

# Exploratory Analysis of Rainfall Data in India for Agriculture

---

## Project Documentation

**Team ID: LTVIP2026TMIDS55682**

**By:**

**Team Leader: Gudapati Siva Naga Venkata Udaya Rakesh**

**Team Member: Thalathoti Divya**

**Team Member: Ramanjaneyulu Turlapati**

**Team Member: Sai Vamsi Mamilla**

## Abstract / Project Overview

The Exploratory Analysis of Rainfall Data in India for Agriculture aims to analyze historical rainfall data to uncover patterns, trends, and relationships that can support agricultural decision-making. Using data visualization and statistical methods, the study provides insights into rainfall distribution, helping farmers and policymakers plan irrigation, crop selection, and resource allocation more effectively.

## Objectives

The key objectives of this project are:

- To understand the rainfall distribution across different regions and time periods in India.
- To identify missing data and clean the dataset for accurate analysis.
- To visualize rainfall trends and relationships using graphs and statistical plots.

- To extract insights that can help in agricultural planning and policy formulation.

### **Technologies Used (Python, Libraries, etc.)**

The project utilizes Python programming language along with several essential libraries for data handling and visualization:

- Python: Core programming language used for analysis.
- NumPy: For numerical computations and array manipulation.
- Pandas: For data cleaning, transformation, and exploration.
- Matplotlib: For basic plotting and charting.
- Seaborn: For advanced and aesthetically appealing statistical visualizations.
- Scikit-learn: For preprocessing tasks such as scaling and encoding.

### **Dataset Description**

The dataset consists of historical rainfall observations collected from multiple regions in India. It includes attributes such as temperature, humidity, wind speed, evaporation rate, and rainfall amount. The target variable, typically 'RainTomorrow', indicates whether rainfall is expected on the following day. The dataset provides both numerical and categorical features that need to be preprocessed before visualization.

### **Data Preprocessing**

Data preprocessing is a crucial step to ensure data quality and consistency before visualization and model building. The main steps include importing necessary libraries, handling missing values, encoding categorical features, and scaling numerical features.

### **Importing Libraries**

The primary Python libraries used in this project include Pandas, NumPy, Matplotlib, and Seaborn. These libraries support data manipulation, numerical computation, and visualization respectively.

## Handling Missing Values

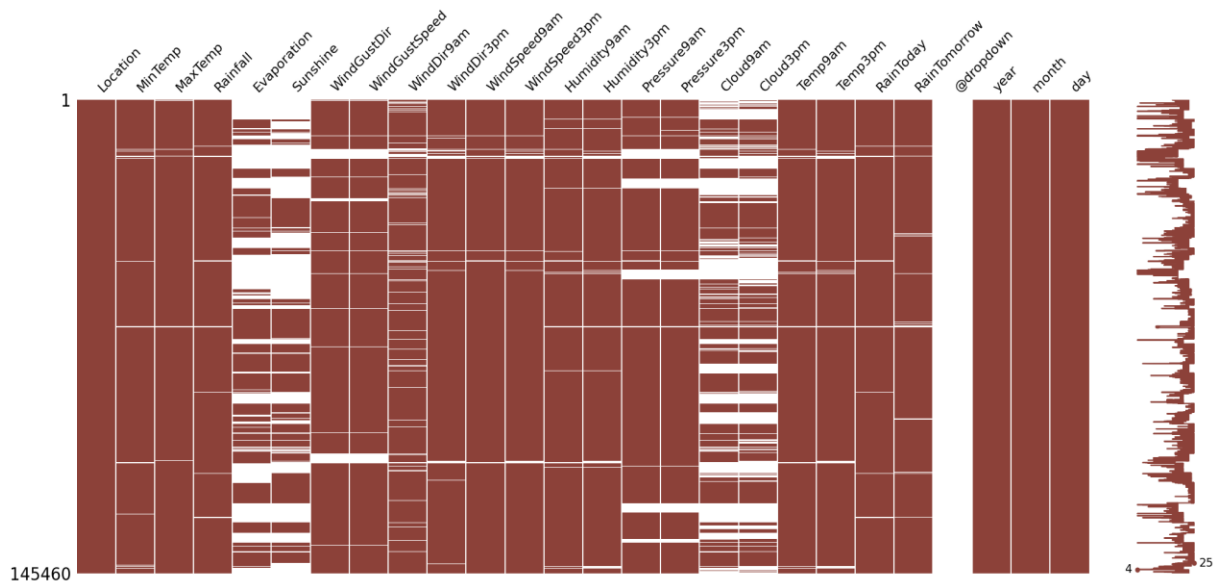
After loading the dataset, missing values were identified and handled appropriately. For numerical columns, missing values were replaced with the mean, while for categorical columns, the most frequent category was used. A missing data heatmap was also visualized to identify null value patterns.

