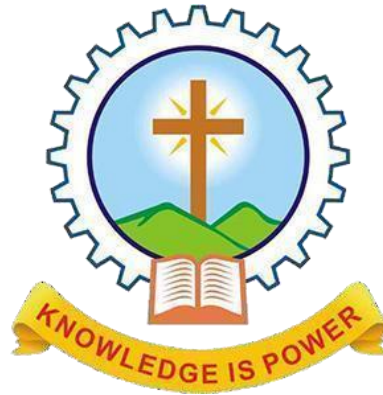


MAR ATHANASIUS COLLEGE OF ENGINEERING
(Affiliated to APJ Abdul Kalam Technological University, TVM)
KOTHAMANGALAM



Department of Computer Applications

Mini Project Report

RAINFALL PREDICTION USING
MACHINE LEARNING

Done by

Dona Jince

Reg No:MAC23MCA-2024

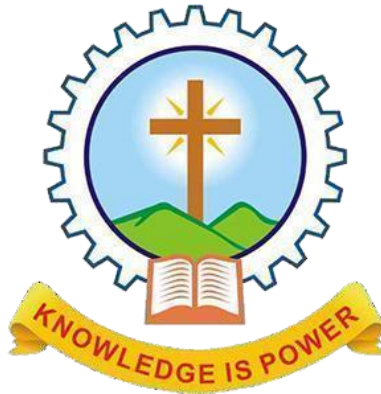
Under the guidance of
Prof. Sonia Abraham

2023-2025

MAR ATHANASIOUS COLLEGE OF ENGINEERING
(Affiliated to APJ Abdul Kalam Technological University, TVM)
KOTHAMANGALAM

MAR ATHANASIUS COLLEGE OF ENGINEERING
(Affiliated to APJ Abdul Kalam Technological University, TVM)
KOTHAMANGALAM

CERTIFICATE



Rainfall Prediction Using Machine Learning

Certified that this is the bonafide record of project work done by

Dona Jince
Reg No: MAC23MCA-2024

during the third semester, in partial fulfilment of requirements for award of the degree

Master of Computer Applications

of

APJ Abdul Kalam Technological University Thiruvananthapuram

Faculty Guide
Prof. Sonia Abraham

Head of the Department
Prof. Biju Skaria

Project Coordinator
Prof. Sonia Abraham

Internal Examiners

ACKNOWLEDGEMENT

With heartfelt gratitude, I extend my deepest thanks to the almighty for his unwavering grace and blessings that have made this journey possible. May his guidance continue to illuminate my path in the years ahead.

I am immensely thankful to Prof. Biju Skaria, Head of the Department of Computer Applications and Prof. Sonia Abraham, our dedicated project coordinator and my mini project guide, for their invaluable guidance and timely advice, which played a pivotal role in shaping this project. Their guidance, constant supervision, and provision of essential information were instrumental in the successful completion of the mini project.

I extend my profound thanks to all the professors in the department and the entire staff at MACE for their unwavering support and inspiration throughout my academic journey. My sincere appreciation goes to my beloved parents, whose guidance has been a beacon in every step of my path.

I am also grateful to my friends and individuals who generously shared their expertise and assistance, contributing significantly to the fulfillment of this endeavor.

ABSTRACT

The "Rainfall Prediction Using Machine Learning" project focuses on improving weather forecasting accuracy, particularly for predicting rainfall within a 24-hour period. This project employs three machine learning algorithms—Artificial Neural Networks (ANN), Support Vector Machines (SVM), and Random Forest (RF)—to accomplish this goal. Machine learning techniques are used due to their ability to identify complex patterns in large datasets, making them more effective than traditional methods in meteorological forecasting.

The prediction models are trained using the WeatherAUS dataset, which contains historical weather data such as temperature, humidity, wind speed, and rainfall records from multiple locations across Australia. Data preprocessing, including missing value imputation and handling class imbalance using SMOTE, ensures the dataset's quality. Feature selection is applied to identify the most relevant variables, while data standardization is performed to enhance model performance by ensuring consistency across features.

Each machine learning algorithm offers distinct advantages: ANN captures non-linear relationships between weather variables, SVM works well with high-dimensional data, and RF is robust against overfitting by aggregating predictions from multiple decision trees. The models are evaluated on various performance metrics, including accuracy, precision, recall, and F1-score, with Random Forest emerging as the most accurate, achieving 90.1% testing accuracy.

The primary goal of this project is to provide a reliable and accessible tool for predicting rainfall. A user-friendly web interface was developed to allow real-time predictions based on weather parameters, making it useful for farmers, water resource managers, and the general public.

In summary, this project addresses the challenges of traditional weather forecasting by applying machine learning models to improve short-term rainfall predictions. The integration of these models into a real-time prediction system has the potential to support better decision-making in agriculture, water management, and disaster preparedness.

LIST OF TABLES

2.1 Reference papers summary	5
------------------------------------	---

LIST OF FIGURES

3.1	Snapshot of WeatherAUS.csv	7
3.2	Types of Attributes	8
3.3	Count of each class labels.....	8
3.4	Description of features.....	9
3.5	Missing values in the dataset	9
3.6	Dropped rows with RainTomorrow null	10
3.7	Converting RainToday and RainTomorrow in to Numerical.....	10
3.8	Filling missing values using Mode	10
3.9	Converting categorical to numerical features.....	11
3.10	Multiple Imputation using Chained Equations	11
3.11	SMOTE.....	11
3.12	Feature Importance using f_classif	12
3.13	10 Selected features.....	12
3.14	Normalization using minmaxscalaer	12
3.15	RainTomorrow distribution before smote	15
3.16	RainTomorrow distribution after smote.....	15
3.17	Missing value pattern	16
3.18	Outlier detection.....	16
3.19	Correlation between attributes	17
3.20	Pairplot	18
3.21	ANN model	19

3.22	SVM working.....	20
3.23	RF working.....	21
3.24	ANN	22
3.25	SVM	23
3.26	RF	24
3.27	Project pipeline.....	25
4.1	Code snippet for dataset splitting.....	33
4.2	Code snippet for ANN model training.....	33
4.3	Code snippet for SVM model training.....	33
4.4	Code snippet for RF model training.....	34
5.1	ROC Curve	36
5.2	Precision-Recall curve.....	36
5.3	Code snippet for Rainfall Prediction	37
5.4	Rainfall Prediction	37
6.1	Home Page	39
6.2	Index Page.....	39
6.3	User input	40
6.4	Alert while submitting the input	40
6.5	Prediction	41
7.1	Git history	42
7.2	Git history of sprint 3.....	43

CONTENTS

ACKNOWLEDGEMENT	i
ABSTRACT	ii
LIST OF TABLES	iii
LIST OF FIGURES	iv
1 Introduction	1
2 Supporting Literature	2
2.1 Literature Review	2
2.1.1 Summary Table	5
2.2 Findings and Proposals	6
3 System Analysis	7
3.1 Analysis of Dataset	7
3.1.1 About the Dataset	7
3.1.2 Explore the dataset	8
3.2 Data Pre-processing	10
3.2.1 Data Cleaning	10
3.2.2 Analysis of Feature Variables	13
3.2.3 Analysis of Class Variables.....	14
3.3 Data Visualization	15
3.4 Analysis of Algorithms.....	19
3.4.1 Activity Diagrams	22
3.5 Project Plan.....	25
3.5.1 Project Pipeline	25

3.6	Feasibility Analysis	26
3.6.1	Technical Feasibility	26
3.6.2	Economic Feasibility	27
3.6.3	Operational Feasibility	28
3.7	System Environment.....	29
3.7.1	Software Environment	29
3.7.2	Hardware Environment	31
4	System Design	32
4.1	Model Building.....	32
4.1.1	Model Planning	32
4.1.2	Training.....	34
4.1.3	Testing.....	34
5	Results and Discussion	35
6	Model Deployment	38
7	Git History	42
8	Conclusion	44
9	Future Work	45
10	Appendix	46
10.1	Minimum Software Requirements.....	46
10.2	Minimum Hardware Requirements	47
11	References	48

1. INTRODUCTION

Rainfall prediction is an essential aspect of meteorology with significant implications for agriculture, water resource management, and disaster preparedness. Accurate and timely forecasts of rainfall can help farmers plan their irrigation schedules, assist authorities in managing water resources efficiently, and provide early warnings for floods and other natural disasters. Traditional statistical methods for rainfall prediction often fall short in capturing the complex, nonlinear relationships among various meteorological variables. As a result, there is a growing interest in leveraging advanced machine learning techniques to improve the accuracy and reliability of rainfall predictions.

This project focuses on developing a machine learning model to predict whether it will rain in the next 24 hours based on current weather data. The primary objective is to build and compare predictive models using three sophisticated machine learning algorithms: Artificial Neural Networks (ANN), Support Vector Machines (SVM), and Random Forest (RF). These algorithms have been chosen due to their proven efficacy in handling complex data patterns and their high accuracy in previous studies.

The dataset for this project is sourced from Kaggle's "weatheraus.csv" dataset, which includes extensive historical meteorological data from various regions in Australia. The dataset comprises features such as temperature, humidity, wind speed, atmospheric pressure, and historical rainfall records. To ensure high-quality input data, preprocessing steps will involve cleaning the data, handling missing values, normalizing features, and selecting relevant features.

A significant component of this project is the development of a user-friendly web interface that allows users to input current weather parameters and receive predictions on whether it will rain in the next 24 hours. This interface will make the model accessible to a broader audience, including farmers and the general public, enabling them to make informed decisions based on the predictions.

By integrating ANN, SVM, and RF, this project aims to enhance the accuracy and reliability of rainfall prediction models. The comprehensive evaluation and comparison of these algorithms will provide valuable insights into their performance, ultimately contributing to better decision-making in agriculture, water resource management, and disaster preparedness. The development of this advanced rainfall prediction model represents a significant step forward in leveraging machine learning for practical and impactful applications in meteorology.

2. SUPPORTING LITERATURE

2.1 Literature Review

Paper 1: M. M. Hassan *et al.*, "Machine Learning-Based Rainfall Prediction: Unveiling Insights and Forecasting for Improved Preparedness," in *IEEE Access*, vol. 11, pp. 132196-132222, 2023, doi: 10.1109/ACCESS.2023.3333876.

The paper "Machine Learning-Based Rainfall Prediction" by M. M. Hassan et al. presents an in-depth exploration of machine learning techniques to improve the accuracy and reliability of rainfall prediction models. Using a dataset that covers ten years of weather data from various Australian stations, the study includes crucial parameters like temperature, humidity, wind speed, and past rainfall records. The authors emphasize the importance of data preprocessing steps such as handling missing values, conducting outlier and correlation analysis, and performing feature selection. These steps aim to refine the quality of the input data, which is essential for maximizing the models' predictive capabilities.

The methodology of the study adopts a structured approach to developing and evaluating machine learning models. After the data preprocessing phase, the authors implement several machine learning algorithms, comparing their results to identify the most effective model for rainfall prediction. The algorithms tested include Naive Bayes, Decision Tree, Support Vector Machine (SVM), Random Forest, Logistic Regression, and Artificial Neural Network (ANN). The ANN model achieved the highest accuracy, starting at 90% before feature selection and improving to 91% after feature selection, demonstrating the positive impact of feature selection on model performance. This high level of accuracy is largely attributed to ANN's capability to capture and model complex, nonlinear patterns in the data, making it an excellent choice for rainfall prediction.

The study by M. M. Hassan et al. makes a notable contribution to the field of rainfall forecasting by systematically evaluating a range of machine learning algorithms and underscoring the effectiveness of ANNs. The findings suggest that machine learning techniques hold considerable potential for enhancing the precision of weather predictions, which is beneficial for sectors that depend on accurate rainfall forecasts, such as agriculture, water management, and disaster preparedness. Despite some limitations, the study's methodical approach and insightful recommendations offer a solid foundation for future research. The authors propose that integrating feature selection as a standard procedure in model development can further improve predictive accuracy, making it an essential area for ongoing research and optimization in machine learning-based weather forecasting.

Paper 2: Sarasa-Cabezuelo, A. Prediction of Rainfall in Australia Using Machine Learning. *Information* **2022**, *13*, 163. <https://doi.org/10.3390/info13040163>

The paper "Prediction of Rainfall in Australia Using Machine Learning" by Sarasa-Cabezuelo (2022) delves into the application of machine learning algorithms to enhance the precision of rainfall forecasting. Using a dataset containing meteorological data from various Australian cities, the study incorporates key variables like temperature, humidity, wind speed, and historical rainfall records. Data preprocessing is a critical part of the methodology, with steps such as addressing missing values, normalization, and feature selection to ensure a robust dataset that can optimize model performance.

Sarasa-Cabezuelo adopts a comparative approach by evaluating several machine learning algorithms, including k-Nearest Neighbors (k-NN), Decision Trees, Random Forest, and Neural Networks. Among these, Neural Networks achieved the highest accuracy (0.84) for short-term rainfall forecasts. This success is attributed to Neural Networks' ability to capture complex, nonlinear patterns within the data, providing more precise predictions than k-NN, Decision Trees, and Random Forest. This thorough comparison of models not only highlights Neural Networks as a promising approach for rainfall prediction but also provides valuable insights into the relative strengths and weaknesses of each algorithm.

The study's primary advantage lies in its comprehensive evaluation of multiple machine learning models, shedding light on their effectiveness for rainfall prediction. Neural Networks' superior accuracy underscores the potential of machine learning techniques for enhancing weather forecasting accuracy. Sarasa-Cabezuelo suggests that future research could incorporate additional meteorological variables, such as atmospheric pressure and solar radiation, to further boost model accuracy. Additionally, the study proposes the development of a real-time prediction system, possibly utilizing cloud computing to handle computational demands. By exploring advanced techniques like deep learning models and hybrid approaches that combine multiple algorithms, future work can aim to further enhance prediction accuracy and robustness, paving the way for more effective rainfall forecasting solutions.

Paper 3: Ghosh, S., Gourisaria, M.K., Sahoo, B. *et al.* A pragmatic ensemble learning approach for rainfall prediction. *Discov Internet Things* 3, 13 (2023). <https://doi.org/10.1007/s43926-023-00044-3>

The paper "A Pragmatic Ensemble Learning Approach for Rainfall Prediction" by Soumili Ghosh et al. delves into the application of ensemble learning methods to enhance the precision and robustness of rainfall forecasts. The research utilizes a comprehensive meteorological dataset, containing parameters like temperature, humidity and historical rainfall data collected from various cities across Australia. The data undergoes extensive preprocessing steps, including handling missing values, normalization, and feature selection, to ensure high-quality input for the machine learning models. This preprocessing is vital as it improves the data's consistency and relevance, which is essential for boosting model performance and ensuring accurate predictions.

