

## **Assessment Report**

on

# "Model to predict whether it will rain tomorrow using classification algorithms and weather data."

submitted as partial fulfillment for the award of

## BACHELOR OF TECHNOLOGY DEGREE

**SESSION 2024-25** 

in

CSE(AI)-B

Ву

Ayush Gupta - 202401100300084

Ashish Yadav-202401100300076

Ayush Prasad - 202401100300086

Himanshu Singh - 202401100300125

## Under the supervision of

"Shivansh Prasad"

KIET Group of Institutions, Ghaziabad

May, 2025

1. Introduction

With the increasing impact of weather on agriculture, travel, and daily planning, accurate rainfall prediction has become a vital component of meteorology. This project focuses on using **supervised machine learning** to build a model that predicts whether it will rain tomorrow. By leveraging historical weather data such as temperature, humidity, wind, and rainfall patterns, we aim to help users and organizations make informed decisions.

#### 2. Problem Statement

To **predict whether it will rain tomorrow** (Yes/No) using past weather data. The classification model can assist in preparing for weather-related risks and enhance decision-making in agriculture, event planning, and infrastructure management.

## 3. Objectives

- Preprocess the weather dataset for use in a machine learning pipeline.
- Train a **Logistic Regression** model to predict rainfall occurrence.
- Evaluate the model's performance using standard classification metrics.
- Visualize the results with a confusion matrix heatmap for better understanding.

#### 4. Methodology

#### 4.1 Data Collection

• The dataset used includes weather observations from various locations in Australia. It contains features like temperature, humidity, wind speed, pressure, and rainfall measurements.

#### 4.2 Data Preprocessing

- Handling Missing Values: Numerical values imputed with mean; categorical with mode.
- **Encoding**: Used **one-hot encoding** for categorical variables (e.g., WindDirection, Location).
- Scaling: Applied StandardScaler to normalize features.
- Splitting: The dataset is divided into 80% training and 20% testing sets.

#### 4.3 Model Building

• A **Logistic Regression** classifier is trained on the preprocessed data to predict the binary outcome: **RainTomorrow** (Yes/No).

#### 4.4 Model Evaluation

- Evaluated using the following metrics:
  - Accuracy
  - Precision
  - o Recall
  - o F1-Score
- A confusion matrix is created and visualized using **Seaborn heatmap**.

## 5. Data Preprocessing

- Numerical Features: Missing values filled with column mean.
- Categorical Features: Missing values filled with mode, then one-hot encoded.
- **Standardization**: All numerical features scaled to have zero mean and unit variance.
- Train-Test Split: Dataset split into 80% training and 20% testing.

#### 6. Model Implementation

A **Logistic Regression** model is selected due to its effectiveness in binary classification and ease of implementation. The model is trained on the processed training set and then used to predict the **RainTomorrow** label on the test set.

### 7. Evaluation Metrics

- Accuracy: Proportion of total correct predictions.
- **Precision**: Correctly predicted rain days out of all predicted rain days.
- Recall: Actual rain days correctly predicted.
- **F1 Score**: Harmonic mean of precision and recall.
- Confusion Matrix: Heatmap visualization to interpret prediction outcomes.

#### 8. Results and Analysis

- The model showed reasonable prediction capability on the test dataset.
- The **confusion matrix** indicated the trade-off between predicting rain (true positives) and missing actual rain days (false negatives).
- **Precision and recall** revealed the effectiveness of the model in identifying rainy days with minimized false alarms.

#### 9. Conclusion

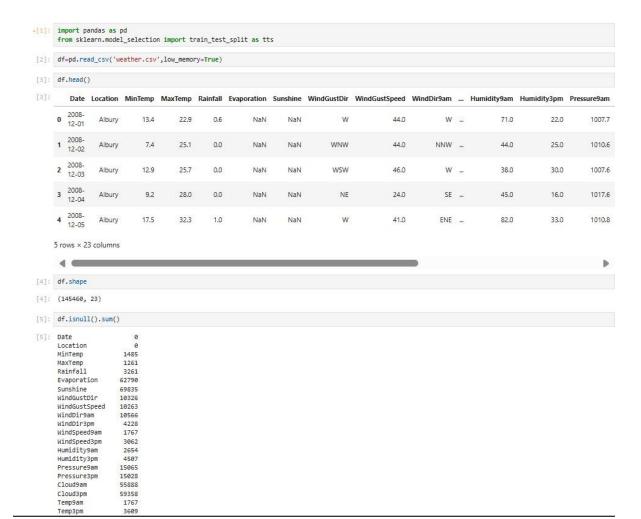
The Logistic Regression model effectively predicted the likelihood of rain on the next day using historical weather data. This project highlights the potential of machine learning in **weather forecasting** and supports automation in climate-sensitive domains.

