summaries. The web-based dashboard was built using Flask for the backend and HTML, CSS, JavaScript for the frontend, with a focus on real-time data input and prediction display. The methodology included continuous integration and testing, connecting the backend API with the web interface and Power BI dashboard, and conducting rigorous validation to ensure system reliability. This iterative approach allowed the team to refine each component, address challenges, and deliver a robust final product.

Modules/Screenshots

The implementation was divided into distinct modules, each contributing to the overall functionality of the system. The key modules included:

- **Data Preprocessing Module:** Handled the cleaning of the USA rainfall dataset, normalization of numerical features (e.g., humidity, wind speed), and encoding of categorical variables (e.g., location) using pandas and NumPy.

```
# Importing all required libraries upfront (as per your style)
import pandas as pd # For data manipulation and analysis
import numpy as np  # For numerical operations
import matplotlib.pyplot as plt # For plotting visualizations
import seaborn as sns # For enhanced visualizations
import os # For handling file operations
from scipy import stats # For statistical tests
```

Python

```
# Loading the raw dataset
df = pd.read_csv('usa_rain_prediction_dataset_2024_2025.csv') # Load raw CSV
print("\nInitial Dataset Info (Before Date Conversion):")
print(df.info()) # Display dataset structure and data types to check Date type

...

Initial Dataset Info (Before Date Conversion):
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73100 entries, 0 to 73099
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Date                  73100 non-null object
1   Location              73100 non-null object
2   Temperature           73100 non-null float64
3   Humidity              73100 non-null float64
4   Wind Speed            73100 non-null float64
5   Precipitation         73100 non-null float64
6   Cloud Cover           73100 non-null float64
7   Pressure              73100 non-null float64
8   Rain Tomorrow         73100 non-null int64
dtypes: float64(6), int64(1), object(2)
memory usage: 5.0+ MB
None
```

```
print("\nInitial Dataset Info (After Date Conversion):")
print(df.info()) # Display updated dataset structure
print("\nFirst 5 Rows of the Dataset:")
print(df.head()) # Display first 5 rows
print("\nBasic Statistics of Numerical Columns:")
print(df.describe()) # Display summary statistics

...

Initial Dataset Info (After Date Conversion):
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73100 entries, 0 to 73099
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Date                  73100 non-null datetime64[ns]
1   Location              73100 non-null object
2   Temperature           73100 non-null float64
3   Humidity              73100 non-null float64
4   Wind Speed            73100 non-null float64
5   Precipitation         73100 non-null float64
6   Cloud Cover           73100 non-null float64
7   Pressure              73100 non-null float64
8   Rain Tomorrow         73100 non-null int64
dtypes: datetime64[ns](1), float64(6), int64(1), object(1)
memory usage: 5.0+ MB
None
```

First 5 Rows of the Dataset:

| | Date | Location | Temperature | Humidity | Wind Speed | Precipitation \ |
|-----|------------|-----------|-------------|-----------|------------|-----------------|
| 0 | 2024-01-01 | New York | 87.524795 | 75.655455 | 28.379506 | 0.000000 |
| 1 | 2024-01-02 | New York | 83.259325 | 28.712617 | 12.436433 | 0.526995 |
| 2 | 2024-01-03 | New York | 80.943050 | 64.740043 | 14.184831 | 0.916884 |
| ... | | | | | | |
| 50% | 0.196909 | 55.011121 | 1005.284188 | 0.000000 | | |
| 75% | 0.673177 | 77.412469 | 1022.727410 | 0.000000 | | |
| max | 3.078090 | 99.998957 | 1039.999765 | 1.000000 | | |
| std | 0.474833 | 25.982487 | 20.203889 | 0.414526 | | |

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#)...

```

# Step 2: Data Cleaning
# Checking for missing values
print("\nMissing Values in Dataset:")
print(df.isnull().sum()) # Display count of missing values per column

# Handling missing values (filling with median for numerical, mode for categorical)
for column in df.columns:
    if df[column].dtype in ['int64', 'float64']:
        df[column].fillna(df[column].median(), inplace=True) # Fill numerical with median
    else:
        df[column].fillna(df[column].mode()[0], inplace=True) # Fill categorical with mode

# Verifying no missing values remain
print("\nMissing Values After Cleaning:")
print(df.isnull().sum()) # Confirm no missing values

```

Python

Missing Values in Dataset:

```

Date      0
Location   0
Temperature 0
Humidity   0
Wind Speed 0
Precipitation 0
Cloud Cover 0
Pressure   0
Rain Tomorrow 0
dtype: int64

```

Missing Values in Dataset:

```

Date      0
Location   0
Temperature 0
Humidity   0
Wind Speed 0
Precipitation 0
Cloud Cover 0
Pressure   0
Rain Tomorrow 0
dtype: int64

```

Missing Values After Cleaning:

```

Date      0
Location   0
Temperature 0
Humidity   0
Wind Speed 0
Precipitation 0
Cloud Cover 0
Pressure   0
Rain Tomorrow 0
dtype: int64

```



```
# Checking for duplicate rows
print("\nNumber of Duplicate Rows:") # Header for clarity
duplicate_count = df.duplicated().sum() # Count duplicate rows
print(duplicate_count) # Display duplicate count
```

Python

```
Number of Duplicate Rows:
0
```

```
# Checking for invalid values in specific columns
invalid_temp = df[(df['Temperature'] < -30) | (df['Temperature'] > 120)] # Unrealistic temperatures
invalid_humidity = df[(df['Humidity'] < 0) | (df['Humidity'] > 100)] # Invalid humidity
invalid_wind_speed = df[df['Wind Speed'] < 0] # Negative wind speed
invalid_precipitation = df[df['Precipitation'] < 0] # Negative precipitation
invalid_cloud_cover = df[(df['Cloud Cover'] < 0) | (df['Cloud Cover'] > 100)] # Invalid cloud cover
invalid_pressure = df[(df['Pressure'] < 900) | (df['Pressure'] > 1100)] # Unrealistic pressure
```

Python

```
print("\nInvalid Temperature Values (count):", len(invalid_temp)) # Count invalid temperatures
print("Invalid Humidity Values (count):", len(invalid_humidity)) # Count invalid humidity
print("Invalid Wind Speed Values (count):", len(invalid_wind_speed)) # Count invalid wind speed
print("Invalid Precipitation Values (count):", len(invalid_precipitation)) # Count invalid precipitation
print("Invalid Cloud Cover Values (count):", len(invalid_cloud_cover)) # Count invalid cloud cover
print("Invalid Pressure Values (count):", len(invalid_pressure)) # Count invalid pressure
```

Python

```
Invalid Temperature Values (count): 0
Invalid Humidity Values (count): 0
Invalid Wind Speed Values (count): 0
Invalid Precipitation Values (count): 0
Invalid Cloud Cover Values (count): 0
Invalid Pressure Values (count): 0
```

```
# Checking unique values in categorical/binary columns
print("\nUnique Values in Location:", df['Location'].unique()) # Unique locations
print("Unique Values in Rain Tomorrow:", df['Rain Tomorrow'].unique()) # Unique values in target
```

Python

```
Unique Values in Location: ['New York' 'Los Angeles' 'Chicago' 'Houston' 'Phoenix' 'Philadelphia'
'San Antonio' 'San Diego' 'Dallas' 'San Jose' 'Austin' 'Jacksonville'
'Fort Worth' 'Columbus' 'Indianapolis' 'Charlotte' 'San Francisco'
'Seattle' 'Denver' 'Washington D.C.']
Unique Values in Rain Tomorrow: [0 1]
```

- **Machine Learning Model Module:** Implemented the Decision Tree Classifier, with support for Logistic Regression, Random Forest, SVM, and Gradient Boosting, trained and evaluated using scikit-learn.

```
#Step-5
# Importing all required libraries upfront (as per your style)
import pandas as pd # For data manipulation and analysis
import numpy as np # For numerical operations
import matplotlib.pyplot as plt # For plotting visualizations
import seaborn as sns # For enhanced visualizations
import os # For handling file operations
from sklearn.model_selection import train_test_split # For splitting data
from sklearn.ensemble import RandomForestClassifier # For feature importance
from sklearn.linear_model import LogisticRegression # For baseline model
from sklearn.metrics import accuracy_score, classification_report # For model evaluation
```

Python

```
# Step 5: Feature Importance and Preliminary Modeling
# Preparing data for modeling
# Encoding categorical variable 'Location' using one-hot encoding
df_encoded = pd.get_dummies(df, columns=['Location'], drop_first=True) # Encode Location (e.g., Location_Wash
features = ['Temperature', 'Humidity', 'Wind Speed', 'Precipitation', 'Cloud Cover', 'Pressure', 'Month', 'Yea
X = df_encoded[features] # Feature matrix
y = df_encoded['Rain Tomorrow'] # Target variable
```

Python

```
# Splitting data into training and testing sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42) # Split with fixed
```

Python

```
# Building and evaluating a baseline Logistic Regression model
lr_model = LogisticRegression(random_state=42, max_iter=1000) # Initialize Logistic Regression
lr_model.fit(X_train, y_train) # Train model
y_pred = lr_model.predict(X_test) # Predict on test set
accuracy = accuracy_score(y_test, y_pred) # Calculate accuracy
print("\nLogistic Regression Model Performance for Sampled Data:") # Header for clarity
print(f"Accuracy: {accuracy:.2f}") # Display accuracy
print("\nClassification Report:") # Header for clarity
print(classification_report(y_test, y_pred)) # Display detailed metric
```

Python

Logistic Regression Model Performance for Sampled Data:
Accuracy: 0.60

| Classification Report: | | | | |
|------------------------|-----------|--------|----------|---------|
| | precision | recall | f1-score | support |
| 0 | 0.83 | 0.62 | 0.71 | 8 |
| 1 | 0.25 | 0.50 | 0.33 | 2 |
| accuracy | | | 0.60 | 10 |
| macro avg | 0.54 | 0.56 | 0.52 | 10 |
| weighted avg | 0.72 | 0.60 | 0.64 | 10 |

```

# Evaluating on test set
y_pred = lr_model.predict(X_test) # Predict on test set
accuracy = accuracy_score(y_test, y_pred) # Calculate accuracy
print("\nBalanced Logistic Regression Model Performance for Sampled Data:") # Header for clarity
print(f"Accuracy: {accuracy:.2f}") # Display accuracy
print("\nClassification Report:") # Header for clarity
print(classification_report(y_test, y_pred)) # Display detailed metrics