The study employs an ensemble learning approach, combining multiple machine learning models to improve predictive accuracy. Ensemble methods like Bagging (Bootstrap Aggregating) and Boosting (using weak learners to form a strong predictor) are applied, as these approaches capitalize on the strengths of individual models, reduce overfitting, and enhance generalization. Specifically, the research evaluates Random Forest, Gradient Boosting, and AdaBoost as key ensemble algorithms. The results indicate that ensemble models generally outperform individual models, offering increased accuracy and robustness in predictions. Among the techniques tested, Gradient Boosting and Random Forest demonstrated notable improvements in rainfall prediction, though the Artificial Neural Network (ANN) surpassed these methods in terms of overall accuracy.

The study concludes that combining multiple models through ensemble learning leads to more reliable and accurate rainfall predictions compared to standalone models, despite the computational demands associated with these methods. The ensemble approach, therefore, presents a valuable avenue for future research and practical applications in weather forecasting. The authors propose that further investigation into optimizing ensemble methods could enhance model performance, particularly if integrated into real-time weather forecasting systems. This research highlights the potential of ensemble learning to address the challenges in weather prediction, contributing to improved decision-making in fields like agriculture, water management, and disaster preparedness.

2.1 Summary Table

Table 2.1 Reference papers summary

TITLE	DATASET	ALGORITHM	ACCURACY
M. M. Hassan <i>et al.</i> , "Machine Learning-Based Rainfall Prediction: Unveiling Insights and Forecasting for Improved Preparedness," in <i>IEEE Access</i> , vol. 11, pp. 132196-132222, 2023, doi: 10.1109/ACCESS.2023.3333 876.	Kaggle: weatheraus .csv : 145461 instances, 22 features	Artificial Neural Network Support Vector Machine Logistic Regression Lons Short Term Memory Random Forest Naïve Bayes	91% 84% 83% 83% 81% 81%
Sarasa-Cabezuelo, A. Prediction of Rainfall in Australia Using Machine Learning. <i>Information</i> 2022 , 13, 163. https://doi.org/10.3390/info13040163	Kaggle: weatheraus .csv : 145461 instances, 22 features	Neural Networks Random Forest K Nearest Neighbour Decision Trees	84% 83% 83% 83%
Ghosh, S., Gourisaria, M.K., Sahoo, B. <i>et al.</i> A pragmatic ensemble learning approach for rainfall prediction. <i>Discov</i> <i>Internet Things</i> 3 , 13 (2023). https://doi.org/10.1007/s43926-023-00044-3	Kaggle: weatheraus .csv : 145461 instances, 22 features	Random Forest Decision Trees Logistic Regression K Nearest Neighbour	87% 82% 79% 78%

Table 2.1 has the summary of all three papers.

2.2 Findings and Proposals

After a comprehensive review of three notable papers in the domain of rainfall prediction using machine learning, it is evident that each paper contributes substantially to the understanding and advancement of predictive models in this field. In the pursuit of an optimal rainfall prediction model, the choice of Artificial Neural Networks (ANN), Support Vector Machines (SVM), and Random Forest (RF) emerges as a sound decision. These models are well-suited for capturing complex, nonlinear relationships within meteorological data, making them robust solutions for handling variability and noise in weather datasets. ANN, in particular, has shown high accuracy in detecting intricate patterns, while RF and SVM offer robustness and generalization, especially in datasets with noise and variability.

The first paper by Hassan et al. presents a rainfall prediction model leveraging ANN with a focus on feature selection to enhance model accuracy, achieving a high precision rate of 91%. This approach highlights the effectiveness of ANN in capturing complex data patterns. The second paper by Sarasa-Cabezuelo compares k-Nearest Neighbors (k-NN), Decision Trees, RF, and ANN, identifying ANN as the most accurate model (Accuracy: 0.84) for short-term rainfall predictions. This comparative approach underscores ANN's potential for weather forecasting and its superior performance over traditional models like k-NN and Decision Trees. The third paper by Ghosh et al. applies an ensemble learning approach, combining multiple machine learning algorithms for improved predictive accuracy and resilience. Their findings reveal the efficacy of ensemble methods in enhancing prediction stability, achieving an accuracy of 85%.

Furthermore, the papers collectively advocate for the incorporation of additional meteorological parameters, such as atmospheric pressure and solar radiation, to improve prediction accuracy. They also suggest exploring cloud computing for real-time predictions and investigating hybrid models, combining multiple algorithms for enhanced robustness and accuracy. The integration of a user-friendly web interface is proposed to make these models accessible for practical applications, allowing users to input weather parameters and receive real-time rainfall forecasts.

In conclusion, the adoption of ANN, SVM, and RF for rainfall prediction is substantiated by empirical findings across these studies, emphasizing their efficacy in handling complex meteorological data. This approach aligns with the evolving needs of predictive weather systems, balancing accuracy, robustness, and accessibility, thereby fostering practical applications in agriculture, disaster management, and resource planning..

3. SYSTEM ANALYSIS

3.1 Analysis of Dataset

3.1.1 About the Dataset

The dataset used in this project is the "WeatherAUS" dataset, which provides comprehensive weather data from various locations in Australia. It includes 145,461 instances and 22 features, recorded over several years. This dataset is pivotal in developing robust rainfall prediction models using machine learning algorithms, as it includes a wide range of meteorological parameters recorded daily.

Features: The dataset includes various meteorological parameters such as temperature, humidity, wind speed, and atmospheric pressure. Each feature provides valuable information that can be leveraged to develop predictive models.

Target Variable: "RainTomorrow," is a boolean indicator of whether rainfall was recorded the following day. It is represented as 1 if rainfall occurred and 0 if not.

Dataset Link: <https://www.kaggle.com/datasets/trisha2094/weatheraus>

	Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	...	Humidity9am	Humidity3pm	Pressure9am	Pressure3pm
0	2008-12-01	Albury	13.4	22.9	0.6	NaN	NaN	W	44.0	W	...	71.0	22.0	1007.7	1007.1
1	2008-12-02	Albury	7.4	25.1	0.0	NaN	NaN	WNW	44.0	NNW	...	44.0	25.0	1010.6	1007.8
2	2008-12-03	Albury	12.9	25.7	0.0	NaN	NaN	WSW	46.0	W	...	38.0	30.0	1007.6	1008.7
3	2008-12-04	Albury	9.2	28.0	0.0	NaN	NaN	NE	24.0	SE	...	45.0	16.0	1017.6	1012.8
4	2008-12-05	Albury	17.5	32.3	1.0	NaN	NaN	W	41.0	ENE	...	82.0	33.0	1010.8	1006.0
Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow										
8.0	NaN	16.9	21.8	No	No										
NaN	NaN	17.2	24.3	No	No										
NaN	2.0	21.0	23.2	No	No										
NaN	NaN	18.1	26.5	No	No										
7.0	8.0	17.8	29.7	No	No										

Figure 3.1 Snapshot of WeatherAUS.csv

Figure 3.1 shows the result produce by listing the dataset files using head() function.

3.1.2 Explore the Dataset

The WeatherAUS dataset has both categorical and numerical features. There are 16 numerical attributes and 7 categorical attributes. Figure 3.2 shows the type of values of each attributes.

```

RangeIndex: 145460 entries, 0 to 145459
Data columns (total 23 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Date                  145460 non-null object
1   Location              145460 non-null object
2   MinTemp               143975 non-null float64
3   MaxTemp              144199 non-null float64
4   Rainfall              142199 non-null float64
5   Evaporation           82670 non-null float64
6   Sunshine              75625 non-null float64
7   WindGustDir           135134 non-null object
8   WindGustSpeed         135197 non-null float64
9   WindDir9am            134894 non-null object
10  WindDir3pm            141232 non-null object
11  WindSpeed9am          143693 non-null float64
12  WindSpeed3pm          142398 non-null float64
13  Humidity9am           142806 non-null float64
14  Humidity3pm           140953 non-null float64
15  Pressure9am           130395 non-null float64
16  Pressure3pm           130432 non-null float64
17  Cloud9am              89572 non-null float64
18  Cloud3pm              86102 non-null float64
19  Temp9am               143693 non-null float64
20  Temp3pm               141851 non-null float64
21  RainToday             142199 non-null object
22  RainTomorrow          142193 non-null object
dtypes: float64(16), object(7)
memory usage: 25.5+ MB

```

Figure 3.2 Types of attributes in the dataset

	count
RainTomorrow	
No	110316
Yes	31877

Figure 3.3 Count of each class in target variable

Figure 3.3 shows the distribution of values of target variable, RainTomorrow.

	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustSpeed	WindSpeed9am	WindSpeed3pm	Humidity9am	Humidity3pm	Pressure9am
count	143975.000000	144199.000000	142199.000000	82670.000000	75625.000000	135197.000000	143693.000000	142398.000000	142806.000000	140953.000000	130395.000000
mean	12.194034	23.221348	2.360918	5.468232	7.611178	40.035230	14.043426	18.662657	68.880831	51.539116	1017.64994
std	6.398495	7.119049	8.478060	4.193704	3.785483	13.607062	8.915375	8.809800	19.029164	20.795902	7.10653
min	-8.500000	-4.800000	0.000000	0.000000	0.000000	6.000000	0.000000	0.000000	0.000000	0.000000	980.50000
25%	7.600000	17.900000	0.000000	2.600000	4.800000	31.000000	7.000000	13.000000	57.000000	37.000000	1012.90000
50%	12.000000	22.600000	0.000000	4.800000	8.400000	39.000000	13.000000	19.000000	70.000000	52.000000	1017.60000
75%	16.900000	28.200000	0.800000	7.400000	10.600000	48.000000	19.000000	24.000000	83.000000	66.000000	1022.40000
max	33.900000	48.100000	371.000000	145.000000	14.500000	135.000000	130.000000	87.000000	100.000000	100.000000	1041.00000

Figure 3.4 Description of features

Figure 3.4 shows the description of some features in the dataset.

	missing_values	percent_missing %	data type
Date	0	0.000000	object
Location	0	0.000000	object
MinTemp	1485	1.020899	float64
MaxTemp	1261	0.866905	float64
Rainfall	3261	2.241853	float64
Evaporation	62790	43.166506	float64
Sunshine	69835	48.009762	float64
WindGustDir	10326	7.098859	object
WindGustSpeed	10263	7.055548	float64
WindDir9am	10566	7.263853	object
WindDir3pm	4228	2.906641	object
WindSpeed9am	1767	1.214767	float64
WindSpeed3pm	3062	2.105046	float64
Humidity9am	2654	1.824557	float64
Humidity3pm	4507	3.098446	float64
Pressure9am	15065	10.356799	float64
Pressure3pm	15028	10.331363	float64
Cloud9am	55888	38.421559	float64
Cloud3pm	59358	40.807095	float64
Temp9am	1767	1.214767	float64
Temp3pm	3609	2.481094	float64
RainToday	3261	2.241853	object
RainTomorrow	3267	2.245978	object

Figure 3.5 Missing values in the dataset

Figure 3.5 shows the count, percentage, type of missing values in each attributes in the dataset.

3.2 Data Preprocessing

3.2.1 Data Cleaning

Data Cleaning is the data pre-processing method we choose. Data cleaning routines attempt to fill in missing values, smooth out noisy data and correct inconsistencies. This dataset has many missing values which has to be addressed. Categorical attributes and numerical values are treated separately. At first, the rows with target variable null is dropped , then the missing values in the categorical features are filled using mode. Then using labelencoder categorical features are converted in to numerical. Then by using Multiple Imputation by Chained Equations the missing values in the entire dataset is treated.After this process there will be no missing values. Figure 3.6 shows the dropping of rows with target variable null.

Treating Missing Values

```
df = df.dropna(subset=['RainTomorrow'])
```

Figure 3.6 Dropped rows with RainTomorrow null

```
df['RainToday'].replace({'No': 0, 'Yes': 1},inplace = True)
df['RainTomorrow'].replace({'No': 0, 'Yes': 1},inplace = True)
```

Figure 3.7 Converting RainToday and RainTomorrow in to numerical

Figure 3.7 shows the conversion of values of RainToday and RainTomorrow in to 0 for 'No' and 1 for 'Yes'.

```
df['Date'] = df['Date'].fillna(df['Date'].mode()[0])
df['Location'] = df['Location'].fillna(df['Location'].mode()[0])
df['WindGustDir'] = df['WindGustDir'].fillna(df['WindGustDir'].mode()[0])
df['WindDir9am'] = df['WindDir9am'].fillna(df['WindDir9am'].mode()[0])
df['WindDir3pm'] = df['WindDir3pm'].fillna(df['WindDir3pm'].mode()[0])
```

Figure 3.8 Filling missing values using mode

Figure 3.8 shows the filling of missing values in the categorical attributes using mode.

```
# Convert categorical features to continuous features with Label Encoding
lencoders = {}
for col in df.select_dtypes(include=['object']).columns:
    lencoders[col] = LabelEncoder()
    df[col] = lencoders[col].fit_transform(df[col])
```

Figure 3.9 Converting categorical to numerical features

Figure 3.9 shows the the conversion of categorical in to numerical features.

```
warnings.filterwarnings("ignore")
MiceImputed = df.copy(deep=True)
mice_imputer = IterativeImputer()
MiceImputed.iloc[:, :] = mice_imputer.fit_transform(df)
```

Figure 3.10 Multiple Imputation by Chained Equations

Figure 3.10 shows the usage of Multiple Imputation by Chained Equations for treating the missing values in the entire dataset.