#### **Future Improvements:**

- Handle class imbalance (e.g., SMOTE or resampling).
- Experiment with advanced models like Random Forest, Logistic Regression.
- Incorporate time-series components or real-time weather feeds for enhanced accuracy.

#### 10. References

- Kaggle Weather Dataset
- Scikit-learn Documentation
- Pandas Documentation
- Seaborn Visualization Library
- Research papers on weather prediction using machine learning



```
[10]: from sklearn.preprocessing import LabelEncoder
[11]: df.dtypes
[11]: Location
                                 object
         MinTemp
MaxTemp
Rainfall
                                float64
float64
float64
         WindGustDir
WindGustSpeed
                                object
float64
                                 object
         WindDir9am
         WindDir3pm
WindSpeed9am
WindSpeed3pm
                                object
float64
float64
                                float64
float64
float64
         Humidity9am
Humidity3pm
         Pressure9am
         Pressure3pm
Temp9am
Temp3pm
                                float64
float64
float64
         RainToday
RainTomorrow
dtype: object
                                 object
object
         le1=LabelEncoder()
         le2=LabelEncoder()
         le3=LabelEncoder()
         le4=LabelEncoder()
         le5=LabelEncoder()
         le6=LabelEncoder()
         df['Location']=le1.fit_transform(df['Location'])
df['WindGustDir']=le2.fit_transform(df['WindGustDir'])
df['WindDir9am']=le3.fit_transform(df['WindDir9am'])
df['WindDir3pm']=le4.fit_transform(df['WindDir3pm'])
df['RainToday']=le5.fit_transform(df['RainToday'])
         df['RainTomorrow']=le6.fit_transform(df['RainTomorrow'])
[13]: df.head()
            Location MinTemp MaxTemp Rainfall WindGustDir WindGustSpeed WindDir9am WindDir3pm WindSpeed9am WindSpeed3pm Humidity9am Humidity3pm
                    2
                               13.4
                                            22.9
                                                         0.6
                                                                           13
                                                                                              44.0
                                                                                                                 13
                                                                                                                                  14
                                                                                                                                                    20.0
                                                                                                                                                                         24.0
                                                                                                                                                                                            71.0
                                                                                                                                                                                                              22.0
         0
                                                                                              44.0
                                                                                                                                                      4.0
              2 7.4
                                            25.1 0.0
                                                                        14
                                                                                                                                  15
                                                                                                                                                                          22.0
                                                                                                                                                                                            44.0
                                                                                                                                                                                                              25.0
                               12.9
                                             25.7
                                                                                              46.0
                                                                                                                 13
                                                                                                                                  15
                                                         0.0
                                                                           15
                                                                                                                                                     19.0
                                                                                                                                                                          26.0
                                                                                                                                                                                            38.0
                                                                                                                                                                                                              30.0
         3
                    2
                                9.2
                                            28.0
                                                        0.0
                                                                           4
                                                                                              24.0
                                                                                                                  9
                                                                                                                                   0
                                                                                                                                                     11.0
                                                                                                                                                                          9.0
                                                                                                                                                                                            45.0
                                                                                                                                                                                                              16.0
         4
                    2
                               17.5
                                             32.3
                                                         1.0
                                                                           13
                                                                                              41.0
                                                                                                                                   7
                                                                                                                                                      7.0
                                                                                                                                                                          20.0
                                                                                                                                                                                            82.0
                                                                                                                                                                                                              33.0
                                                                                                                  1
          4
                                                                                                                                                                                                                 •
```

```
[61]: accuracy_score(Y_test,y_pred)*100
[61]: 84,09413813648655
[63]: user_input_dict = {
                   "Imputuate = (
"Date': '2023-01-01',
'Location': 'Sydney',
'MinTemp': 15.0,
'MaxTemp': 28.0,
'Rainfall': 2.0,
'Evaporation': 5.0,
'Sunchipe': 9.0
                   'Evaporation': 5.0,
'Sunshine': 8.0,
'WindoustDir': 'N',
'WindoustSpeed': 40.0,
'Windoir9am': 'NE',
'WindDir9am': 'E',
'WindSpeed3am': 15.0,
'Windspeed3am': 20.0,
'Humidity9am': 65.0,
'Pressure9am': 1012.0,
'Pressure3am': 1010.0,
'Cloud9am': 3.0,
                    'Cloud9am': 3.0,
'Cloud3pm': 4.0,
                    'Temp9am': 20.0,
'Temp3pm': 27.0,
                   'RainToday': 'No'
[65]: user_df = pd.DataFrame([user_input_dict])
    user_df=user_df.drop(columns=['Date','Evaporation','Sunshine','Cloud9am','Cloud3pm'])
    user_df['tocation']=le1.fit_transform(user_df['Location'])
    user_df['WindGustDir']=le2.fit_transform(user_df''WindGustDir'])
    user_df['WindDiraym']=le3.fit_transform(user_df''WindDiraym'])
    user_df['WindDiraym']=le4.fit_transform(user_df''WindDiraym'])
    user_df['RainToday']=le5.fit_transform(user_df''RainToday'])
[67]: user_df.head()
[67]: Location MinTemp MaxTemp Rainfall WindGustDir WindGustSpeed WindDir9am WindDir9am WindSpeed9am WindSpeed3pm Humidity9am Humidity9am Humidity9am
           0 0 15.0 28.0 2.0 0 40.0 0 0 15.0
            1
                                                                                                                                                                                                                                                                              D
[69]: ans= lr.predict(user_df)
            if ans==1 :