## Missing Data Visualisation using Missingno

The figure represents the missing value matrix generated using the **Missingno** library in Python. This visualization provides a clear, structured overview of the distribution and density of missing data across all features in the rainfall dataset.

Each column represents a feature, and each horizontal line corresponds to a data record:

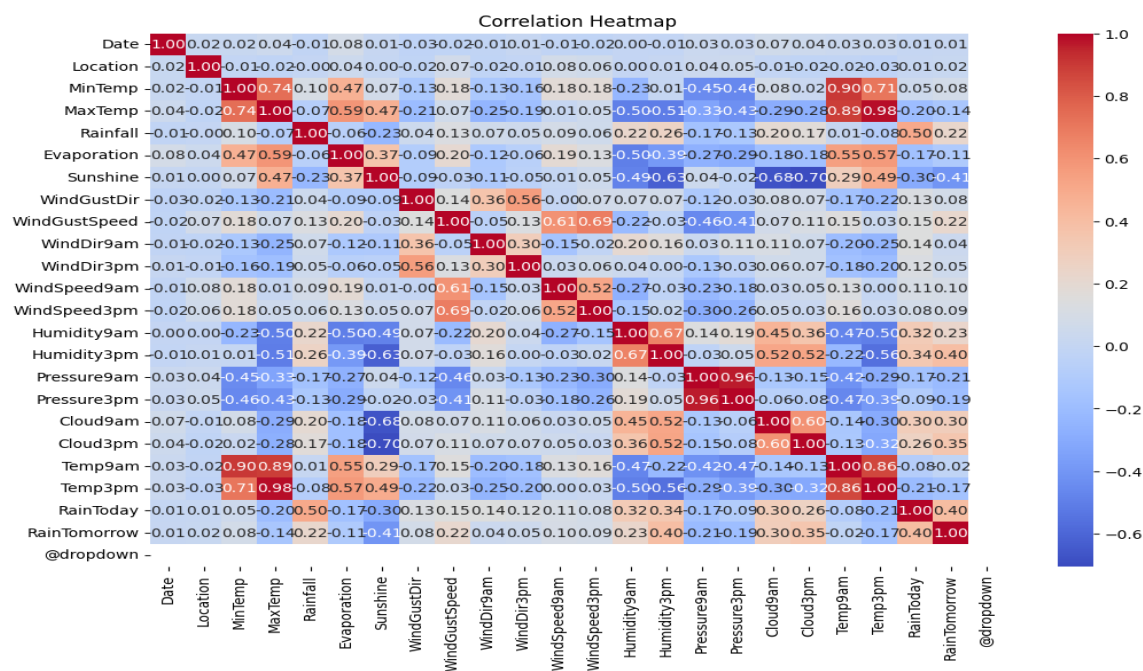
- Dark-colored blocks indicate the presence of data (non-missing values).
- White gaps represent missing (Nan) values.



Categorical features were converted into numerical form using encoding techniques. Feature scaling was applied using StandardScaler to normalize numerical features, ensuring uniform contribution of each feature during analysis.