Handling class imbalance

Dataset contains value No more than Yes for target variable RainTomorrow, so it is imbalanced and has to be balanced. Used smote for balancing the dataset. Figure 3.11 shows applying smote to the dataset.

```
X = MiceImputed.drop('RainTomorrow', axis=1)
y = MiceImputed['RainTomorrow']
# Apply SMOTE
smote = SMOTE(random_state=42)
X_resampled, y_resampled = smote.fit_resample(X, y)
# Create a new DataFrame with the resampled data
resampled_df = pd.concat([pd.DataFrame(X_resampled), pd.DataFrame(y_resampled)], axis=1)
```

Figure 3.11 SMOTE

Feature Selection

Filter method is used for feature selection to train the rainfall prediction models. Features are selected using `f_classif`.

```
X = resampled_df.drop('RainTomorrow', axis=1)
y = resampled_df['RainTomorrow']
# Select the 10 best features from the feature set (excluding 'RainTomorrow')
selector = SelectKBest(f_classif, k=10)
X_selected = selector.fit_transform(X, y)
```

Figure 3.12 Feature Importance using `f_classif`

```
Selected columns: Index(['Sunshine', 'WindGustSpeed', 'Humidity9am', 'Humidity3pm',
                        'Pressure9am', 'Pressure3pm', 'Cloud9am', 'Cloud3pm', 'Temp3pm',
                        'RainToday'],
                        dtype='object')
```

Figure 3.13 10 selected features

Figure 3.12 shows the application of `f_classif` to select 10 best features and Figure 3.13 shows the 10 selected features.

Normalizing data

Normalization of data is done using `MinMaxScaler`

```
#Normalizing data
r_scaler = preprocessing.MinMaxScaler()
r_scaler.fit(modified_data_selected)
modified_data = pd.DataFrame(r_scaler.transform(modified_data_selected),
                             index=modified_data_selected.index,
                             columns=modified_data_selected.columns)
modified_data['RainTomorrow'] = y.values
```

Figure 3.14 Normalization using `minmaxscaler`

Figure 3.14 shows the normalization of features using `minmaxscaler`

3.2.2 Analysis of Features

Analysis of features describe their relevance to rainfall prediction and potential influence on the machine learning model:

1. Date: Represents the day the measurements were recorded. Although not directly useful for prediction, it may be helpful for identifying seasonal patterns in rainfall.
2. Location: Indicates the weather station's location, which is valuable for capturing regional climate variations that affect rainfall. Different locations may have unique weather trends.
3. MinTemp: The minimum temperature recorded during the day. Lower temperatures might correlate with rainfall as temperature affects atmospheric conditions.
4. MaxTemp: The maximum temperature recorded during the day. High temperatures can influence evaporation rates and humidity, impacting the likelihood of rainfall.
5. Rainfall: Records the amount of rainfall (in mm) for the day, providing context for past rainfall, which can be a predictor for future rainfall events.
6. Evaporation: Reflects the amount of water evaporated (in mm) over 24 hours. Higher evaporation rates may indicate drier conditions, which could be inversely related to rainfall probability.
7. Sunshine: Measures the hours of sunlight received during the day. Less sunshine often correlates with cloudy or rainy conditions, making this feature relevant for rainfall prediction.
8. WindGustDir: The direction of the strongest wind gust over 24 hours. Specific wind directions can signal approaching weather systems that may bring rain.
9. WindGustSpeed: The speed of the strongest wind gust (km/h). Strong wind gusts can signify stormy weather conditions, potentially leading to rainfall.
10. WindDir9am: The direction of the wind at 9 a.m., which may indicate early atmospheric conditions relevant for rainfall prediction.
11. WindDir3pm: The wind direction at 3 p.m., helping to capture weather changes throughout the day that might influence evening or overnight rain.
12. WindSpeed9am: The average wind speed before 9 a.m., which can affect evaporation rates and indicate changing weather conditions.
13. WindSpeed3pm: The average wind speed before 3 p.m., similar to WindSpeed9am, providing additional insight into atmospheric conditions as they develop.

14. Humidity9am: The humidity percentage at 9 a.m. High humidity often indicates moisture in the air, which could increase the likelihood of rainfall.
15. Humidity3pm: The humidity percentage at 3 p.m. This feature is critical for short-term rainfall prediction as high afternoon humidity might signal impending rain.
16. Pressure9am: The atmospheric pressure at 9 a.m. Low pressure often correlates with rainy weather, so this feature is a valuable predictor.
17. Pressure3pm: Atmospheric pressure at 3 p.m., allowing the model to track pressure changes throughout the day to detect potential rain events.
18. Cloud9am: Fraction of the sky covered by clouds at 9 a.m. (measured in oktas). Higher cloud cover suggests overcast conditions, which can indicate rain.
19. Cloud3pm: Fraction of the sky covered by clouds at 3 p.m. This provides an afternoon snapshot of cloudiness, useful for predicting rain later in the day.
20. Temp9am: The temperature at 9 a.m., helpful for understanding morning atmospheric conditions.
21. Temp3pm: The temperature at 3 p.m., capturing afternoon conditions that may affect the likelihood of rainfall later in the day.
22. RainToday: A boolean indicator of whether rain was recorded today (1 for yes, 0 for no). This is directly relevant, as recent rainfall may suggest weather patterns likely to bring more rain.

Each feature contributes uniquely to understanding atmospheric conditions and has the potential to influence rainfall prediction models. The diversity and completeness of these meteorological parameters make the dataset highly suitable for machine learning applications focused on predicting rainfall.

3.2.3 Analysis of Class Variables

The class variable in the "WeatherAUS" dataset is RainTomorrow, which is a binary variable indicating whether it rained the following day. This variable has two possible values:

- 1: Indicates that rainfall was recorded the next day (Yes).
- 0: Indicates that no rainfall was recorded the next day (No).

3.3 Data Visualization

Data visualization serves as a fundamental technique in data analysis, offering a graphical representation of information through charts, graphs, and maps. Its significance lies in transforming complex datasets into visually intuitive formats, allowing for the identification of trends, patterns, and outliers. This visual representation enhances data exploration, enabling easy comparison between data points and facilitating effective communication of insights. The diverse types of visualizations, such as charts for trends, tables for structured data, graphs for relationships, maps for geographic data, and dashboards for comprehensive displays, cater to various analytical needs. Utilizing tools like Tableau, Excel, or programming libraries like Matplotlib, data visualization aids in making data-driven decisions by providing accessible insights into massive amounts of information.

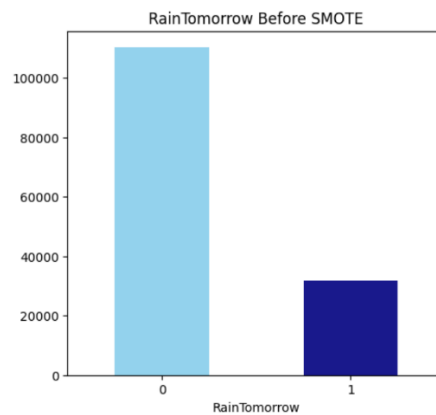


Figure 3.15 RainTomorrow distribution before smote

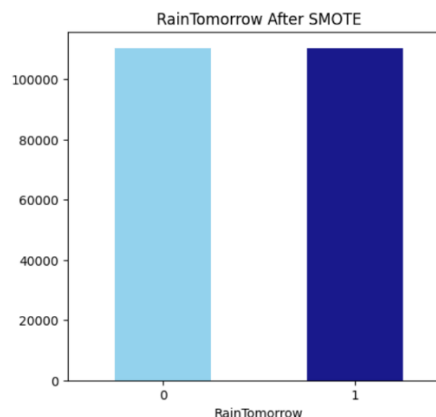


Figure 3.16 RainTomorrow distribution after smote

Figure 3.15 and 3.16 shows the distribution of target variable RainTomorrow before and after smote.

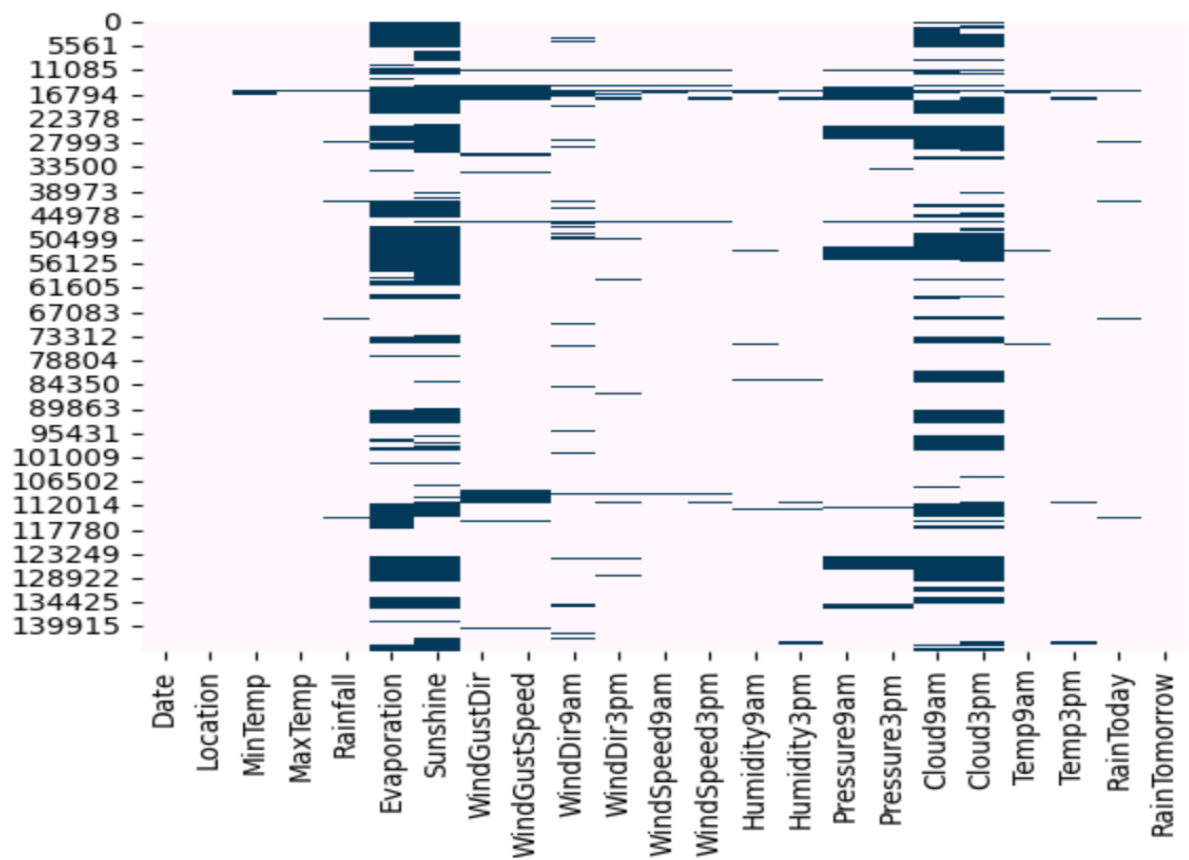


Figure 3.17 Missing value pattern

Figure 3.17 shows the missing value pattern in the dataset before data cleaning.

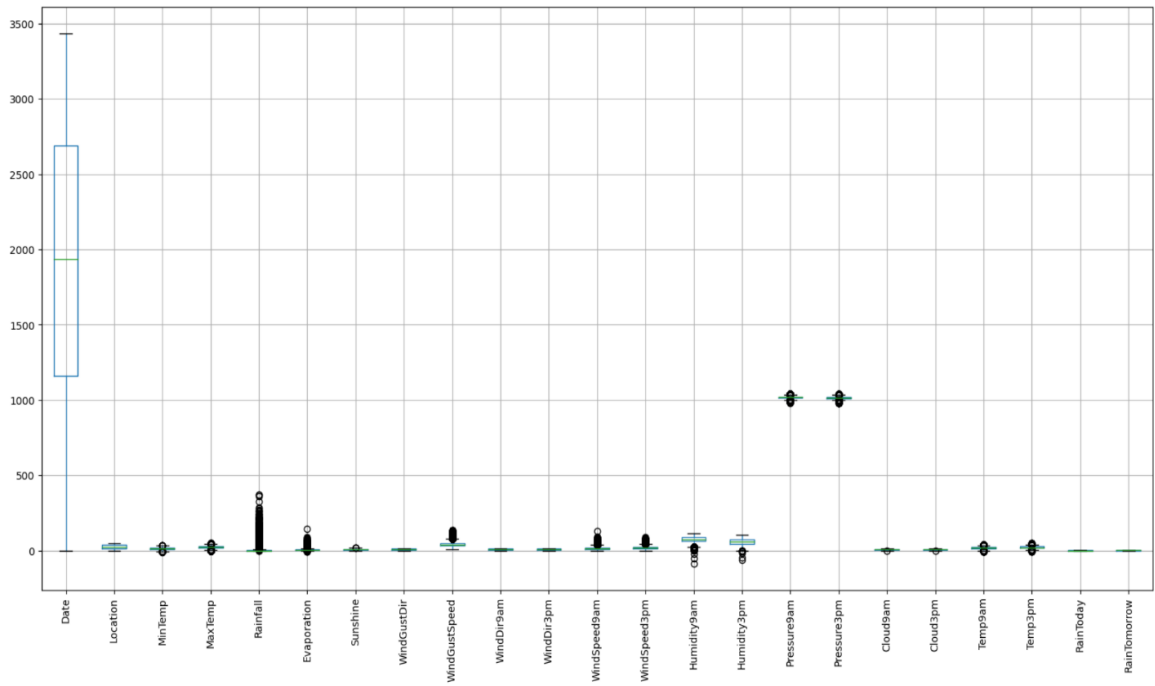


Figure 3.18 Outlier Detection

Figure 3.18 shows that there are no outliers in the dataset.