    print('Yes')
             else:
            print('No')
             No
```

```
[6]: df=df.drop(columns=['Date','Evaporation','Sunshine','Cloud9am','Cloud3pm'])
      df.head()
[6]: laxTemp Rainfall WindGustDir WindGustSpeed WindDir9am WindDir3pm WindSpeed9am WindSpeed3pm Humidity9am Humidity3pm Pressure9am Pressure9pm
                                                44.0
         22.9
                  0.6
                                 W
                                                                W
                                                                           WNW
                                                                                             20.0
                                                                                                             24.0
                                                                                                                            71.0
                                                                                                                                           22.0
                                                                                                                                                       1007.7
                                                                                                                                                                     1007.1
        25.1
                  0.0
                              WNW
                                                44.0
                                                             NNW
                                                                           WSW
                                                                                             4,0
                                                                                                             22.0
                                                                                                                            44.0
                                                                                                                                           25.0
                                                                                                                                                       1010.6
                                                                                                                                                                     1007.8
         25.7
                                                46.0
                                                                                             19.0
                                                                                                             26.0
                                                                                                                            38.0
                                                                                                                                                                     1008.7
     28.0
                  0.0
                                NE
                                                24.0
                                                               SE
                                                                           Е
                                                                                             11.0
                                                                                                              9.0
                                                                                                                            45.0
                                                                                                                                           16.0
                                                                                                                                                       1017.6
                                                                                                                                                                     1012.8
                  1.0
                                 W
                                                41.0
                                                              ENE
                                                                                              7.0
                                                                                                             20.0
                                                                                                                            82.0
                                                                                                                                           33.0
                                                                                                                                                       1010.8
                                                                                                                                                                     1006.0
         32.3
                                                                            NW
       4
                                                                                                                                                                       D
[7]: for i in df.columns:
          # print(i)
if df[i].dtypes=='object':
    df[i] = df[i].fillna(df[i].mode()[0])
          0+[1] = 0+[1].+111na(0+[1].mode()]@
else:
    df[i]=df[i].fillna(df[i].median())
[8]: df.isnull().sum()
[8]: Location
      MinTemp
MaxTemp
Rainfall
      WindGustDir
WindGustSpeed
WindDir9am
      WindDir3pm
WindSpeed9am
WindSpeed3pm
Humidity9am
      Humidity3pm
Pressure9am
Pressure3pm
      Temp9am
Temp3pm
RainToday
      RainTomorrow
dtype: int64
[9]: df.head()
[9]: Location MinTemp MaxTemp Rainfall WindGustDir WindGustSpeed WindDir9am WindDir3pm WindSpeed9am WindSpeed3pm Humidity9am Humidity3pm
                                             0.6
                                                           W
      0 Albury
                        13.4
                                   22.9
                                                                           44.0
                                                                                          W
                                                                                                                       20.0
                                                                                                                                        24.0
                                                                                                                                                       71.0
                                                                                                     WNW
                                                                                                                                                                      22.0
                                                                                                                                                                     25.0
      1 Albury
                         7.4
                                   25.1
                                            0.0
                                                         WNW
                                                                           44.0
                                                                                        NNW
                                                                                                     wsw
                                                                                                                        4.0
                                                                                                                                        22.0
                                                                                                                                                       44.0
                                   25.7
                                                                                                                                                       38.0
           Albury
                        12.9
                                            0.0
                                                         WSW
                                                                           46.0
                                                                                          W
                                                                                                     WSW
                                                                                                                        19.0
                                                                                                                                        26.0
                                                                                                                                                                      30.0
                                   28.0
           Albury
                        9.2
                                            0.0
                                                          NE
                                                                           24.0
                                                                                          SE
                                                                                                                        11.0
                                                                                                                                         9.0
                                                                                                                                                       45.0
                                                                                                                                                                      16.0
```