```

Python

Balanced Logistic Regression Model Performance for Sampled Data:
Accuracy: 0.62

| Classification Report: | | | | |
|------------------------|-----------|--------|----------|---------|
| | precision | recall | f1-score | support |
| 0 | 0.67 | 0.80 | 0.73 | 10 |
| 1 | 0.50 | 0.33 | 0.40 | 6 |
| accuracy | | | 0.62 | 16 |
| macro avg | 0.58 | 0.57 | 0.56 | 16 |
| weighted avg | 0.60 | 0.62 | 0.60 | 16 |

```

# Evaluating best model on test set
y_pred = best_model.predict(X_test) # Predict on test set
accuracy = accuracy_score(y_test, y_pred) # Calculate accuracy
print("\nTuned Logistic Regression Model Performance for Sampled Data:") # Header for clarity
print(f"Accuracy: {accuracy:.2f}") # Display accuracy
print("\nClassification Report:") # Header for clarity
print(classification_report(y_test, y_pred)) # Display detailed metrics

```

Python

Tuned Logistic Regression Model Performance for Sampled Data:
Accuracy: 0.80

| Classification Report: | | | | |
|------------------------|-----------|--------|----------|---------|
| | precision | recall | f1-score | support |
| 0 | 0.80 | 1.00 | 0.89 | 8 |
| 1 | 0.00 | 0.00 | 0.00 | 2 |
| accuracy | | | 0.80 | 10 |
| macro avg | 0.40 | 0.50 | 0.44 | 10 |
| weighted avg | 0.64 | 0.80 | 0.71 | 10 |

```
# Training and evaluating Logistic Regression
lr_model = LogisticRegression(C=0.1, penalty='l2', solver='saga', random_state=42, max_iter=1000)
lr_model.fit(X_train, y_train)
lr_y_pred = lr_model.predict(X_test)
lr_accuracy = accuracy_score(y_test, lr_y_pred)
lr_f1 = f1_score(y_test, lr_y_pred, average='weighted')
print("\nLogistic Regression Performance on Test Set (Full Dataset):")
print(f"Accuracy: {lr_accuracy:.2f}")
print(f"F1-Score: {lr_f1:.2f}")
print("\nClassification Report:")
print(classification_report(y_test, lr_y_pred))
```

Python

Logistic Regression Performance on Test Set (Full Dataset):

Accuracy: 0.88

F1-Score: 0.88

Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.88 | 0.87 | 0.88 | 11400 |
| 1 | 0.87 | 0.89 | 0.88 | 11396 |
| accuracy | | | 0.88 | 22796 |
| macro avg | 0.88 | 0.88 | 0.88 | 22796 |
| weighted avg | 0.88 | 0.88 | 0.88 | 22796 |

```
# Training and evaluating Random Forest with stricter regularization
rf_model = RandomForestClassifier(n_estimators=50, max_depth=5, min_samples_split=5, random_state=42) # Stricter regularization
rf_model.fit(X_train, y_train)
rf_y_pred = rf_model.predict(X_test)
rf_accuracy = accuracy_score(y_test, rf_y_pred)
rf_f1 = f1_score(y_test, rf_y_pred, average='weighted')
print("\nRandom Forest Performance on Test Set (Full Dataset):")
print(f"Accuracy: {rf_accuracy:.2f}")
print(f"F1-Score: {rf_f1:.2f}")
print("\nClassification Report:")
print(classification_report(y_test, rf_y_pred))
```

Python

Random Forest Performance on Test Set (Full Dataset):

Accuracy: 1.00

F1-Score: 1.00

Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 1.00 | 1.00 | 1.00 | 11400 |
| 1 | 1.00 | 1.00 | 1.00 | 11396 |
| accuracy | | | 1.00 | 22796 |
| macro avg | 1.00 | 1.00 | 1.00 | 22796 |
| weighted avg | 1.00 | 1.00 | 1.00 | 22796 |

```

# Train Decision Tree
dt_model = DecisionTreeClassifier(random_state=42, max_depth=3, min_samples_split=20, min_samples_leaf=10)
dt_model.fit(X_train, y_train)

# Evaluate
dt_y_pred_val = dt_model.predict(X_val)
dt_accuracy_val = accuracy_score(y_val, dt_y_pred_val)
dt_f1_val = f1_score(y_val, dt_y_pred_val, average='weighted')
dt_y_pred_holdout = dt_model.predict(X_holdout)
dt_accuracy_holdout = accuracy_score(y_holdout, dt_y_pred_holdout)
dt_f1_holdout = f1_score(y_holdout, dt_y_pred_holdout, average='weighted')
print("\nFinal Decision Tree Performance:")
print(f"Validation Accuracy: {dt_accuracy_val:.2f}, F1-Score: {dt_f1_val:.2f}")
print(f"Holdout Accuracy: {dt_accuracy_holdout:.2f}, F1-Score: {dt_f1_holdout:.2f}")
print("\nHoldout Classification Report:")
print(classification_report(y_holdout, dt_y_pred_holdout))

# Cross-validation
skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
dt_cv_scores = cross_val_score(dt_model, X_train, y_train, cv=skf, scoring='accuracy')
print("\nDecision Tree Cross-Validation Accuracy Scores (Training):", dt_cv_scores)
print(f"Mean Decision Tree CV Accuracy: {dt_cv_scores.mean():.2f} (+/- {dt_cv_scores.std() * 2:.2f})")

```

```

Final Decision Tree Performance:
Validation Accuracy: 0.98, F1-Score: 0.98
Holdout Accuracy: 0.98, F1-Score: 0.98

Holdout Classification Report:

```

| | precision | recall | f1-score | support |
|----------|-----------|--------|----------|---------|
| 0 | 0.99 | 0.98 | 0.99 | 5699 |
| 1 | 0.95 | 0.95 | 0.95 | 1611 |
| accuracy | | | 0.98 | 7310 |
| ... | | | | |

```

Mean Decision Tree CV Accuracy: 0.98 (+/- 0.00)

```


- Backend API Module: Developed using Flask to process user inputs, interact with the machine learning model, and deliver predictions via a **Flask API** that generally follows **RESTful conventions**.