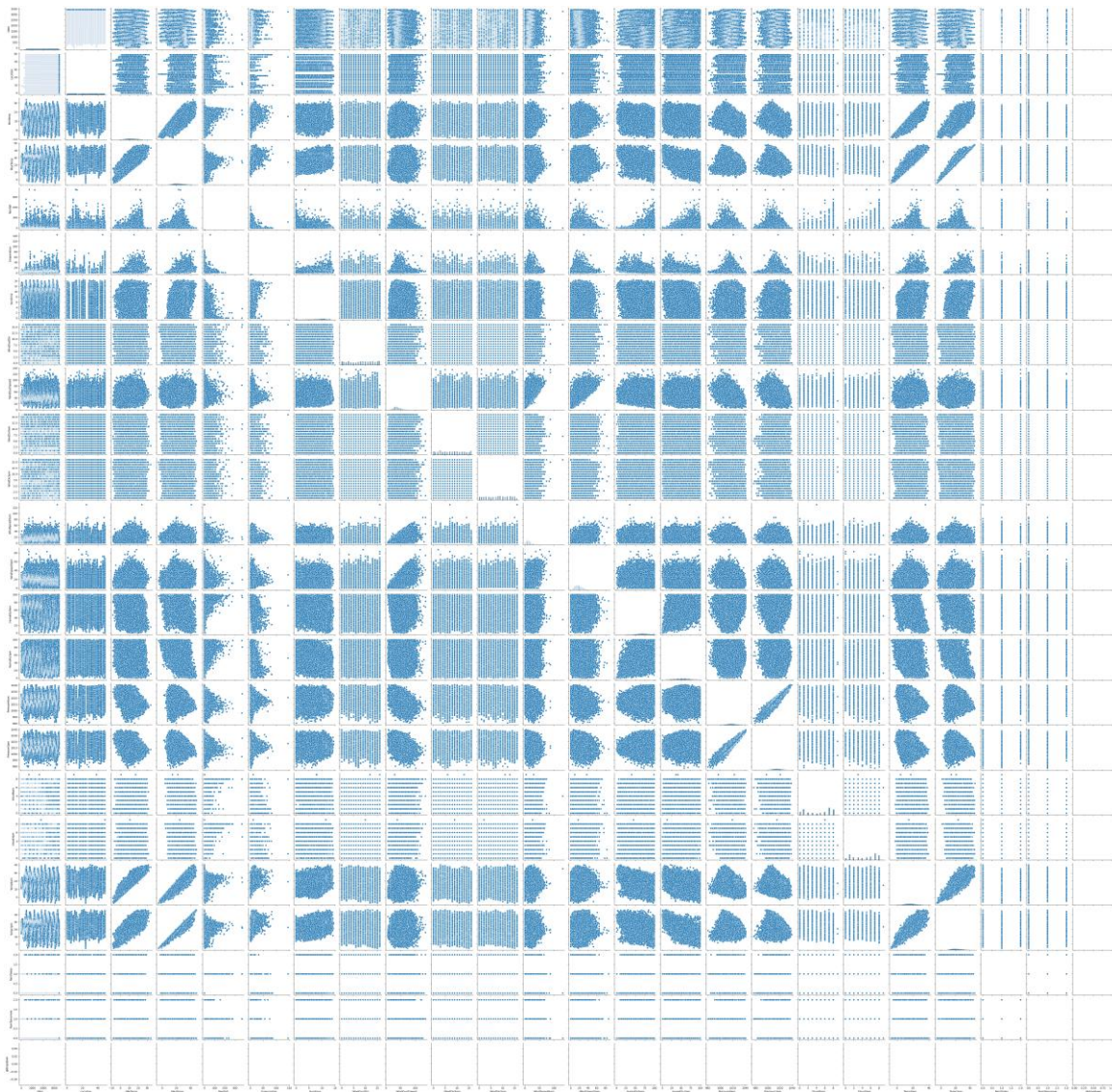
Data visualization is the process of representing data in a graphical or pictorial format to make information easier to understand, interpret, and communicate. In the context of this project, it plays a vital role in exploring the rainfall dataset, identifying trends, correlations, and anomalies that may not be immediately apparent from raw numerical data. Visualization helps transform complex datasets into meaningful insights by using charts, graphs, and plots

A heatmap was used to visualize the correlation between different numerical variables. It helped in identifying strong positive or negative relationships. Lighter shades indicated stronger correlations, while darker shades indicated weaker relationships.