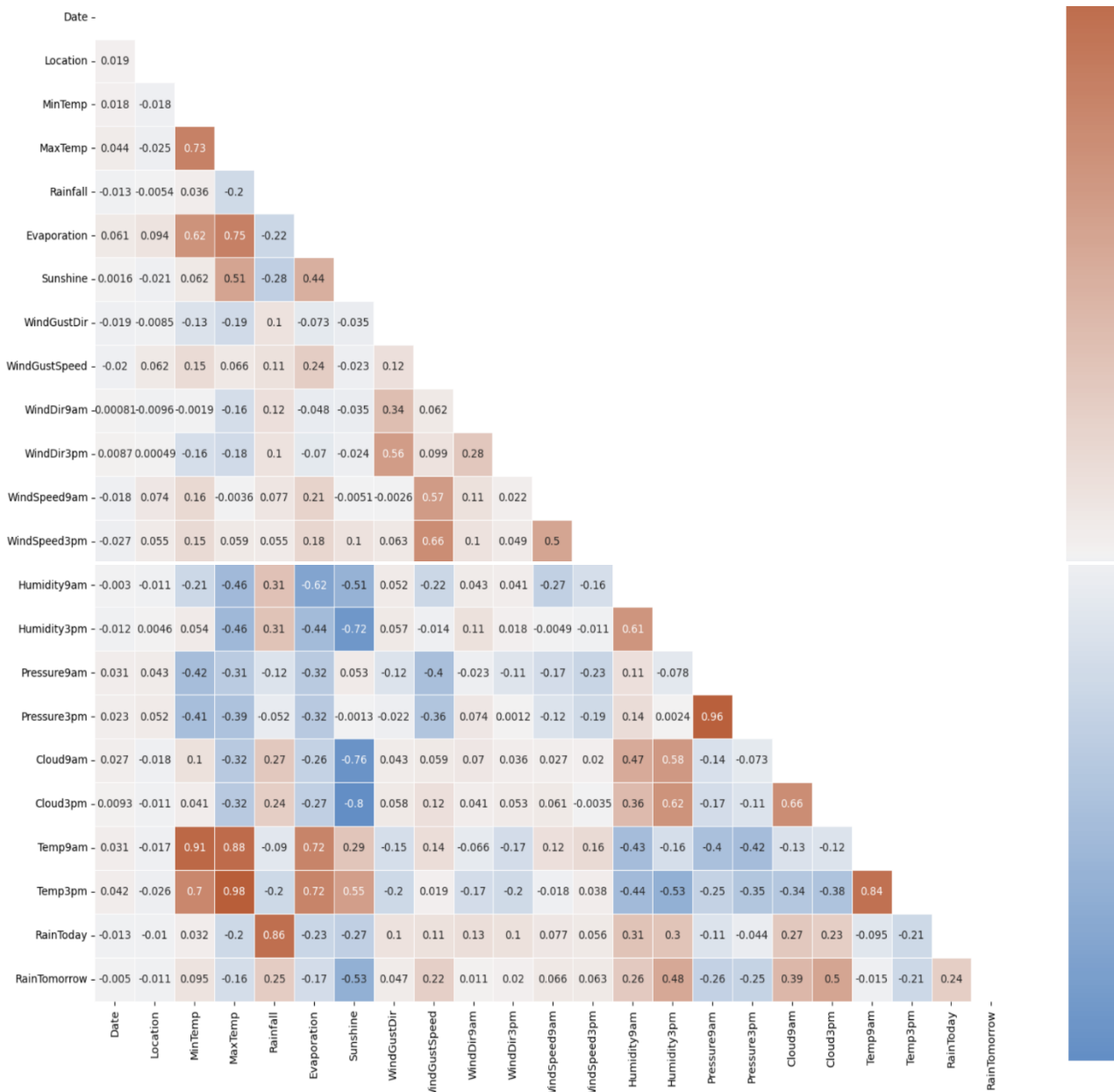


Figure 3.19 Correlation between attributes

Figure 3.19 shows the correlation between attributes. The following feature pairs have a strong correlation with each other:

- MaxTemp and Temp3pm
- Pressure3pm and Pressure9am
- Temp9am and MinTemp
- Temp9am and MaxTemp

But in no case is the correlation value equal to a perfect “1”. We are therefore not removing any functionality

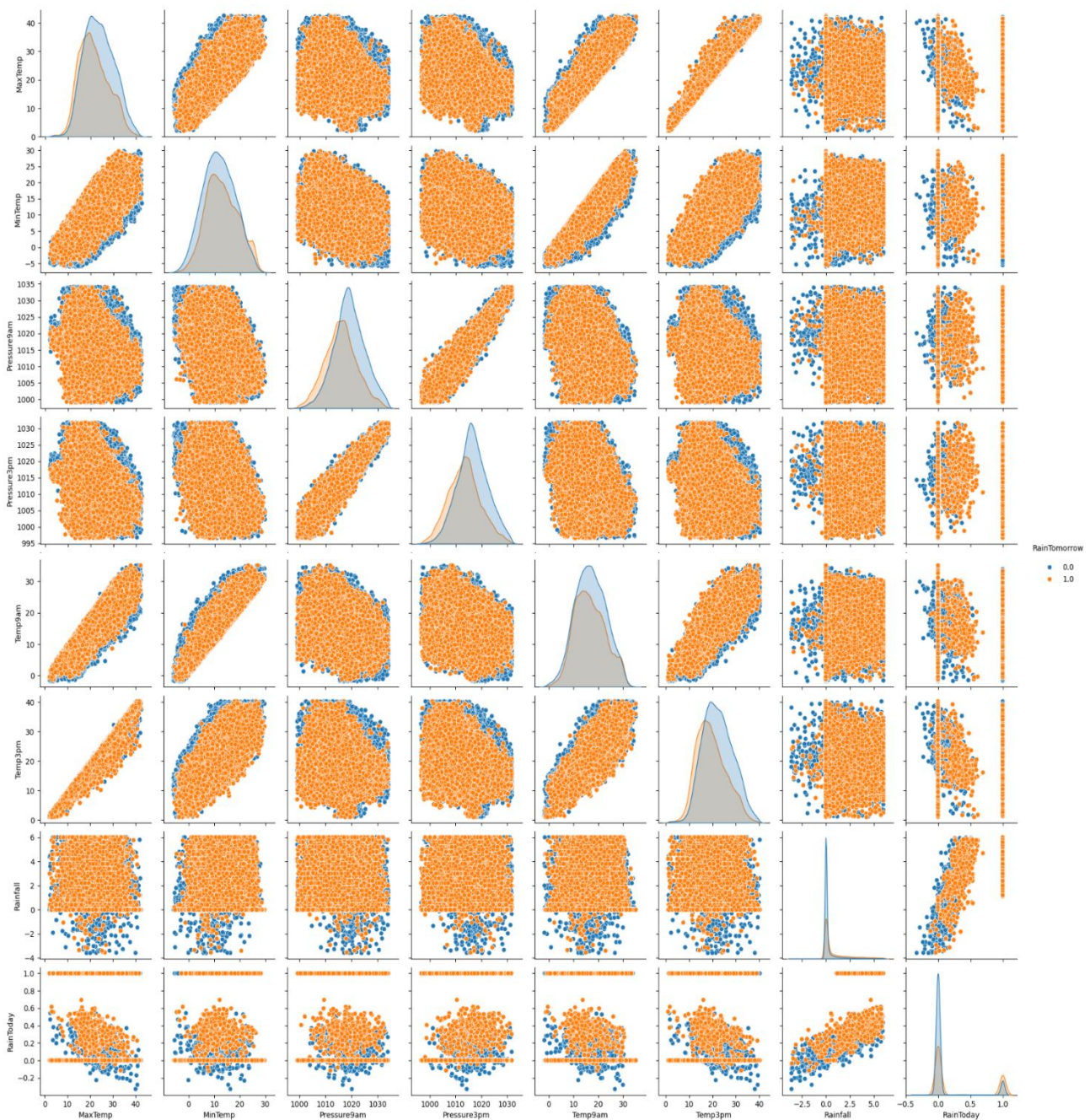


Figure 3.20 Pairplot

Figure 3.20 shows a pairplot. Each of the paired plots shows very clear distinct clusters of RainTomorrow's "yes" and "no" clusters. There is very minimal overlap between them.

3.4 Analysis of Algorithms

Algorithms used in Rainfall Prediction Using Machine Learning are Artificial Neural Networks, Support Vector Machines and Random Forest.

Artificial Neural Networks (ANN)

Artificial Neural Networks (ANNs) are highly flexible machine learning models designed to capture complex and nonlinear relationships in data, making them particularly well-suited for rainfall prediction tasks. In this project, ANN will be used to model intricate dependencies between meteorological variables such as temperature, humidity, wind speed, and atmospheric pressure, which collectively influence rainfall patterns. ANNs work by passing input features through multiple layers of interconnected nodes (neurons), each of which performs a weighted sum of inputs and applies an activation function to introduce nonlinearity. By training the ANN on the "WeatherAUS" dataset, the model learns subtle patterns in weather conditions that indicate whether it will rain the following day ("RainTomorrow"). Through backpropagation, the network adjusts weights iteratively to minimize prediction error, improving accuracy over time. ANNs are known for their robustness against noise and ability to generalize well to new data, making them ideal for capturing the variable nature of weather. The model's ability to learn from historical data and adapt to new patterns positions it as a promising approach for accurate rainfall prediction.

Pseudocode:

Initialize weights and biases

For each epoch:

For each input sample:

Perform forward propagation through all layers

Calculate the loss

Perform backpropagation to update weights and biases

End For

End For

Use the trained model to make predictions

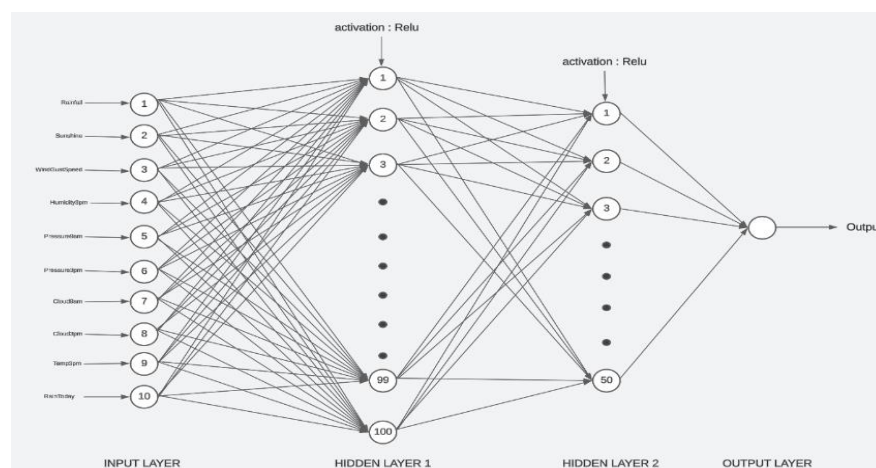


Figure 3.21 ANN model

Figure 3.21 shows an ANN model.

Support Vector Machines

Support Vector Machine (SVM) is a powerful classification algorithm designed to find the optimal hyperplane that best separates data points into distinct classes. In the context of this project, SVM will be applied to classify days into two classes based on whether it will rain the following day ("RainTomorrow" - yes or no), using a set of meteorological features such as temperature, humidity, and wind speed. SVM is particularly useful for datasets with high dimensionality and can handle both linear and nonlinear relationships by using kernel functions (e.g., radial basis function, polynomial). For rainfall prediction, SVM's strength lies in its ability to create a clear separation between classes by maximizing the margin, ensuring robust and reliable predictions. Additionally, SVM's effectiveness in minimizing overfitting, especially with noisy or limited data, makes it suitable for dealing with real-world meteorological datasets. By leveraging the power of SVM, this project aims to develop a model that accurately classifies weather data into rain and no-rain categories, offering reliable predictive capabilities for practical applications like agriculture and disaster management.

Pseudocode:

Transform the input data using a kernel function (if needed)

Find the optimal hyperplane that maximizes the margin

Identify support vectors

For each new input sample:

 Calculate the decision function based on the hyperplane

 Assign the sample to a class based on the sign of the decision function

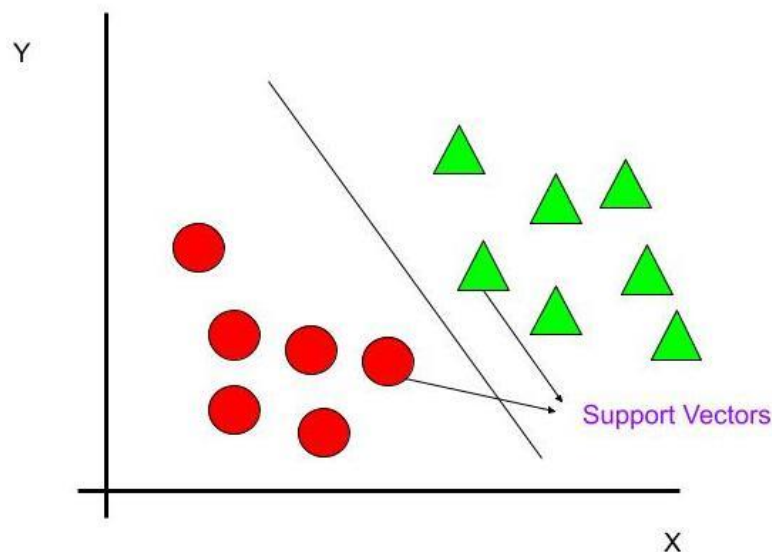


Figure 3.22 SVM working

Figure 3.22 shows how SVM classifies the data samples. In the above figure there are two classes, separated by a hyperline.

Random Forest

Random Forest (RF) is an ensemble learning method that constructs multiple decision trees and combines their outputs to make more stable and accurate predictions. In this project, Random Forest will be used to analyze weather features such as historical rainfall, wind speed, humidity, and atmospheric pressure to determine whether it will rain the following day. Each tree in the forest learns from a random subset of features, which helps to reduce overfitting and increase robustness against noisy or irrelevant data.. Random Forest is especially valuable in this project for its ability to calculate feature importance, allowing us to identify which meteorological parameters most influence the likelihood of rain. This interpretability makes RF a useful model not only for prediction but also for understanding the underlying factors driving rainfall patterns. By leveraging RF's strengths in handling large feature spaces, handling non-linear relationships, and providing reliable performance, the project aims to build a robust and interpretable rainfall prediction model that can aid in decision-making across various fields, including agriculture, water resource management, and weather forecasting.

Pseudocode:

For each tree in the Random Forest:

 Bootstrap sample the data

 Randomly select a subset of features for each split

 Build the decision tree using the selected features and bootstrap sample

 Repeat until the maximum depth or minimum sample size is reached

End For

For each new input sample

 Pass the sample through each decision tree in the forest

 Record the output (class or regression value) from each tree

 Aggregate the outputs:

- For classification: Use majority voting to determine the final class label
- For regression: Compute the average of the outputs to get the final prediction

 Return the aggregated result as the final prediction

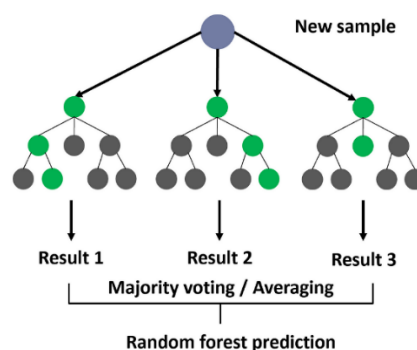


Figure 3.23 RF working

Figure 3.23 shows how RF works.