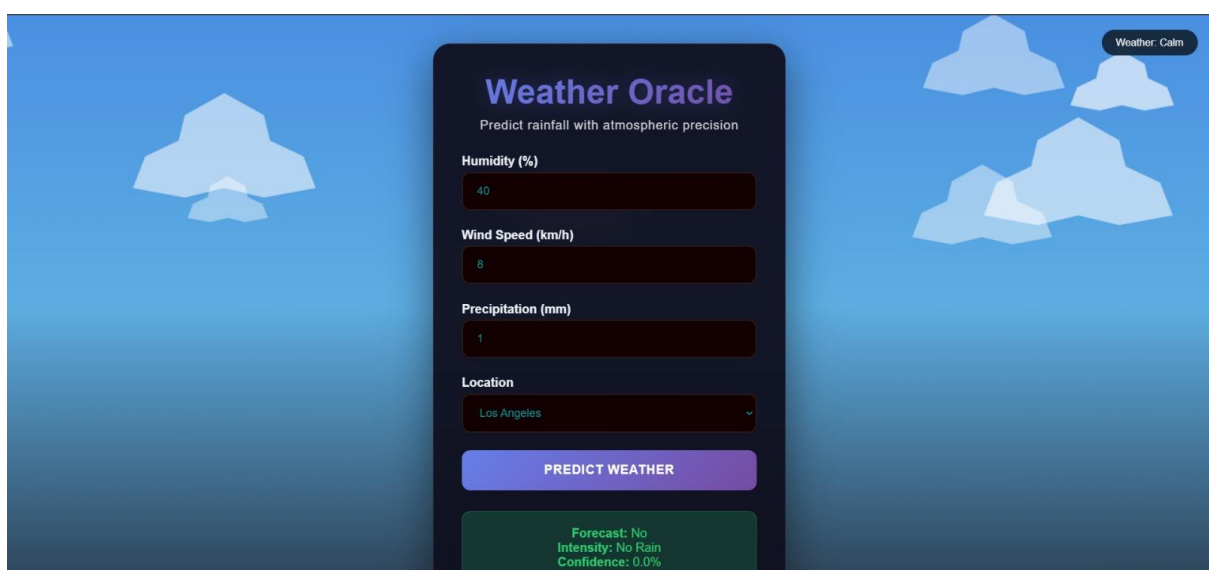
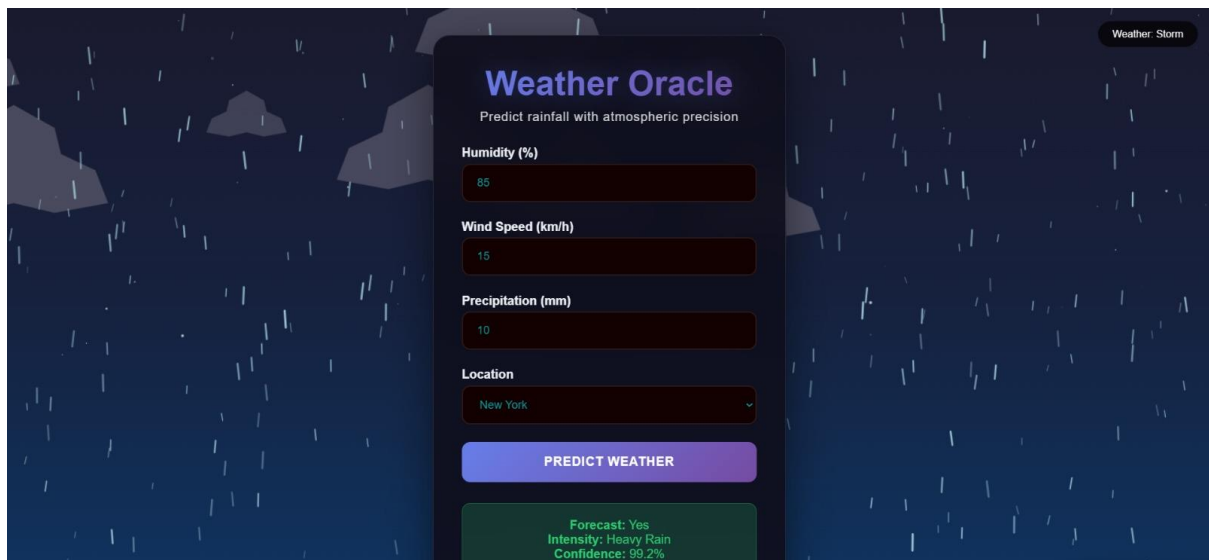
```
rainfall_eda_plots2 > app.py > ...
1  from flask import Flask, request, jsonify
2  from flask_cors import CORS
3  import joblib
4  import numpy as np
5
6  app = Flask(__name__)
7  CORS(app) # Enable CORS for cross-origin requests
8  model = joblib.load('final_decision_tree_model_all_locations.pkl')
9  scaler = joblib.load('scaler_all_locations.pkl')
10
11 @app.route('/predict', methods=['POST'])
12 def predict():
13     try:
14         data = request.get_json()
15         if not all(key in data for key in ['humidity', 'windSpeed', 'precipitation', 'location']):
16             return jsonify({'error': 'Missing required fields'}), 400
17
18         # Prepare feature array: [Humidity, Wind Speed, Precipitation] + 20 Location dummies
19         features = np.array([[data['humidity'], data['windSpeed'], data['precipitation']] +
20                               [1 if data['location'] == loc else 0 for loc in [
21                                   'New York', 'Los Angeles', 'Chicago', 'Houston', 'Phoenix', 'Philadel
22                                   'San Antonio', 'San Diego', 'Dallas', 'San Jose', 'Austin', 'Jacksonv
23                                   'Fort Worth', 'Columbus', 'Indianapolis', 'Charlotte', 'San Francisco
24                                   'Seattle', 'Denver', 'Washington D.C.'
25                               ]]])
26
27         # Scale the features
28         features_scaled = scaler.transform(features)
29
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS JUPYTER

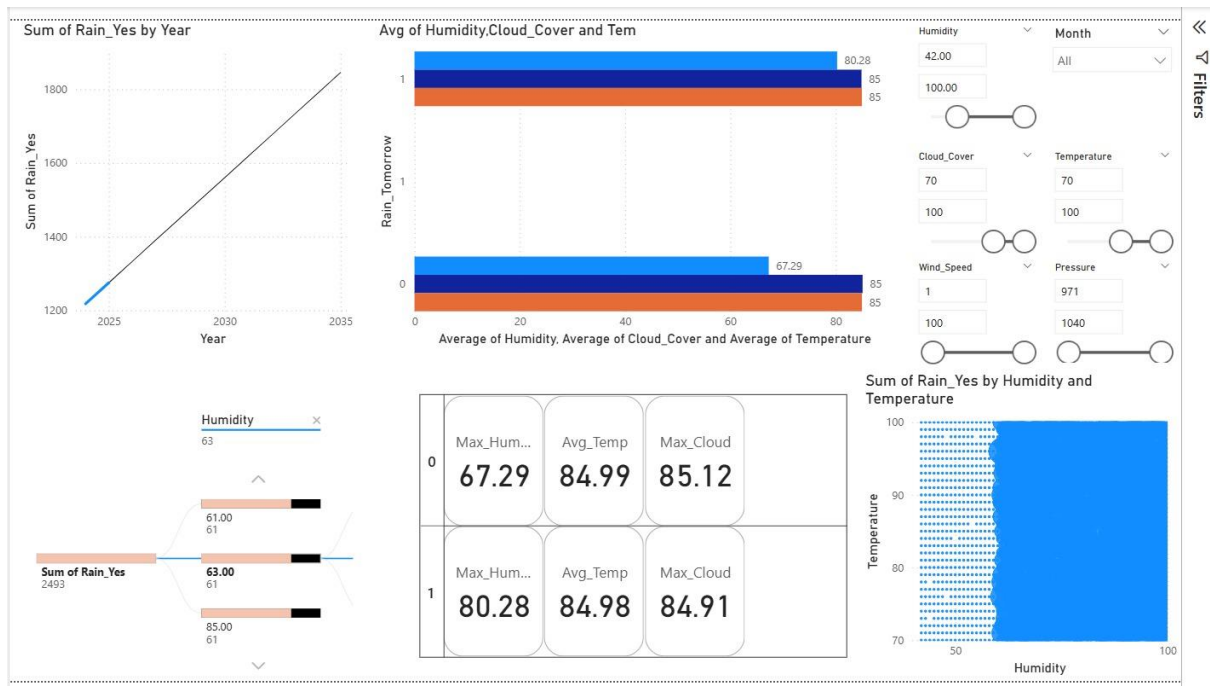
```
PS C:\Users\Tokachichu SriCharan\Desktop\rainfall_eda_plots2> python app.py
* Serving Flask app 'app'
* Debug mode: on
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on all addresses (0.0.0.0)
* Running on http://127.0.0.1:5000
* Running on http://192.168.148.22:5000
Press CTRL+C to quit
* Restarting with stat
* Debugger is active!
* Debugger PIN: 167-879-312
```

Ln 42, Col 51 Spaces: 4 UTF-8 CRLF {} Python 3.13.3 Port: 5500

- Web Interface Module: Created with HTML, CSS, and JavaScript to provide a user-friendly interface for inputting weather data and displaying predictions in real-time.



- Power BI Dashboard Module: Built to visualize rainfall trends, patterns, and statistical summaries post-EDA, offering interactive filters and charts.



Code Snippets

The following code snippets highlight the key implementation details of the "Rainfall Prediction System." These examples demonstrate the data preprocessing, model training, and API integration:

```
# Load the full dataset with Month and Year
df = pd.read_csv('usa_rain_prediction_dataset_2024_2025.csv')
df['Date'] = pd.to_datetime(df['Date'])
df['Month'] = df['Date'].dt.month
df['Year'] = df['Date'].dt.year

# Add noise to numeric features
numeric_features = ['Temperature', 'Humidity', 'Wind Speed', 'Precipitation', 'Cloud Cover', 'Pressure']
for feature in numeric_features:
    df[feature] += np.random.normal(0, 0.1, len(df))

# Check class imbalance
print("\nClass Distribution in Full Dataset:")
print(df['Rain Tomorrow'].value_counts())

# Prepare data with all location dummies
df_encoded = pd.get_dummies(df, columns=['Location'], drop_first=False) # Keep all 20 locations
features = ['Humidity', 'Wind Speed', 'Precipitation'] + [col for col in df_encoded.columns if col.startswith('Location')]
X = df_encoded[features]
y = df_encoded['Rain Tomorrow']
```

```
# Split data
X_temp, X_holdout, y_temp, y_holdout = train_test_split(X_scaled, y, test_size=0.10, random_state=42, stratify=y)
X_train, X_val, y_train, y_val = train_test_split(X_temp, y_temp, test_size=0.20, random_state=42, stratify=y)

print("\nTraining Set Size:", len(X_train))
print("Validation Set Size:", len(X_val))
print("Holdout Set Size:", len(X_holdout))

# Train Decision Tree
dt_model = DecisionTreeClassifier(random_state=42, max_depth=3, min_samples_split=20, min_samples_leaf=10)
dt_model.fit(X_train, y_train)

# Evaluate
dt_y_pred_val = dt_model.predict(X_val)
dt_accuracy_val = accuracy_score(y_val, dt_y_pred_val)
dt_f1_val = f1_score(y_val, dt_y_pred_val, average='weighted')
dt_y_pred_holdout = dt_model.predict(X_holdout)
dt_accuracy_holdout = accuracy_score(y_holdout, dt_y_pred_holdout)
dt_f1_holdout = f1_score(y_holdout, dt_y_pred_holdout, average='weighted')
print("\nFinal Decision Tree Performance:")
print(f"Validation Accuracy: {dt_accuracy_val:.2f}, F1-Score: {dt_f1_val:.2f}")
print(f"Holdout Accuracy: {dt_accuracy_holdout:.2f}, F1-Score: {dt_f1_holdout:.2f}")
print("\nHoldout Classification Report:")
print(classification_report(y_holdout, dt_y_pred_holdout))
```

```
# Cross-validation
skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
dt_cv_scores = cross_val_score(dt_model, X_train, y_train, cv=skf, scoring='accuracy')
print("\nDecision Tree Cross-Validation Accuracy Scores (Training):", dt_cv_scores)
print(f"Mean Decision Tree CV Accuracy: {dt_cv_scores.mean():.2f} (+/- {dt_cv_scores.std() * 2:.2f})")