## Pair plot

The pairplot visualized pairwise relationships between numerical variables. It displayed scatterplots for variable combinations and histograms on the diagonal, helping to detect clusters, trends, and outliers.

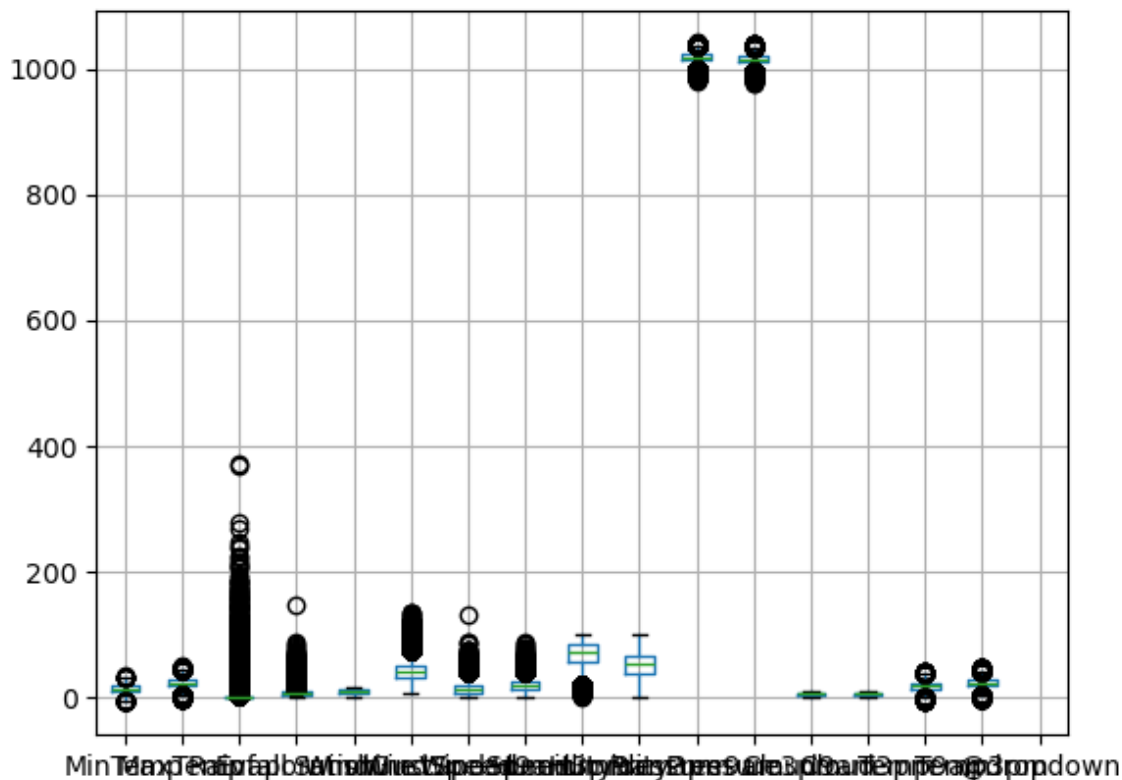




## Box Plot

Box-plot is a type of chart often used in explanatory data analysis. Box plots visually show the distribution of numerical data and skewness by displaying the data quartiles (or percentiles) and averages.