3.4.1 Activity Diagrams

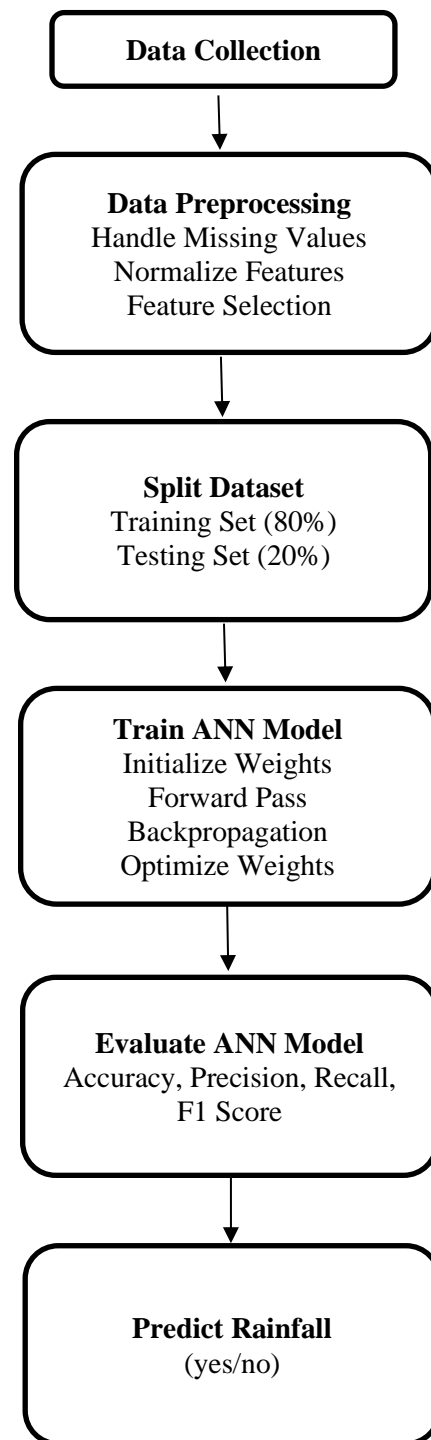


Figure 3.24 ANN

Fig 3.24 shows the activity diagram of ANN model

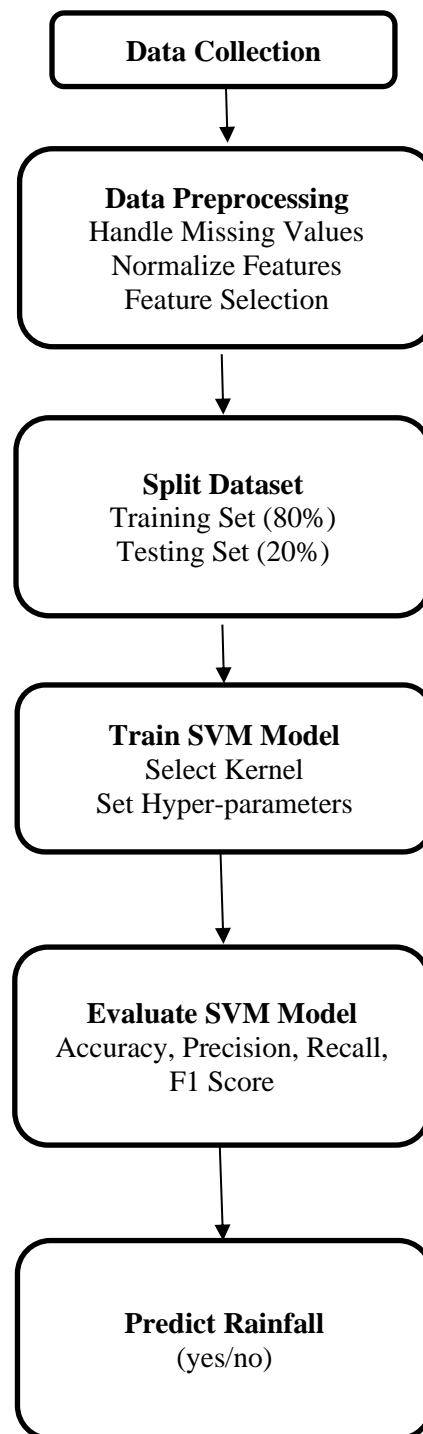


Figure 3.25 SVM

Fig 3.25 shows the activity diagram of SVM model

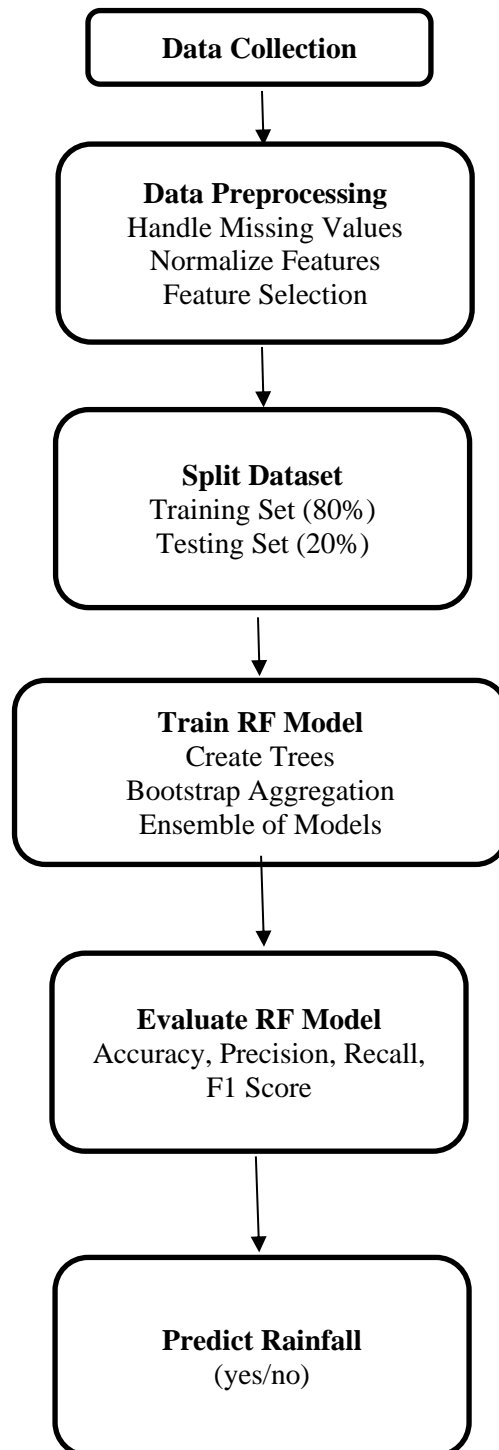


Figure 3.26 RF

Fig 3.26 shows the activity diagram of RF model

3.5 Project Plan

3.5.1 Project Pipeline

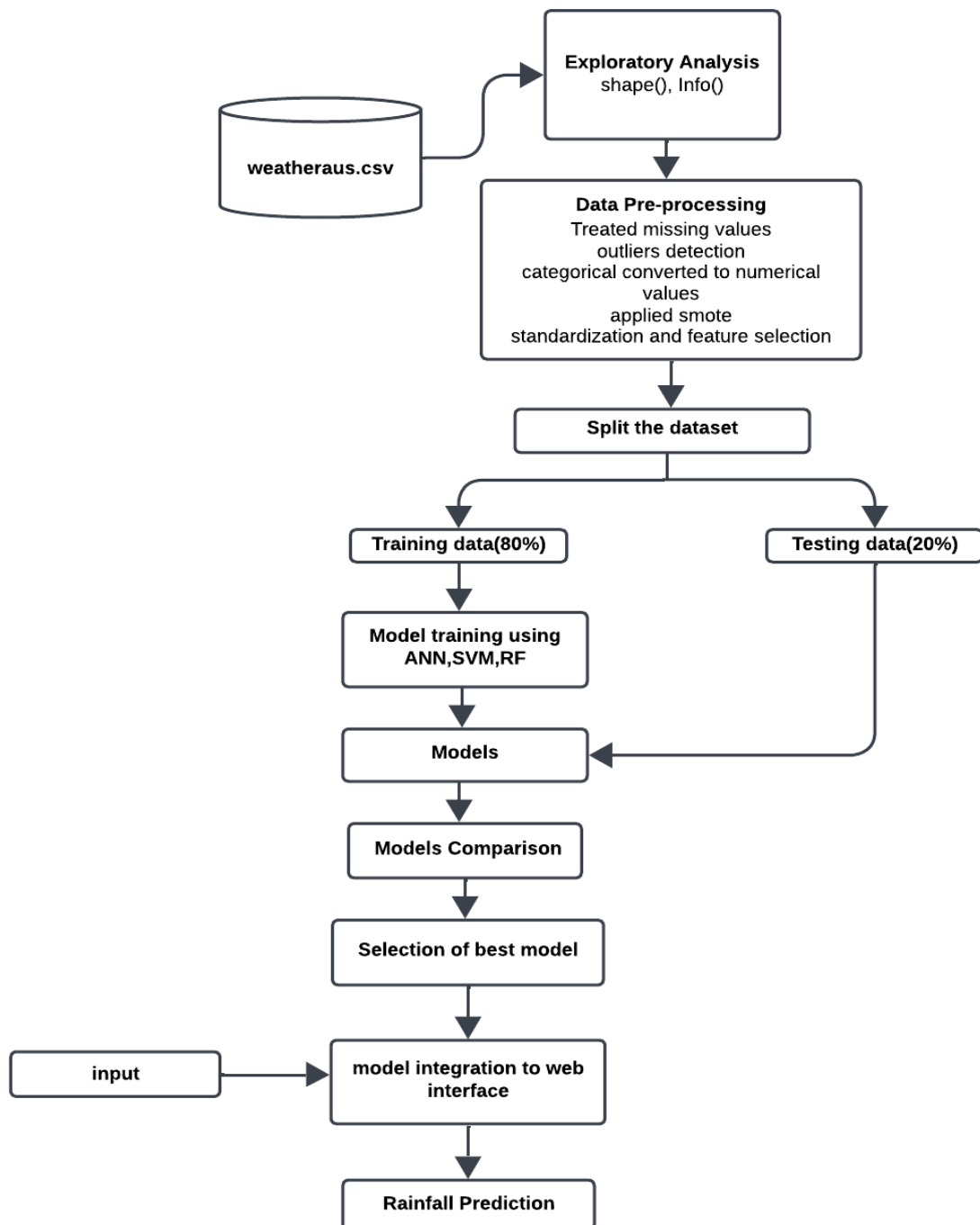


Figure 3.27 Project Pipeline

Figure 3.27 shows the project pipeline.

3.6 Feasibility Analysis

A feasibility study aims to objectively and rationally uncover the strengths and weaknesses of an existing system or proposed system, opportunities and threats present in the natural environment, the resources required to carry through, and ultimately the prospects for success.

Evaluated the feasibility of the system in terms of the following categories:

- Technical Feasibility
- Economic Feasibility
- Operational Feasibility

3.6.1 Technical Feasibility

Technical feasibility assesses whether the technology required for the project is readily available and if the team has the expertise to implement it effectively.

1. Data Processing and Analysis

- **Feasibility:** The project relies on standard data processing libraries like NumPy, Pandas, and scikit-learn for handling and analyzing weather data. These libraries are well-supported, widely used in data science projects, and suitable for handling the large datasets involved, ensuring technical feasibility.
- **Expertise:** Familiarity with these tools, as demonstrated by their use in previous projects, indicates that the team possesses the technical competency to handle data preprocessing, cleaning, and analysis effectively.

2. Machine Learning Algorithms

- **Feasibility:** The project employs Artificial Neural Networks (ANN), Support Vector Machines (SVM), and Random Forest (RF) for rainfall prediction. These algorithms are widely used and have proven effective in predictive modeling tasks, particularly in weather prediction scenarios, ensuring their feasibility for this project.
- **Expertise:** The successful implementation of these algorithms demonstrates the team's understanding of machine learning concepts and their proficiency in applying these methods to real-world datasets.

3. Data Visualization

- **Feasibility:** Visualization libraries such as Matplotlib and Seaborn are used to display data trends, making it easy to explore patterns in the weather data. Both are open-source and technically feasible for data visualization tasks.
- **Expertise:** The team's capability to use these tools effectively for data visualization suggests proficiency, enhancing the project's technical viability.

4. File Handling

- **Feasibility:** Using Pandas for reading and writing CSV files is a standard practice in data handling, ensuring technical feasibility for managing dataset inputs and outputs.
- **Expertise:** Proficiency in file handling and data manipulation demonstrates the team's competence, which is essential for efficiently managing the dataset throughout the project.

3.6.2 Economic Feasibility

Economic feasibility evaluates whether the project is financially viable, considering costs and potential benefits.

1. Development Costs

- **Feasibility:** By utilizing open-source tools and libraries, the project minimizes software acquisition costs. However, potential expenses may include personnel, training, and necessary hardware upgrades if needed for efficient model training.
- **Benefits:** An accurate rainfall prediction system could provide significant value in areas such as agriculture, disaster management, and environmental monitoring, justifying the development costs.

2. Maintenance Costs

- **Feasibility:** Since the project uses open-source libraries, maintenance costs are low, with occasional costs associated with updating models, refining algorithms, or addressing data drift.
- **Benefits:** The long-term benefits, such as accurate predictions that support agricultural planning or disaster preparedness, could justify the modest maintenance expenses involved.

3. Potential ROI

- The project can provide an ROI by offering a practical solution for rainfall prediction, potentially attracting partnerships or funding from agricultural, environmental, or governmental sectors that could benefit from accurate weather forecasting.

3.6.3 Operational Feasibility

Operational feasibility examines how well the project aligns with operational requirements and its potential for seamless integration into existing workflows.

1. User Acceptance

- **Feasibility:** The project aims to provide rainfall predictions that are easy for end-users (e.g., farmers, environmentalists) to understand and apply, increasing its likelihood of user acceptance.
- **Integration:** The web interface, built using Flask, is designed for accessibility and ease of use, making it suitable for users with minimal technical knowledge, thereby enhancing operational feasibility.

2. Scalability

- **Feasibility:** The project should consider scalability to accommodate a larger user base and potential expansion to include weather data from multiple regions. As the data volume and user base grow, the models and data infrastructure may need to be optimized.
- **Mitigation:** Implementing optimizations like feature selection or deploying models on a cloud infrastructure could help maintain efficiency as the project scales.

3. Ease of Use

- **Feasibility:** A user-friendly web interface will ensure ease of access for non-technical users, enhancing the operational feasibility of the project.
- **Training:** With intuitive design and straightforward data input fields, users should find the system easy to navigate, with minimal training requirements. This focus on usability will further improve operational adoption.