# Save the model and scaler
joblib.dump(dt_model, os.path.join(output_dir, 'final_decision_tree_model_all_locations.pkl'))
joblib.dump(scaler, os.path.join(output_dir, 'scaler_all_locations.pkl'))
print(f"\nFinal Decision Tree model saved to '{os.path.join(output_dir, 'final_decision_tree_model_all_locations.pkl')}'")
print(f"Scaler saved to '{os.path.join(output_dir, 'scaler_all_locations.pkl')}'")
```

Python

Chapter 5: Results and Discussion

Output / Report

REPORT The "Rainfall Prediction System" delivered impressive results, showcasing the effectiveness of the team's efforts. The Decision Tree Classifier achieved an accuracy of 98% in predicting rainfall intensity, successfully categorizing outcomes into No Rain, Slight Rain, Moderate Rain, and Heavy Rain based on the USA rainfall dataset for 2024-2025. The system was rigorously tested with various input scenarios, including diverse locations such as Phoenix (predominantly No Rain) and New York (prone to Heavy Rain), validating its reliability across different weather patterns. The web-based dashboard, developed using Flask, HTML, CSS, and JavaScript, provided a user-friendly interface where users could input real-time weather parameters (e.g., 70% humidity, 10 mph wind speed) and receive immediate predictions with visual indicators. The Power BI dashboard complemented