3.7 System Environment

3.7.1 Software Environment

The software environment encompasses the tools, frameworks, and platforms utilized in the development and execution of the forecasting system.

Programming Languages:

Python: The primary programming language employed for its versatility, extensive libraries (NumPy, Pandas, scikit-learn), and strong support in data science and machine learning.

Data Processing and Analysis:

NumPy and Pandas: Utilized for efficient data manipulation, handling, and analysis.

Machine Learning Libraries:

scikit-learn: Used for implementing and training machine learning models, including Artificial Neural Networks, Support Vector Machines and Random Forest.

SciPy: Complements NumPy and provides additional functionality for scientific computing.

Data Visualization:

Matplotlib and Seaborn: Chosen for creating insightful visualizations, aiding in the exploration and communication of data patterns.

File Handling:

Pandas: Applied for reading and writing data to CSV files, ensuring seamless data management.

Development Environment:

Google Colab

Colab is a free Jupyter notebook environment that runs entirely in the cloud. We can write and execute code in Python. Colab supports many machine learning libraries which can be easily loaded in the colab notebook. Visual Studio Code is a streamlined code editor with support for development operations like debugging, task running, and version control. It aims to provide just the tool a developer needs for a quick code-build-debug cycle and leaves more complex work flows to fuller featured IDEs, such as Visual Studio IDE.

HTML and CSS

Hyper Text Markup Language is used for creating web pages. HTML describes the structure of the web page. Here, the user interface of my project is done using HTML. Cascading Style Sheet is used with HTML to style the web pages.

Flask

Flask is a web framework, it's a Python module that lets you develop web applications easily. It's has a small and easy-to-extend core: it's a microframework that doesn't include an ORM (Object Relational Manager) or such features. It does have many cool features like URL routing, template engine. It is a WSGI web app framework.

GitHub

Git is an open-source version control system that was started by Linus Torvalds. Git is similar to other version control systems Subversion, CVS, and Mercurial to name a few. Version control systems keep these revisions straight, storing the modifications in a central repository. This allows developers to easily collaborate, as they can download a new version of the software, make changes, and upload the newest revision. Every developer can see these new changes, download them, and contribute. Git is the preferred version control system of most developers, since it has multiple advantages over the other systems available. It stores file changes more efficiently and ensures file integrity better.

The social networking aspect of GitHub is probably its most powerful feature, allowing projects to grow more than just about any of the other features offered. Project revisions can be discussed publicly, so a mass of experts can contribute knowledge and collaborate to advance a project forward.

3.7.2 Hardware Environment

The hardware environment refers to the physical infrastructure necessary to support the development and deployment of the recommendation system.

Computing Resources:

Processor: 2 GHz, A powerful multi-core processor is recommended to handle the computational demands efficiently. Consider a processor with at least quad- core architecture for optimal performance.

Storage:

Sufficient Disk Space: 512 GB SSD Memory (RAM):

8 GB RAM: Sufficient RAM is crucial for efficiently handling large datasets and performing machine learning tasks.

Internet Connectivity:

High-Speed Internet: A stable and high-speed internet connection is essential for accessing online resources, libraries, and datasets during development.

4. SYSTEM DESIGN

4.1 Model Building

4.1.1 Model Planning

The primary aim of the model planning phase in the rainfall prediction project is to create a structured roadmap for building an accurate predictive model. This phase involves defining clear objectives, selecting appropriate algorithms, and outlining strategies to process and analyze the WeatherAUS dataset effectively for rainfall prediction.

1. Objective

Develop a rainfall prediction model using Machine Learning algorithms that can accurately predict the likelihood of rainfall on the following day, based on various meteorological parameters provided in the WeatherAUS dataset.

2. Approach

Utilize Artificial Neural Networks (ANN), Support Vector Machines (SVM), and Random Forest (RF) as the core algorithms for predictive modeling. Each algorithm will provide unique insights:

- ANN: A neural network model to capture complex non-linear relationships in the data.
- SVM: A classification model to find the optimal hyperplane, effectively distinguishing between rainy and non-rainy days.
- RF: An ensemble-based model that aggregates multiple decision trees for robust, accurate predictions.

3. Data Preparation

Import and preprocess the WeatherAUS dataset. This includes handling missing values, encoding categorical variables, and scaling numerical features as needed to ensure that the dataset is clean and suitable for model training.

4. Exploratory Data Analysis (EDA)

Conduct a thorough analysis of the dataset to identify patterns and relationships between features. Key areas of focus include distribution of rainfall across different locations, identifying correlations between meteorological features, and examining trends over time. Visualizations will aid in understanding these relationships and refining the feature set.

5. Algorithm Implementation

Implement each of the chosen algorithms (ANN, SVM, RF) using scikit-learn and TensorFlow/Keras. The models will be trained on the preprocessed dataset, and hyperparameter tuning will be performed to optimize their performance for accurate rainfall prediction.

6. Evaluation and Comparison

Evaluate each model's performance using metrics such as accuracy, precision, recall, and F1-score. The models will be compared based on these metrics, and the best-performing model will be selected for final deployment.

This structured plan will ensure a methodical approach to building and evaluating the rainfall prediction model, leading to accurate and actionable insights.

Dataset is split in to two sets: 80% of datapoints for training and 20% for testing. `train_test_split` function is used to divide the data into training and testing sets.

```
# Split the data into train and test sets
X_new = modified_data.drop(columns=['RainTomorrow'])
y_remaining = modified_data['RainTomorrow']
X_train, X_test, y_train, y_test = train_test_split(X_new, y_remaining, test_size=0.2, random_state=42)
```

Figure 4.1 Code snippet for Dataset splitting

Figure 4.1 shows the code snippet used to split the dataset in to training and testing data.

```
# 1. Artificial Neural Network (ANN)
ann_model = MLPClassifier(random_state=42, verbose=1, hidden_layer_sizes=(100, 50), activation='relu',
                           solver='adam')
ann_model.fit(X_train, y_train)
ann_predictions = ann_model.predict(X_test)
```

Figure 4.2 Code snippet for ANN model training

Figure 4.2 shows the code snippet for training ANN model.

```
# 2. Support Vector Machine (SVM)
svm_model = SVC(random_state=42, C=1, kernel='linear')
svm_model.fit(X_train, y_train)
svm_predictions = svm_model.predict(X_test)
svm_accuracy = accuracy_score(y_test, svm_predictions)
```

Figure 4.3 Code snippet for SVM model training

Figure 4.3 shows the code snippet for training SVM model.

```
# 3. Random Forest
params_rf = {'max_depth': 17, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 100,
             'random_state': 12345}
model_rf = RandomForestClassifier(**params_rf)
model_rf.fit(X_train, y_train)
rf_predictions = model_rf.predict(X_test)
rf_accuracy = accuracy_score(y_test, rf_predictions)
```

Figure 4.4 Code snippet for RF model training

Figure 4.4 shows the code snippet for training RF model.

4.1.2 Training

The training phase in the rainfall prediction project involves constructing prediction models using Artificial Neural Networks (ANN), Support Vector Machines (SVM), and Random Forest (RF) algorithms. The WeatherAUS dataset is split into training and testing sets to evaluate model performance effectively. After data preprocessing, feature matrices and labels are created for model training. Each algorithm is trained using the entire training set to learn patterns between meteorological parameters and rainfall occurrence.

- **ANN:** The ANN model is trained with multiple layers, using backpropagation to minimize the loss and improve prediction accuracy.
- **SVM:** The SVM model is trained to find the optimal hyperplane that separates rainy and non-rainy days by maximizing the margin between classes.
- **RF:** The Random Forest model is constructed by training multiple decision trees on bootstrapped samples of the training set, aggregating their predictions to enhance accuracy and reduce overfitting.

Throughout the training phase, hyperparameter tuning is performed for each algorithm to identify the optimal configuration that maximizes predictive performance.

4.1.3 Testing

The testing phase evaluates the ability of the ANN, SVM, and RF models to generalize to new data. The test set is used to assess model performance on unseen data, providing insights into the accuracy, precision, recall, and F1-score for each algorithm. This evaluation helps determine each model's robustness in predicting rainfall.

- **ANN Testing:** The ANN model's performance is evaluated based on its accuracy and error rate on the test set, assessing how well it captures non-linear relationships.
- **SVM Testing:** The SVM model's effectiveness is measured by its precision and recall, which determine how well it can classify rainy versus non-rainy days.
- **RF Testing:** The RF model's results are analyzed to check for overfitting, with particular attention to its ability to generalize due to its ensemble nature.

The iterative refinement process involves adjusting model parameters and retraining based on testing results, enhancing each model's performance. This validation ensures that the models align with the project's objectives, providing reliable and accurate rainfall predictions.

5. RESULTS AND DISCUSSION

The results of the rainfall prediction models—Artificial Neural Network (ANN), Support Vector Machine (SVM), and Random Forest (RF)—demonstrate varying levels of performance in terms of accuracy, precision, recall, and F1 score.

The ANN model achieved an accuracy of 87.61% and a training accuracy of 88.01% , indicating a strong alignment between training and testing performance. With a precision of 85.96% and recall of 89.99% , the ANN model demonstrates a balanced trade-off between correctly predicting rainfall events and minimizing false positives. The F1 score of 87.93% confirms the ANN's reliability in handling both precision and recall, making it an effective model for rainfall prediction where a slight bias towards recall is acceptable to ensure rain events are not missed.

The SVM model achieved an accuracy of 80.77% with a training accuracy of 81.17% , suggesting a consistent but slightly lower performance than the ANN model. The precision of 80.75% and recall of 80.97% indicate that the SVM model has a balanced prediction capability but slightly underperforms in identifying rainy days compared to the ANN and RF models. The F1 score of 80.86% shows that while the SVM model is capable, it may not be as robust as ANN and RF for this specific task, particularly in capturing complex relationships in the data.

The Random Forest model stands out with the highest accuracy of 90.13% on the test set, although it exhibits a significantly higher training accuracy of 97.82%, which could suggest a slight tendency towards overfitting. With a precision of 90.19% and recall of 90.13% , the RF model achieves an optimal balance between correctly identifying rainy days and avoiding false positives. The F1 score of 90.16% confirms RF as the most reliable model among the three, effectively handling both precision and recall requirements for rainfall prediction. The ensemble nature of Random Forest enables it to capture complex data patterns while maintaining robustness in prediction.

Among the three models, Random Forest (RF) achieved the highest accuracy, precision, recall, and F1 score, indicating its superior predictive capability for rainfall prediction. Its ensemble approach likely contributes to capturing complex patterns in the meteorological data. However, the noticeable gap between training and testing accuracy suggests potential overfitting, warranting careful tuning or validation for real-world application. The Artificial Neural Network (ANN) also performed well, particularly in recall, indicating its ability to identify rain events with high accuracy. This makes ANN a strong candidate where identifying potential rainy days is a priority, even if it means allowing some false positives. The Support Vector Machine (SVM) model, while effective, demonstrated slightly lower performance in all metrics compared to ANN and RF. Its simpler decision boundary might be less suited to the complex nature of weather data, resulting in its relatively lower accuracy and F1 score.

In conclusion, Random Forest appears to be the best-performing model for this rainfall prediction task, with ANN also showing strong performance and providing a viable alternative when slightly higher recall is preferred. The SVM model, though less effective in this context, still provides valuable insights and may be useful when computational efficiency is a priority.

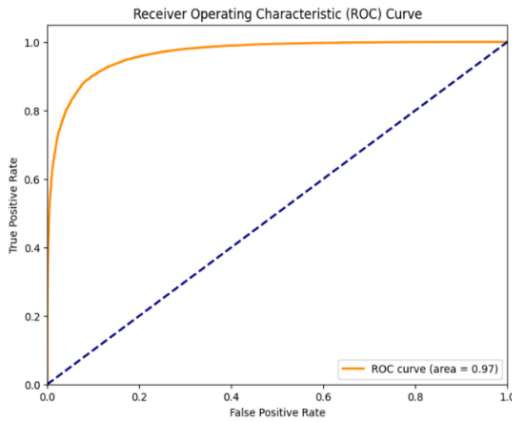


Figure 5.1 ROC Curve

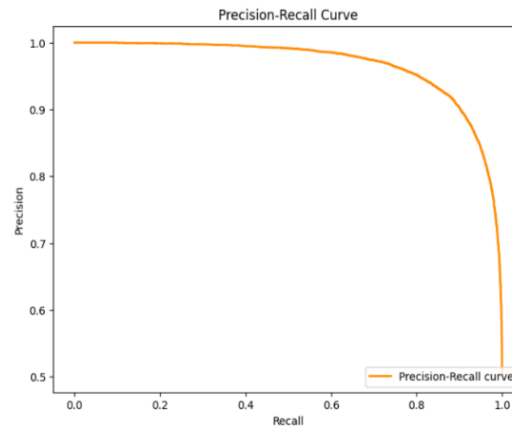


Figure 5.2 Precision-Recall Curve

Figure 5.1 and 5.2 shows the ROC curve and Precision-Recall curve for Random Forest model .

The ROC curve displays the trade-off between the true positive rate (sensitivity) and the false positive rate (1-specificity) across different thresholds. The blue dashed line represents a random classifier, while the orange curve represents the model's performance. The closer the curve follows the left-hand border and the top border, the better the model's performance. The area under the ROC curve (AUC = 0.97) is a performance metric, where a value close to 1 indicates excellent classification capability.

This Precision-Recall Curve shows the relationship between precision (the proportion of positive identifications that were actually correct) and recall (the proportion of actual positives identified correctly).

In the case of imbalanced datasets, the precision-recall curve provides a more informative view of a model's performance. A curve close to the top right corner indicates better performance, meaning high precision and recall.

Both plots suggest that the Random Forest model performs well, with a high degree of accuracy and a good balance between precision and recall.

```
# user input
selected_columns = ['Sunshine', 'WindGustSpeed', 'Humidity9am', 'Humidity3pm',
                    'Pressure9am', 'Pressure3pm', 'Cloud9am', 'Cloud3pm', 'Temp3pm',
                    'RainToday']
input_data = []
for feature in selected_columns:
    value = float(input(f"Enter value for {feature}: "))
    input_data.append(value)

input_df = pd.DataFrame([input_data], columns=selected_columns)

normalized_input = r_scaler.transform(input_df)

prediction = model_rf.predict(normalized_input)[0]

if prediction == 0:
    print("Prediction: No rain tomorrow")
else:
    print("Prediction: Rain tomorrow")
```

Figure 5.3 Code snippet for rainfall prediction

```
Enter value for Sunshine: 4.30
Enter value for WindGustSpeed: 20
Enter value for Humidity9am: 99
Enter value for Humidity3pm: 58
Enter value for Pressure9am: 1021
Enter value for Pressure3pm: 1019.3
Enter value for Cloud9am: 8
Enter value for Cloud3pm: 5.8
Enter value for Temp3pm: 15.6
Enter value for RainToday: 0
Prediction: No rain tomorrow
```

Figure 5.4 Rainfall Prediction

Figure 5.3 and 5.4 shows the code snippet for predicting whether it rains in 24 hours or not and its corresponding output.