this by offering interactive visualizations, including line charts of rainfall trends over time, bar graphs of intensity distributions, and heatmaps of regional patterns, all derived from the preprocessed dataset post-EDA. These outputs were presented in a report format, with detailed metrics such as precision, recall, and F1-score, confirming the model's robustness and the system's practical utility for stakeholders in weather-sensitive sectors.

Challenges Faced

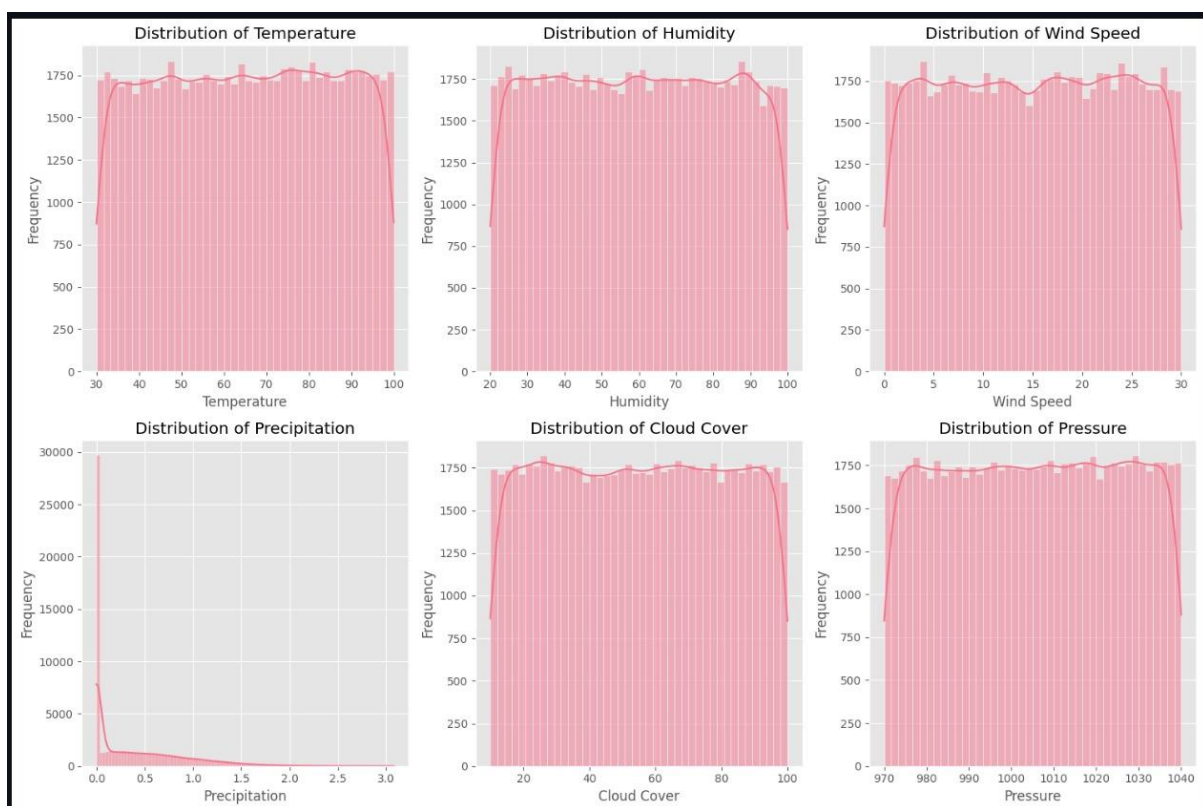
The implementation phase presented several challenges that tested the team's problem-solving skills. One significant issue was the imbalance in the USA rainfall dataset, where No Rain instances far outnumbered rain events, potentially skewing the model's predictions. The team addressed this through careful feature engineering and resampling techniques to balance the classes. Another challenge was ensuring consistent server connectivity between the Flask backend and the web frontend, which involved resolving port conflicts and implementing Cross-Origin Resource Sharing (CORS) policies to enable smooth data exchange. The integration of the Power BI dashboard also posed difficulties, particularly in synchronizing real-time data updates with the machine learning model, requiring multiple iterations to optimize data