6. MODEL DEPLOYMENT

The primary objective of model deployment in this project is to integrate the Random Forest model, developed during the training phase, into a production environment. By deploying this predictive model, the project aims to deliver an accessible solution for forecasting rainfall, catering to the needs of users such as farmers, weather analysts, and other stakeholders.

To facilitate deployment, the project utilizes Flask as the backend framework, integrating the scikit-learn library for model handling and prediction. Flask provides a lightweight and effective solution for deploying machine learning models as APIs, allowing seamless interaction between the model and the user interface.

User Interface (UI):

The rainfall prediction model is made accessible through a simple web interface, designed for straightforward interaction. The interface includes:

- **Input Fields:** Users can enter relevant weather data (e.g., temperature, humidity, wind speed) to receive rainfall prediction.
- **Predict Button:** Upon entering the data, users click the predict button to obtain the forecast result.
- **Output Display:** The prediction outcome (rain/no rain) is displayed on another page, allowing users to quickly interpret the forecast.

This design ensures a user-friendly experience, making the application accessible to individuals with various levels of technical expertise. The web interface is optimized for clarity, with easy navigation and clean design, supporting practical usage in real-time scenarios.

Figure 6,1 shows the home page. After clicking the button “Get Started” it will navigate to the next page, which is index.html , shown in the figure 6.2 , which allows user to input the data to get the prediction result.

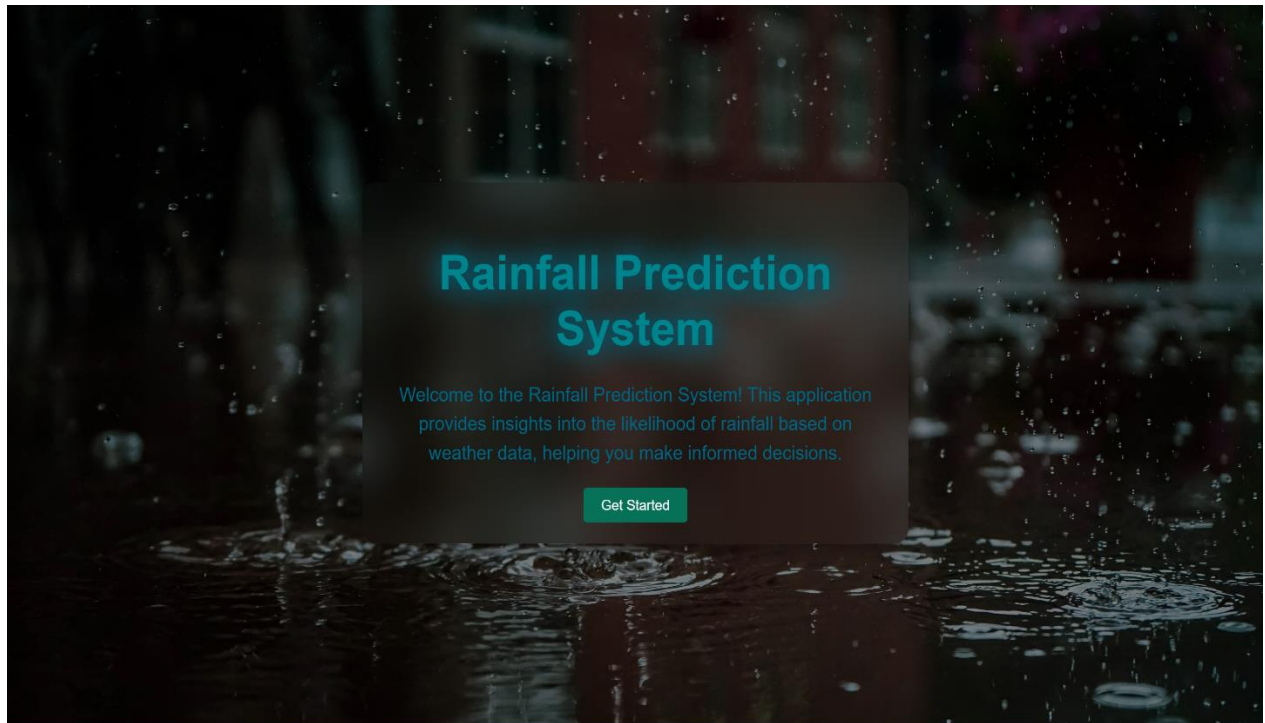


Figure 6.1 Home Page

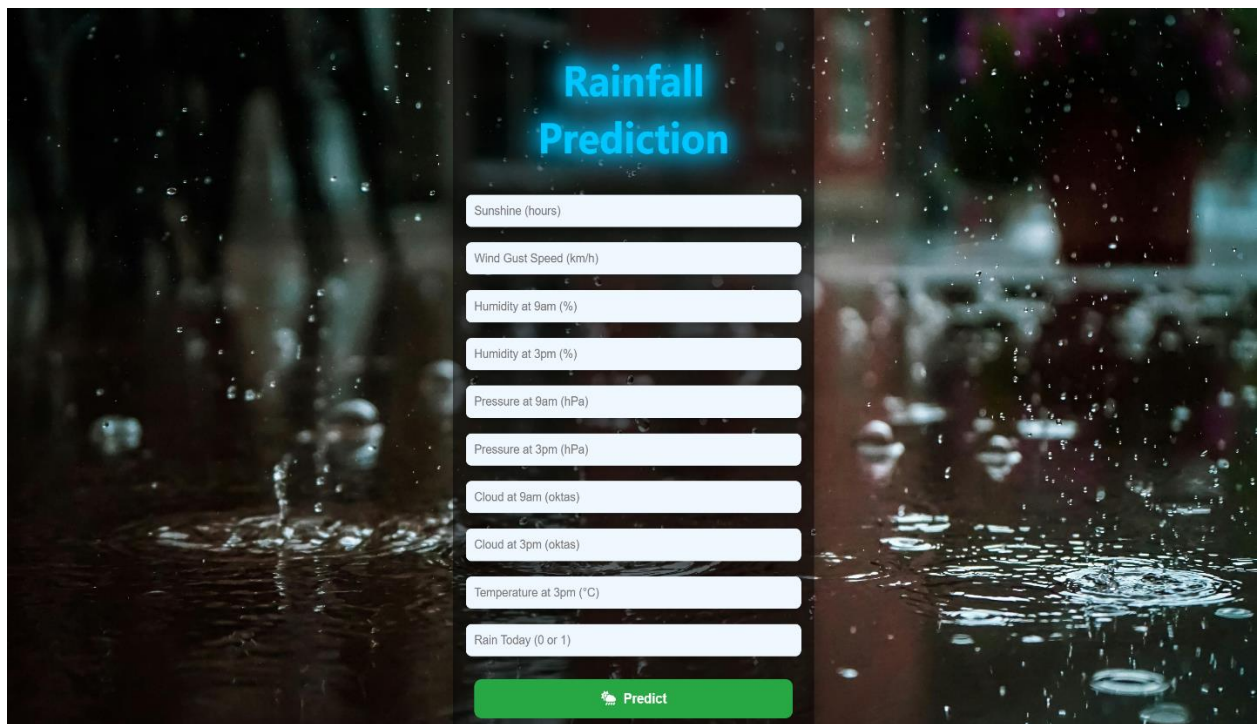
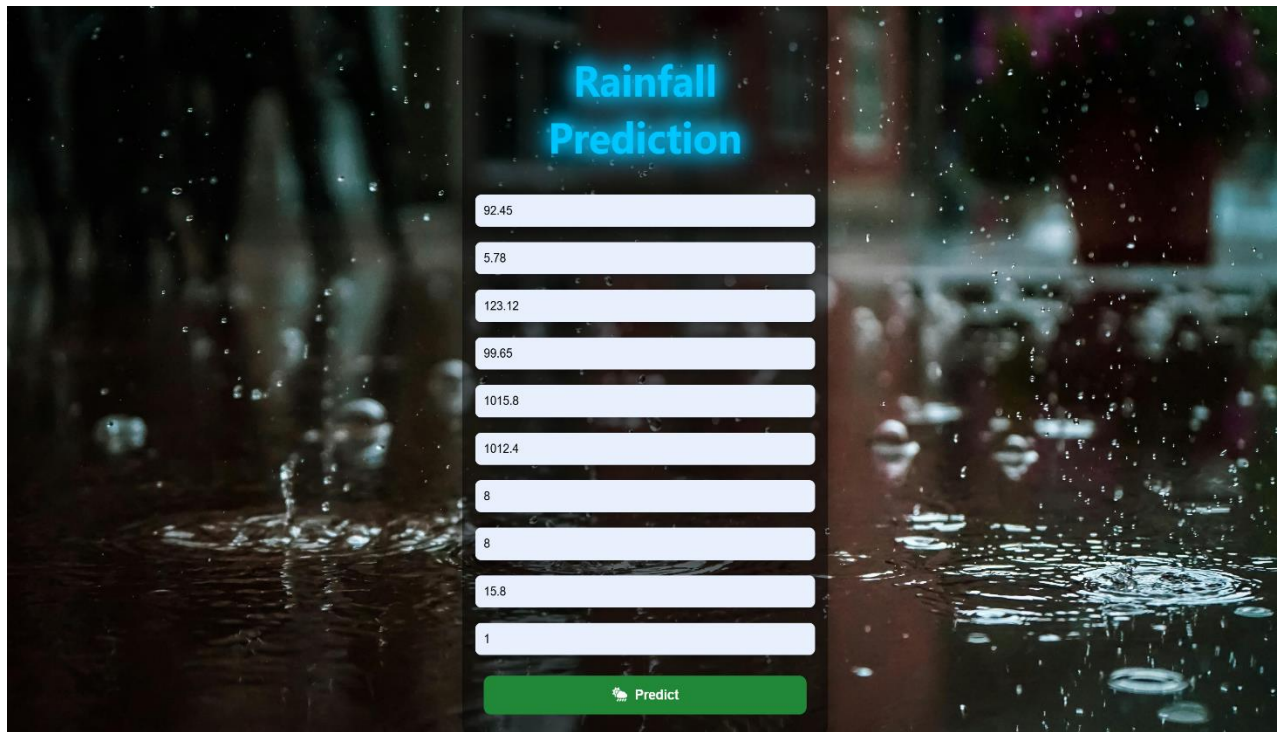


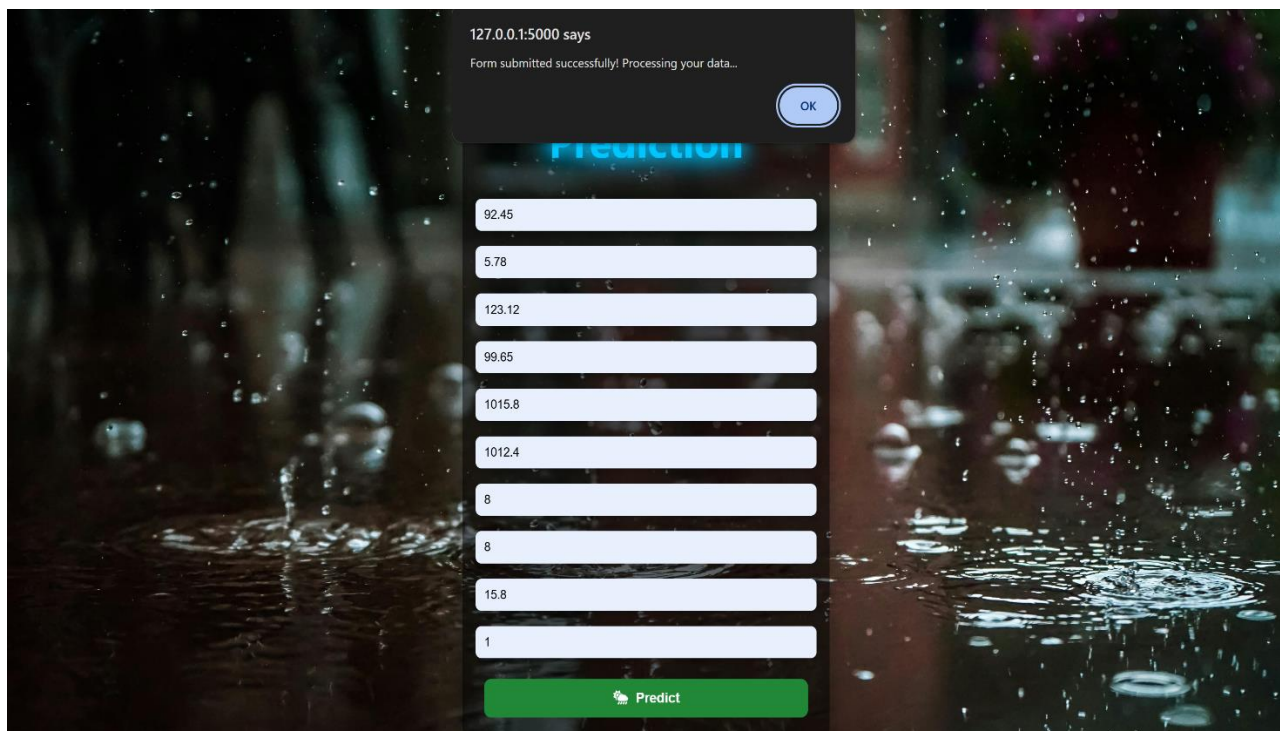
Figure 6.2 Index Page



The screenshot displays a mobile application interface titled "Rainfall Prediction" in blue text. Below the title, there are ten white input fields, each containing a numerical value. The values are: 92.45, 5.78, 123.12, 99.65, 1015.8, 1012.4, 8, 8, 15.8, and 1. At the bottom of the form is a green button with a white icon of a person and the text "Predict". The background of the app is a dark, blurred image of rain falling on a surface.

Figure 6.3 User input

Figure 6.3 shows the user input. After clicking button “Predict” button ,it shows an alert shown in the figure 6.4.