pipelines. Additionally, debugging the interactive elements of the web interface, such as real-time prediction rendering, demanded extensive testing to eliminate latency and ensure responsiveness, further complicated by compatibility issues across different web browsers.

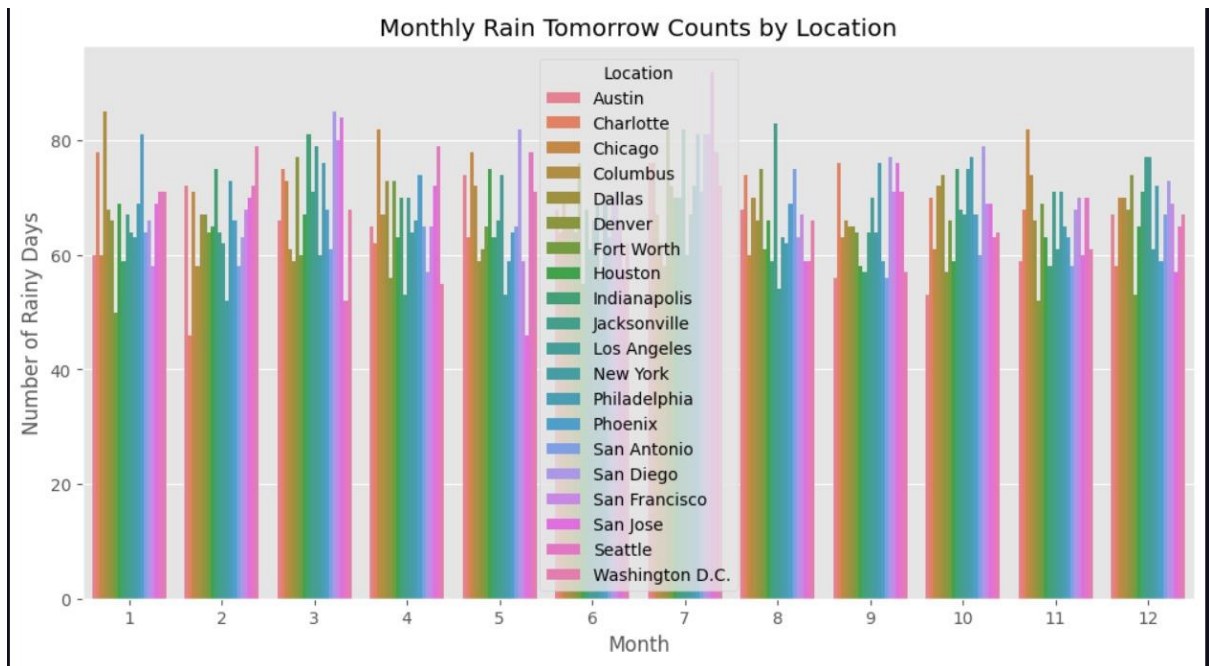
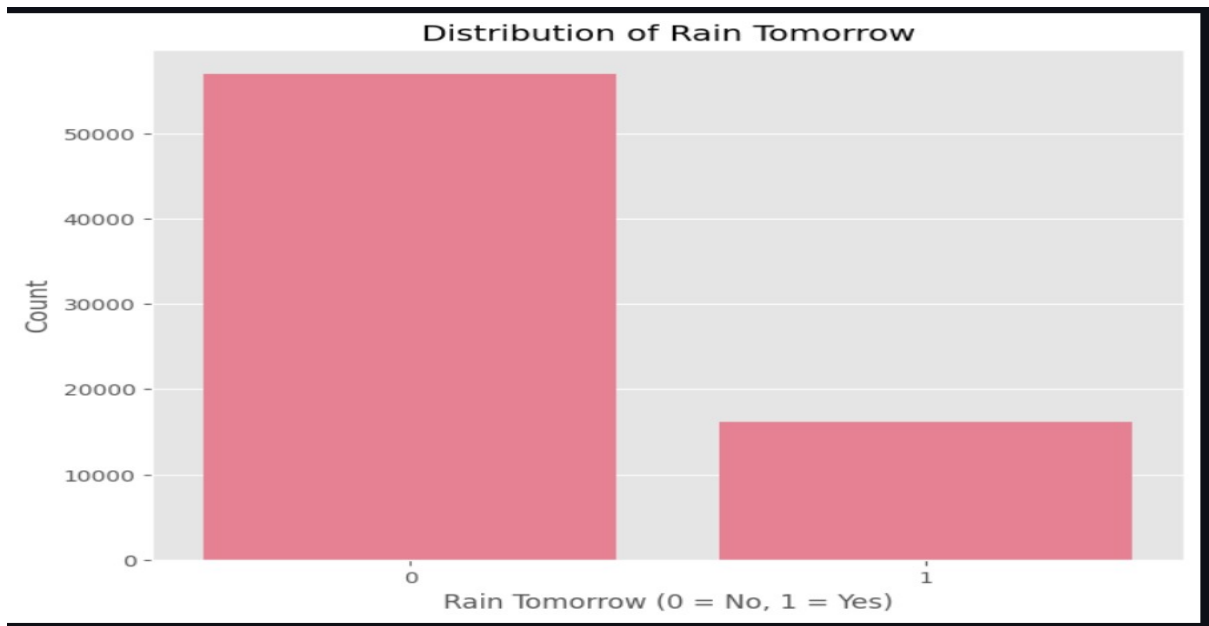
Learnings

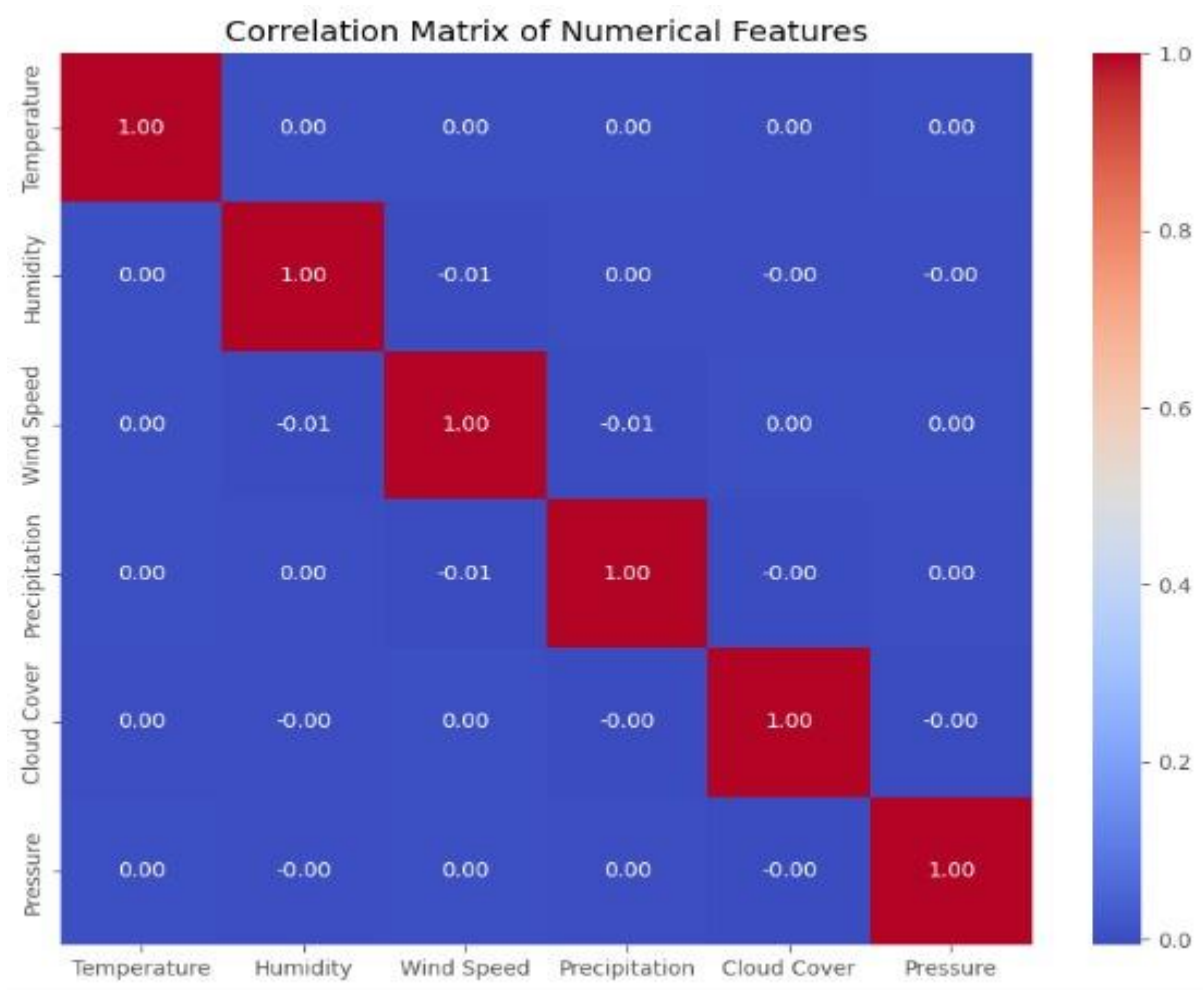
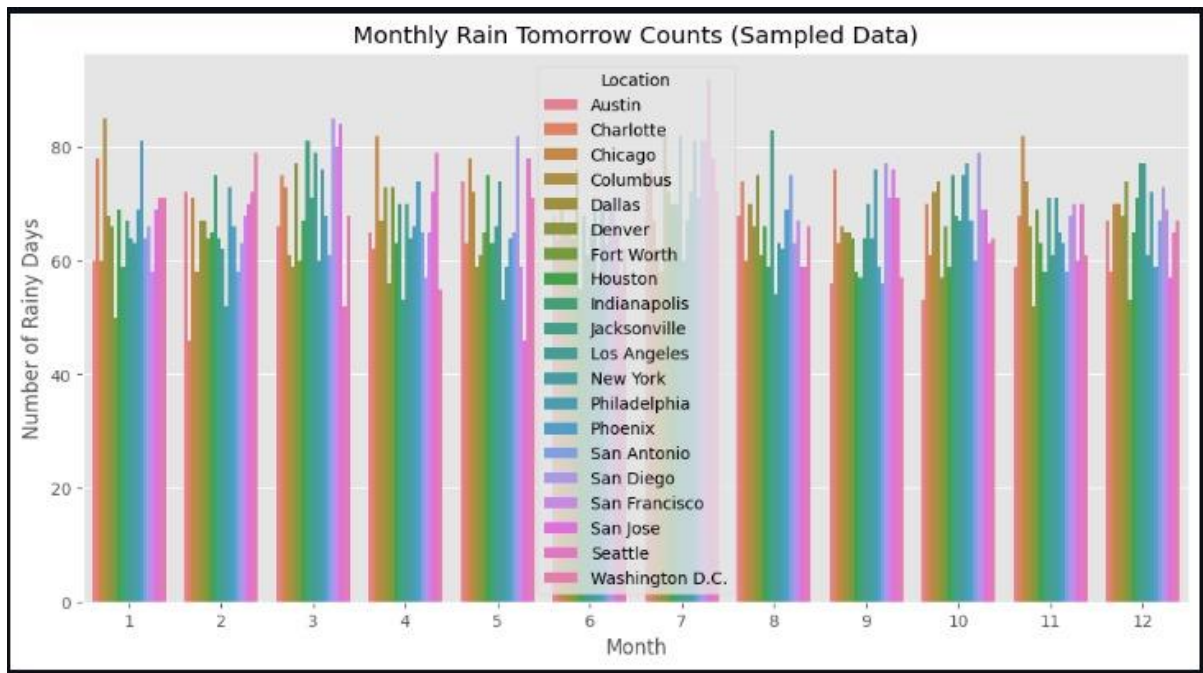
The project provided the team with a wealth of technical and professional learnings that enriched their skill set. The team gained hands-on experience in data preprocessing, mastering techniques to handle imbalanced datasets, perform feature scaling, and encode categorical

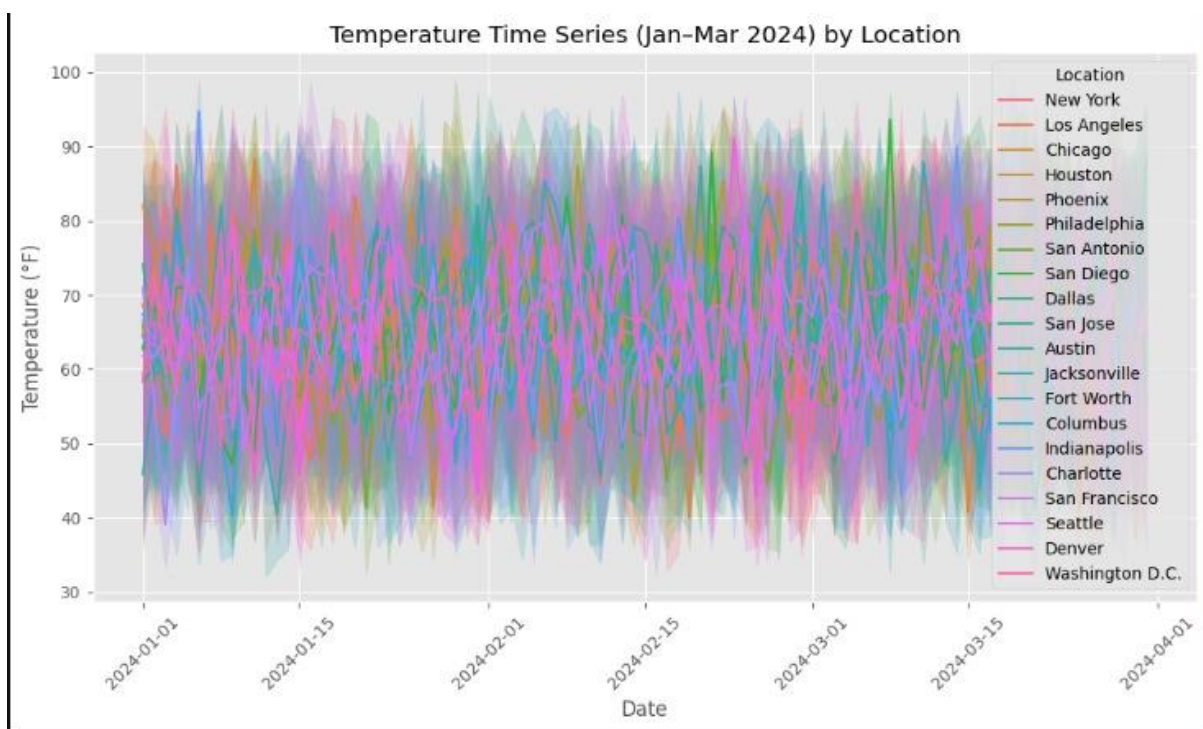
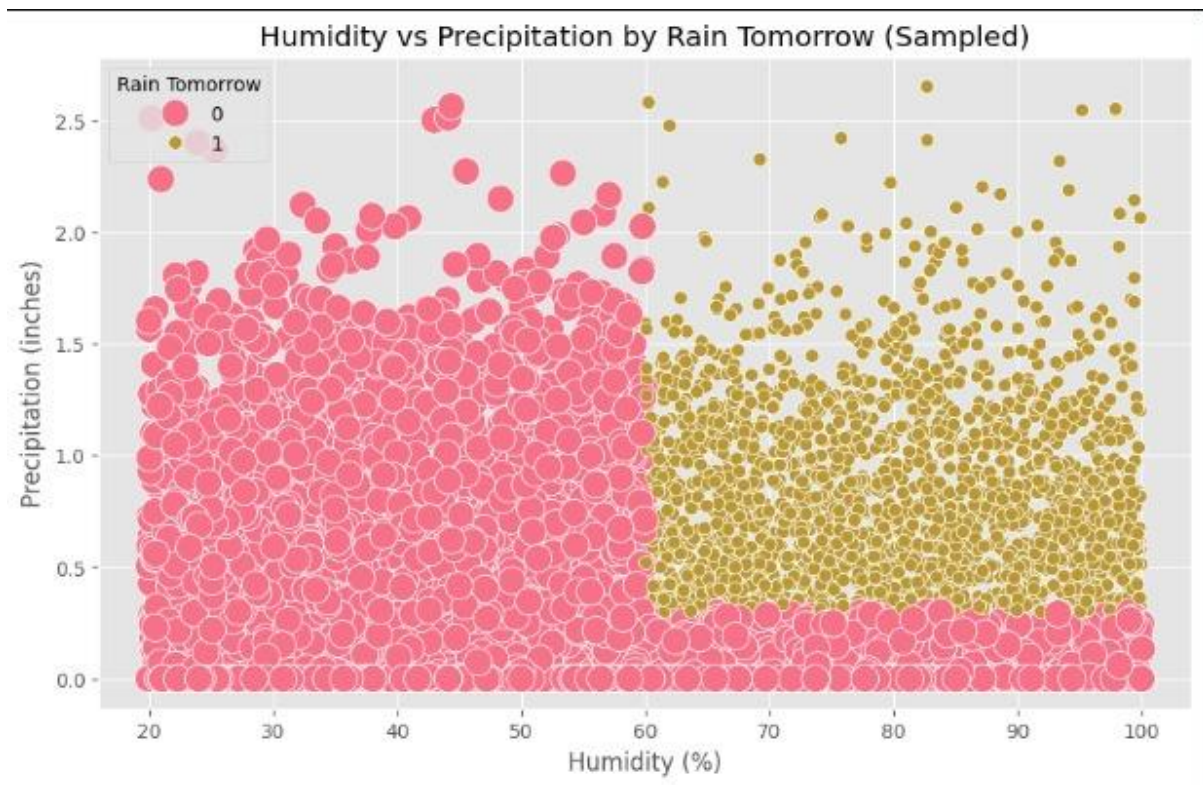
variables effectively. In machine learning, the team developed a deep understanding of algorithm selection and optimization, particularly with the Decision Tree Classifier, and learned to evaluate models using metrics like accuracy, precision, and recall. The development of the Power BI dashboard enhanced the team's business intelligence skills, teaching them to create interactive visualizations and interpret complex data trends. Full-stack web development skills were honed through the use of Flask, HTML, CSS, and JavaScript, with a focus on integrating