The screenshot shows the same mobile application interface as Figure 6.3, but with an alert message displayed at the top. The alert text reads: "127.0.0.1:5000 says" followed by "Form submitted successfully! Processing your data...". There is an "OK" button in the top right corner of the alert. The input fields and the "Predict" button are still visible below the alert. The background remains the same dark, blurred image of rain.

Figure 6.4 Alert while submitting the input

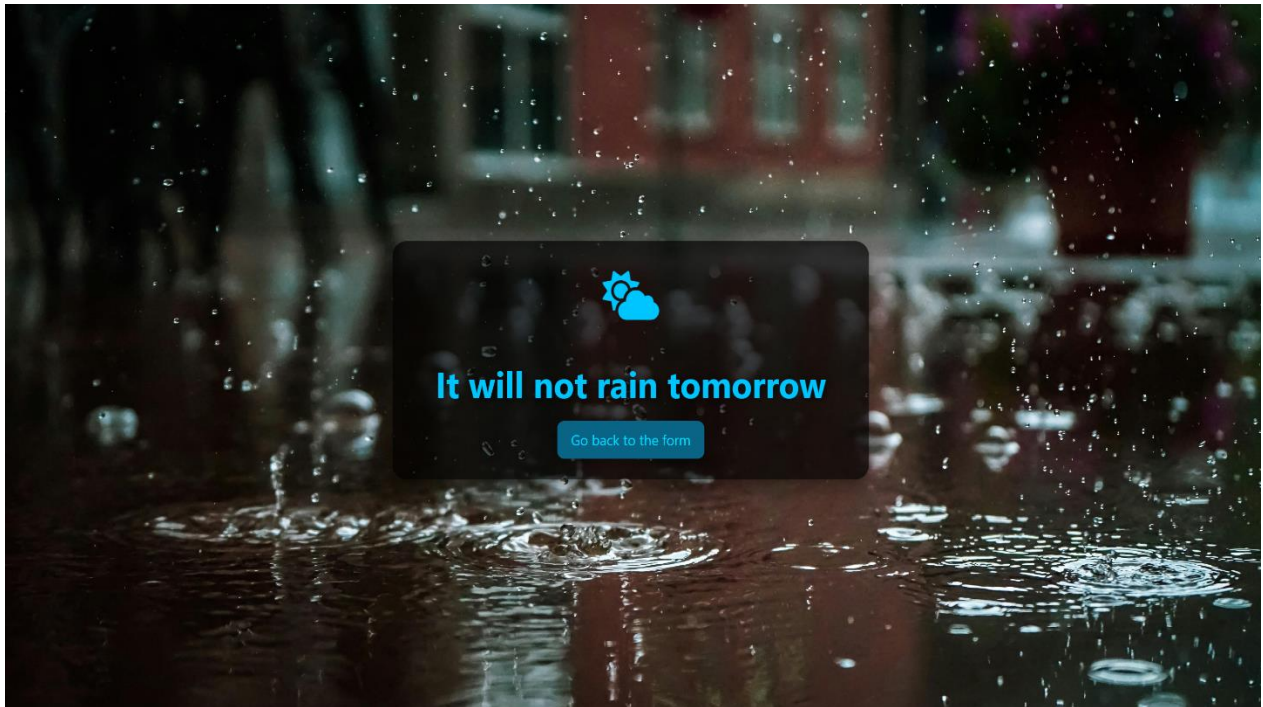


Figure 6.5 Prediction

Figure 6.5 shows the prediction for the input data.

7. GIT HISTORY

Git Repository Rainfall Prediction System contains all colab files, py files, html files and three related research papers. It is maintained for systematic way of project presentation and mainly for future reference. The colab files are created separately for three sprint releases during the project.

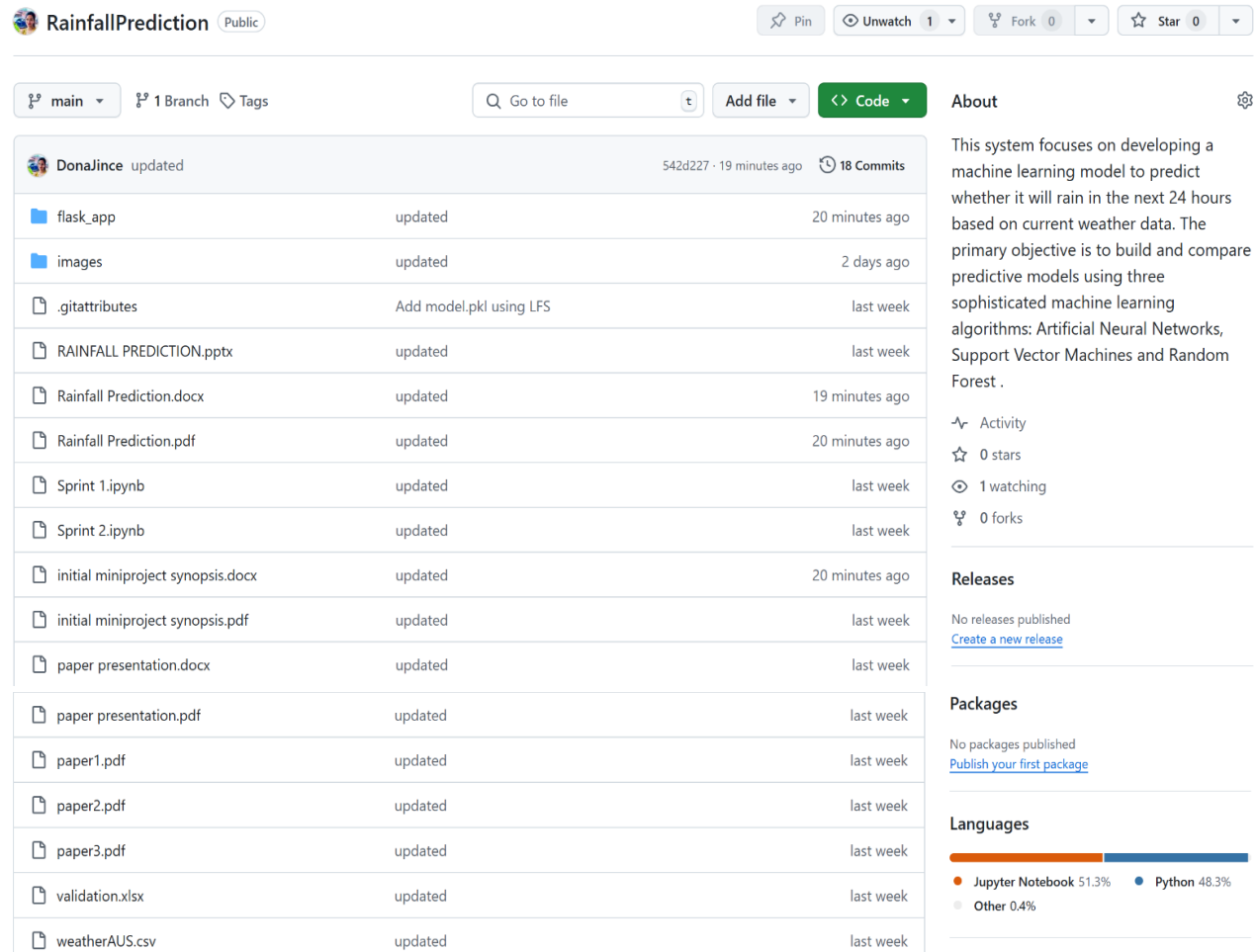





Figure 7.1 Git History

Figure 7.1 shows the git history of Rainfall Prediction System. Figure 7.2 shows the git history of sprint 3.

RainfallPrediction / flask_app / 

 DonaJince updated 9f6b4fe · 24 minutes ago  History







Name	Last commit message	Last commit date
 ..		
 templates	updated	24 minutes ago
 venv	updated	last week
 app.py	updated	24 minutes ago
 minmax_scaler.pkl	updated	last week
 model.pkl	Add model.pkl using LFS	last week

Figure 7.2 Git History of Sprint 3

8. CONCLUSION

The Rainfall Prediction System project represents a significant step forward in utilizing machine learning to address complex environmental forecasting challenges. Through the development and comparative analysis of three machine learning models—Artificial Neural Networks (ANN), Support Vector Machines (SVM), and Random Forests (RF)—the project has successfully demonstrated the capability of data-driven approaches to accurately predict rainfall. Among the models tested, Random Forest emerged as the most effective, delivering the highest accuracy and robustness, which ultimately led to its selection for deployment.

The project's core accomplishment lies in transforming raw weather data into actionable insights, a task achieved through meticulous data preprocessing, feature selection, and model optimization. By employing these techniques, the Rainfall Prediction System can support timely and informed decision-making for various applications, particularly in sectors such as agriculture, water resource management, and disaster preparedness. Predicting rainfall with accuracy can greatly influence crop planning, irrigation scheduling, flood prevention, and other weather-dependent activities, underscoring the practical relevance of this project. Furthermore, the deployment of the Random Forest model into a user-friendly interface exemplifies the project's commitment to accessibility and usability. The interface design allows users to interact seamlessly with the prediction system, making it easy for both experts and non-experts to interpret and utilize the forecast results. The successful deployment marks the transition of this project from a theoretical model to a practical tool, providing end-users with the power of machine learning at their fingertips.

This project also emphasizes the scalability and adaptability of machine learning in handling diverse datasets and complex predictive tasks. Lessons learned from this project, such as the importance of feature selection and data scaling, are instrumental for future iterations and applications. The challenges encountered, such as model tuning and deployment considerations, have provided valuable insights that will inform ongoing improvements.

In summary, the Rainfall Prediction System is more than a proof of concept; it is a practical application of machine learning that holds real-world value. By bridging the gap between predictive accuracy and user accessibility, this project has laid a solid foundation for future innovations in rainfall forecasting. As a data-driven solution, it not only meets current needs but also sets the stage for future advancements in weather prediction technologies, contributing to a more resilient and informed society.

9. FUTURE WORK

To further improve and expand the Rainfall Prediction System, several key areas of future work have been identified:

- **Advanced Machine Learning Models:** Future work could explore more sophisticated models, including deep learning architectures like LSTMs and CNNs, which may capture temporal patterns in weather data more effectively, potentially improving prediction accuracy.
- **Incorporation of Real-Time Data:** Integrating real-time weather data sources could enhance the model's accuracy and relevance, enabling up-to-the-minute predictions. This would be particularly beneficial in regions with rapidly changing weather conditions.
- **Dynamic Model Updating:** To maintain accuracy over time, future implementations could employ dynamic model updating, retraining periodically as new data becomes available. This would help the model adapt to evolving weather patterns caused by climate change.
- **Explainability and Transparency:** Making the model's predictions more interpretable would increase user trust. Adding explainability tools, such as SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations), would allow users to understand the factors influencing each prediction.
- **Scalability and Cloud Integration:** As data volumes grow, optimizing the system for scalability will be crucial. Deploying the model on cloud platforms with distributed computing capabilities can ensure that it remains responsive and efficient, even with increased user demand.
- **User Feedback Mechanisms:** Implementing a feedback loop where users can provide feedback on the model's predictions would allow for continuous improvement. This data could be leveraged to iteratively fine-tune the model, making it more adaptive and responsive to real-world needs.
- **Regional-Specific Models:** Developing models tailored to specific regions could improve accuracy by taking into account localized weather patterns. This would be particularly useful in areas where unique climate conditions influence rainfall.
- **Integration with IoT for Precision Agriculture:** Future iterations could consider integrating with IoT devices on farms, allowing for hyper-local weather data collection. This would enable the system to deliver even more targeted and actionable predictions for agricultural purposes.

In summary, by incorporating advanced modeling techniques, real-time data, explainability, scalability optimizations, user feedback, and regional specificity, the **Rainfall Prediction System** can evolve into a more adaptive, user-centric platform. These improvements will enhance its value to users and its potential as a tool for informed decision-making in weather-sensitive fields.

10. APPENDIX

10.1 Minimum Software Requirements

To deploy and run the Rainfall Prediction System, the following minimum software requirements must be met:

- Python:

Version 3.6 or above is required to execute the project code and its dependencies.

- NumPy:

Essential for numerical operations and array manipulations in Python.

- Pandas:

Necessary for data manipulation, cleaning, and analysis.

- Scikit-learn:

The scikit-learn library is used for machine learning functionalities, including the implementation of the Artificial Neural Networks, Support Vector Machines and Random Forest algorithms.

- Matplotlib and Seaborn:

Required for data visualization and generating plots and charts.

- Pickle:

Used for serialization and deserialization of Python objects, particularly for saving and loading machine learning models.

- Flask

A lightweight web framework to serve the application and provide the user interface. Flask is required to manage HTTP requests and deliver predictions via the web interface.

10.2 Minimum Hardware Requirements

The Rainfall Prediction System is relatively lightweight, but it does involve some computational processes. The minimum hardware requirements to ensure smooth execution are as follows:

Processor (CPU):

A modern multi-core processor (e.g., Intel Core i5 or equivalent) is recommended to handle the computation involved in the forecasting algorithms efficiently.

RAM:

A minimum of 8 GB RAM is recommended to manage the dataset and perform machine learning operations smoothly.

Storage:

Adequate storage space for storing datasets, code, and any additional resources. While the system itself doesn't require extensive storage, having sufficient space for datasets and potential model persistence is essential.

Internet Connection:

An internet connection is necessary for downloading datasets, libraries, and dependencies during the setup phase.

11. REFERENCES

1. Hassan MM, Rony MA, Khan MA, Hassan MM, Yasmin F, Nag A, Zarin TH, Bairagi AK, Alshathri S, El-Shafai W. Machine learning-based rainfall prediction: Unveiling insights and forecasting for improved preparedness. IEEE Access. 2023 Nov 16;11:132196-222.
2. Sarasa-Cabezuelo A. Prediction of rainfall in Australia using machine learning. Information. 2022 Mar 24;13(4):163.
3. Ghosh S, Gourisaria MK, Sahoo B, Das H. A pragmatic ensemble learning approach for rainfall prediction. Discover Internet of Things. 2023 Oct 9;3(1):13.<https://doi.org/10.1007/s43926-023-00044-3>
4. Kaggle, WeatherAUSdataset:
<https://www.kaggle.com/datasets/trisha2094/weatheraus>
5. <https://www.kaggle.com/code/midouazerty/rainfall-prediction-with-6-machine-learn-algo-98#Plotting-Decision-Region-for-all-Models>