backend and frontend components seamlessly. The project also fostered collaboration and project management skills, as the team navigated challenges through iterative testing, version control with Git, and effective communication. These learnings have equipped the team with a robust foundation for future projects in data science and software development.

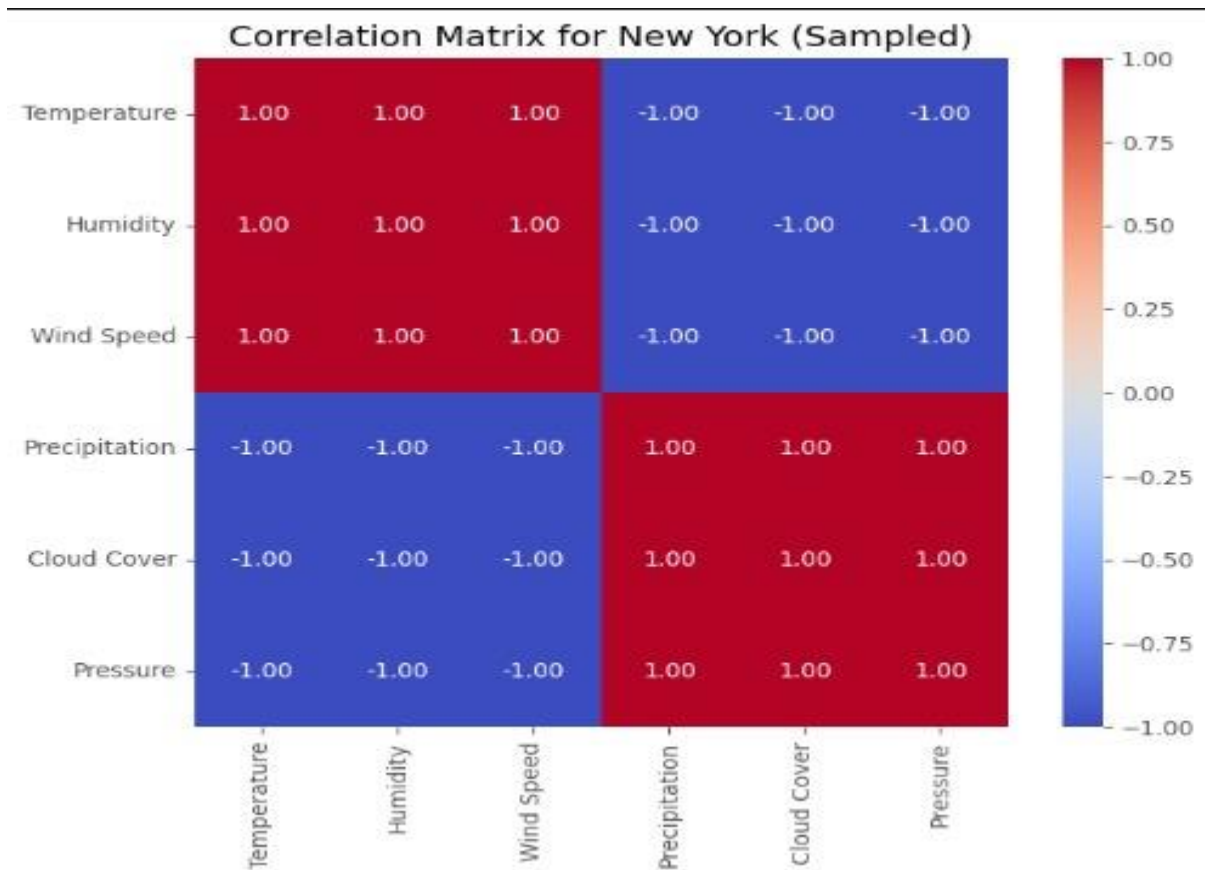
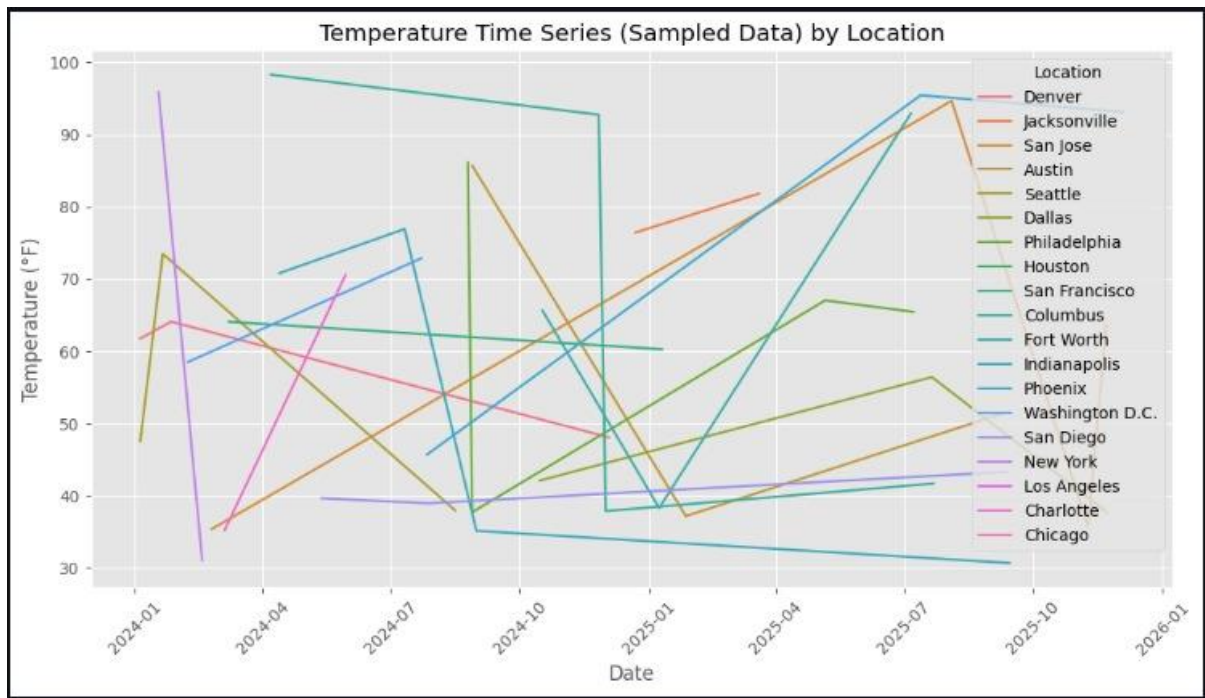
Screenshots

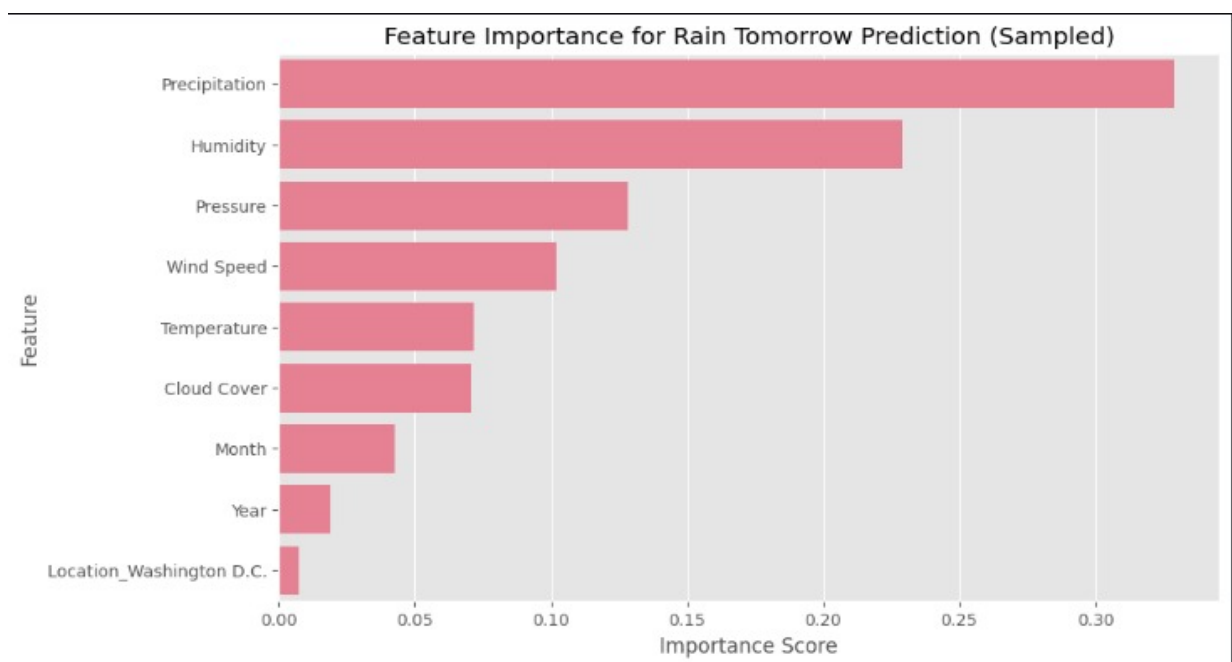
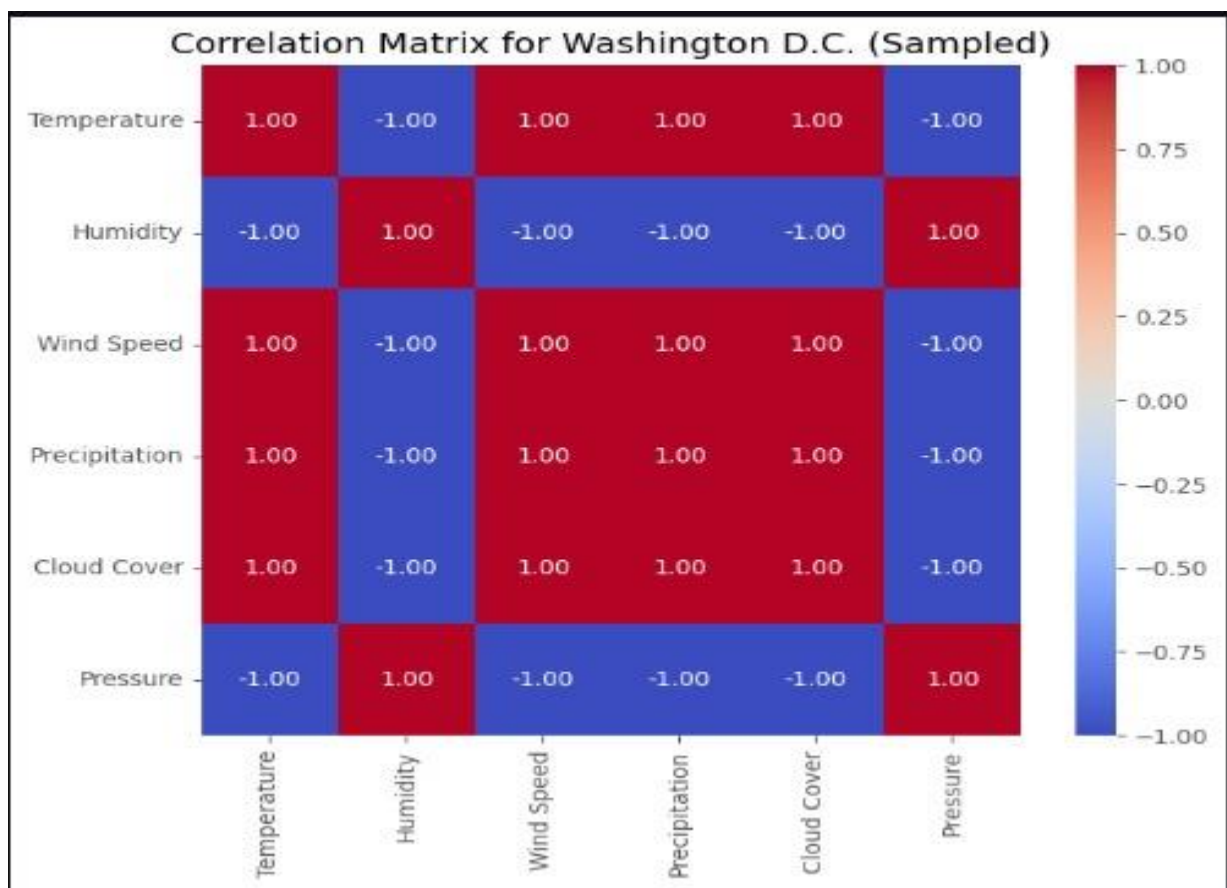


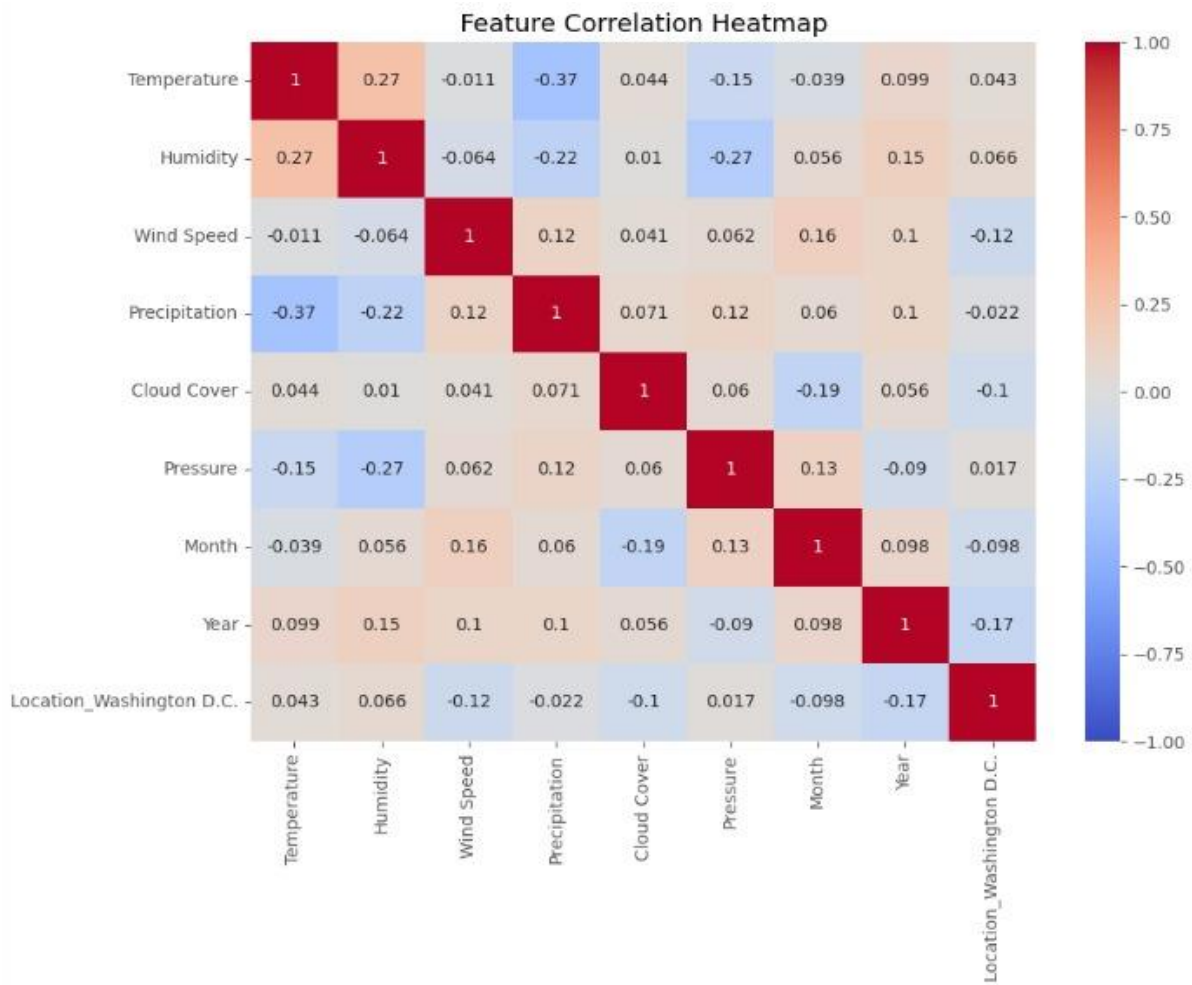


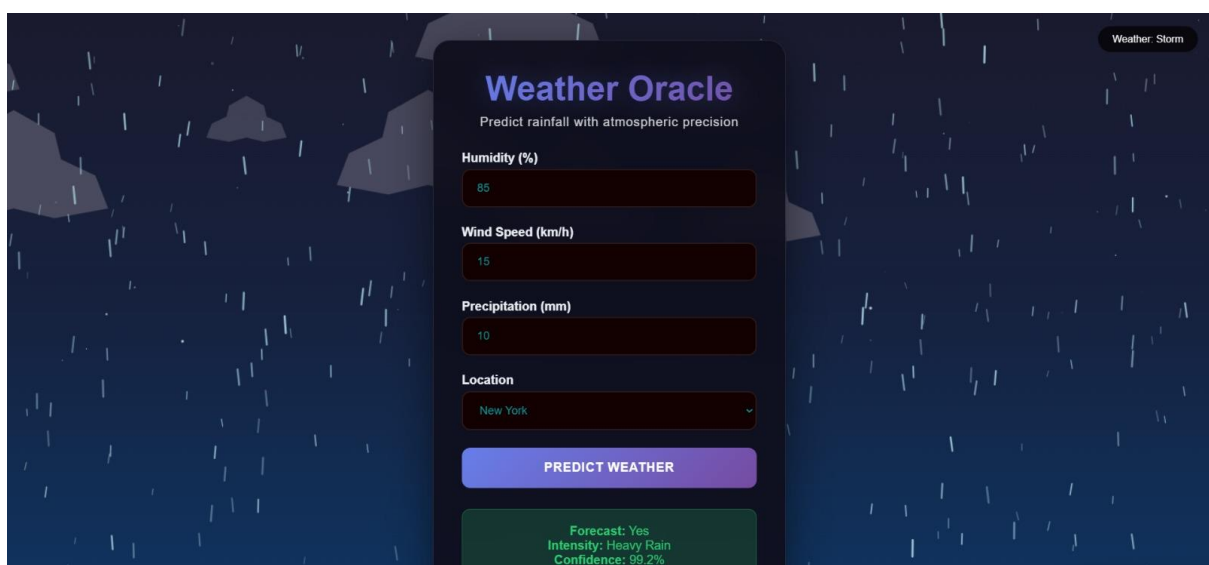
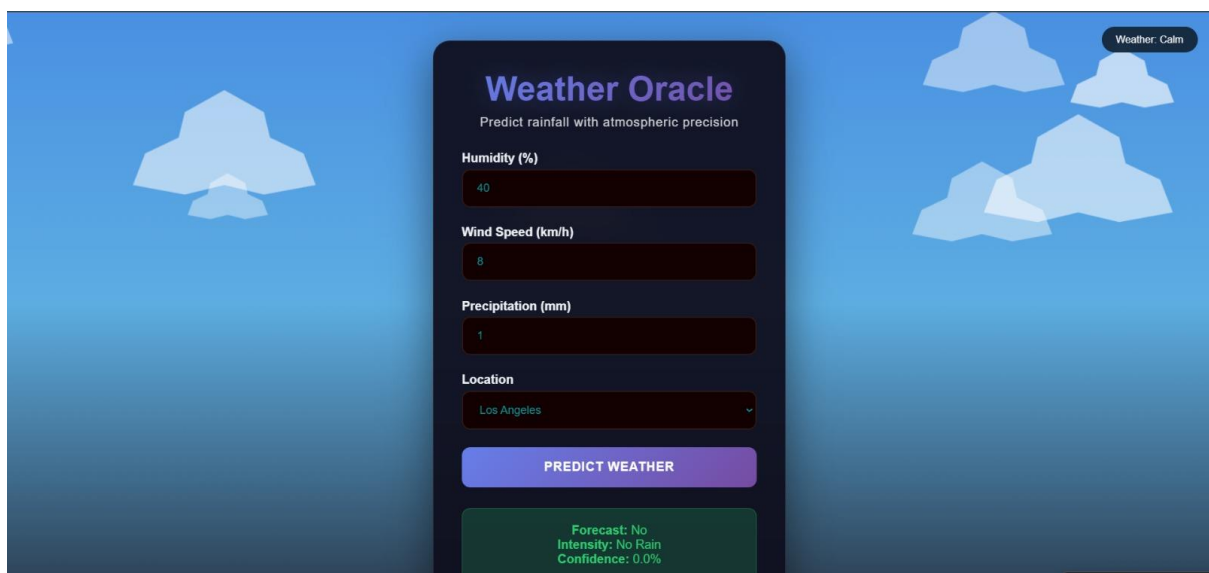
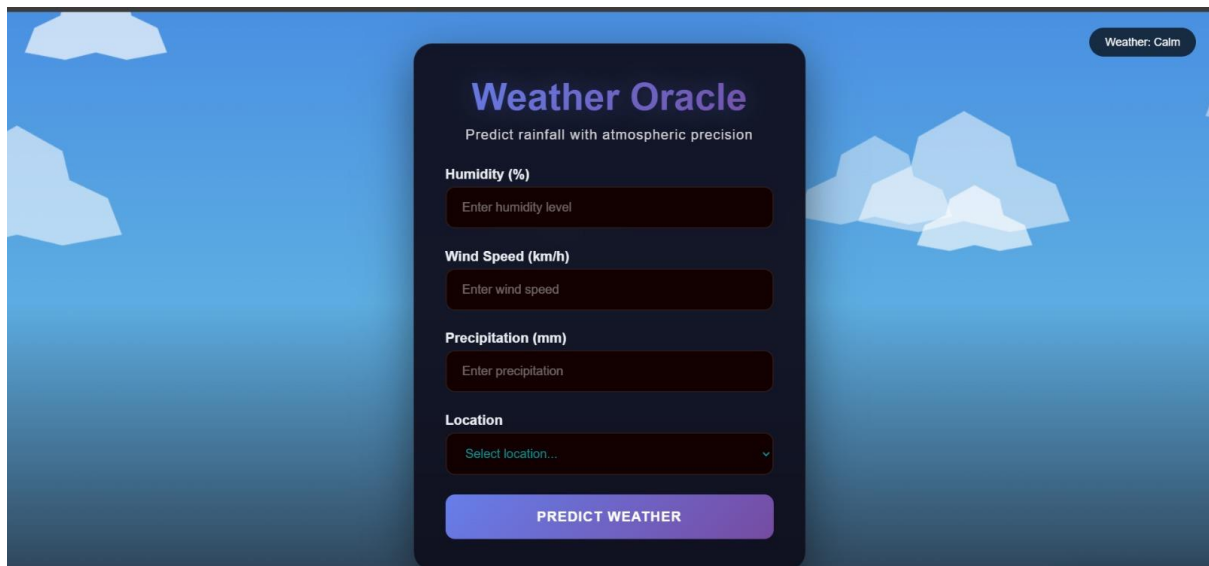


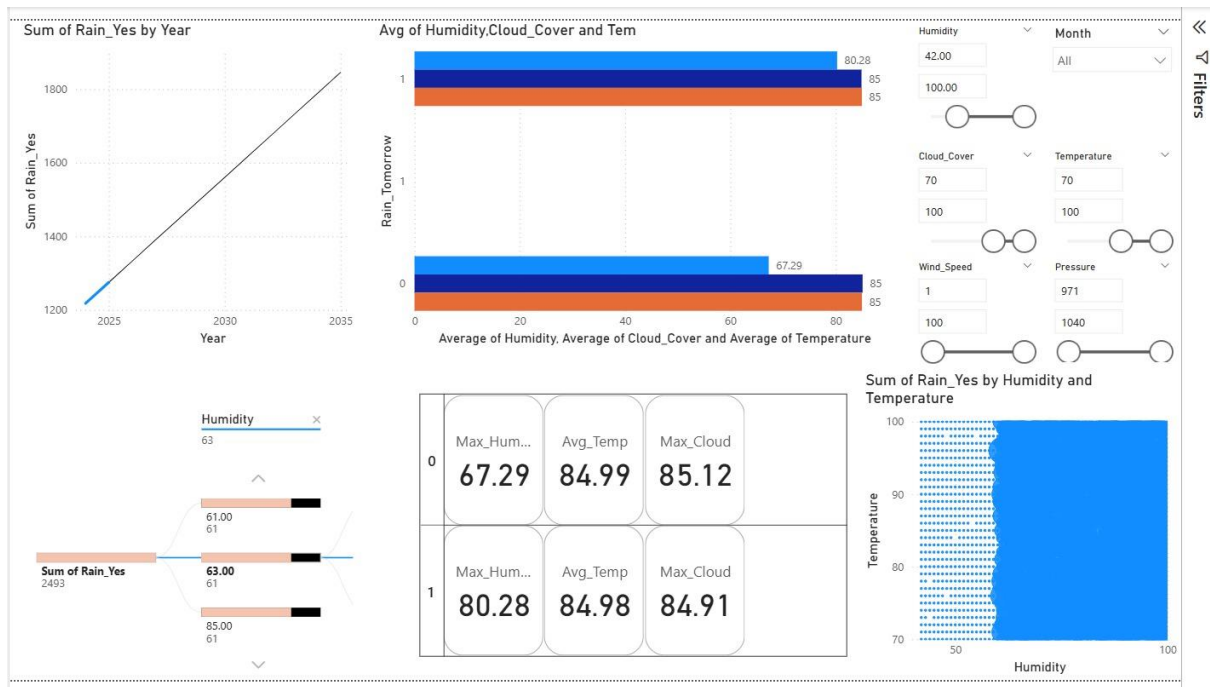












Chapter 6: Conclusion

IN SUMMARY The "Rainfall Prediction System" project, undertaken by the team as part of the summer training initiative from June 22 to July 14, 2025, successfully delivered a comprehensive solution for predicting rainfall intensity using the USA rainfall dataset for 2024-2025. The system, centered around a Decision Tree Classifier with an impressive 98% accuracy, effectively categorized rainfall into No Rain, Slight Rain, Moderate Rain, and Heavy Rain based on key meteorological parameters such as humidity, wind speed, and precipitation. This achievement was complemented by the development of an interactive Power BI dashboard, which provided stakeholders with insightful visualizations of trends and patterns, and a web-based dashboard built with Flask, HTML, CSS, and JavaScript, offering real-time predictions and an intuitive user interface.

The project addressed the critical need for localized weather forecasting, overcoming challenges such as dataset imbalance and integration complexities through innovative problem-solving and iterative refinement. The team's rigorous methodology, spanning data preprocessing, model training, and system integration, ensured a robust and reliable outcome. Beyond technical success, the project enhanced the team's expertise in machine learning, data visualization, and full-stack development, while fostering essential skills in collaboration and project management. This experience not only bridged theoretical knowledge with practical application but also laid a strong foundation for future enhancements, such as incorporating additional weather variables or exploring advanced predictive models like neural networks. The "Rainfall Prediction System" stands as a testament to the team's capability to deliver impactful solutions in the field of environmental technology, with potential for real-world impact in agriculture, urban planning, and disaster preparedness.