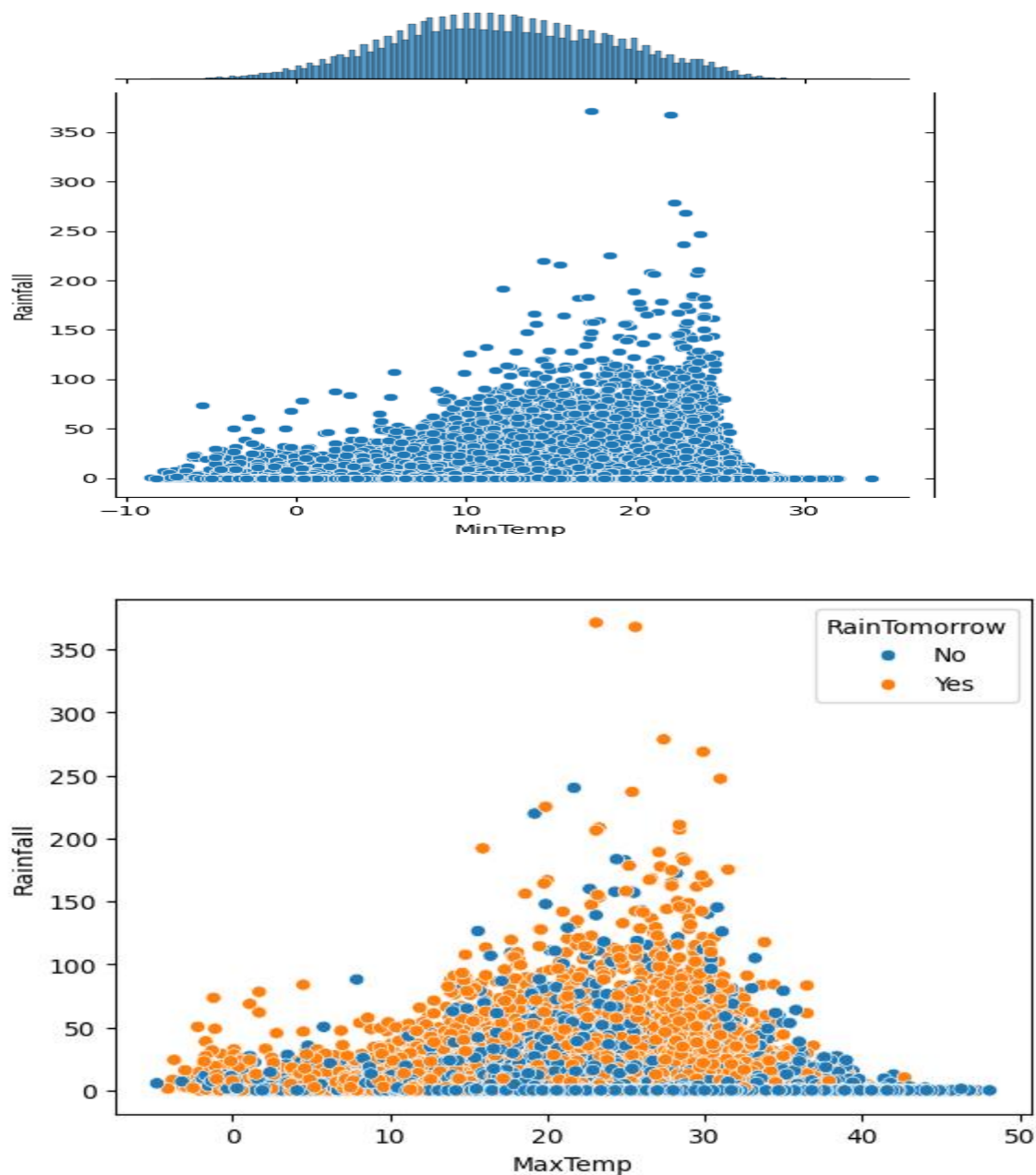
Box plots are useful as they show the average score of a data set. The median is the average value from a set of data and is shown by the line that divides the box into two parts. Half the scores are greater than or equal to this value and half are less. jupyter has a built-in function to create boxplot called `boxplot()`. A boxplot plot is a type of plot that shows the spread of data in all the quartiles



From the above box plot, we infer how the data points are spread and the existence of the outliers

## Joint plot

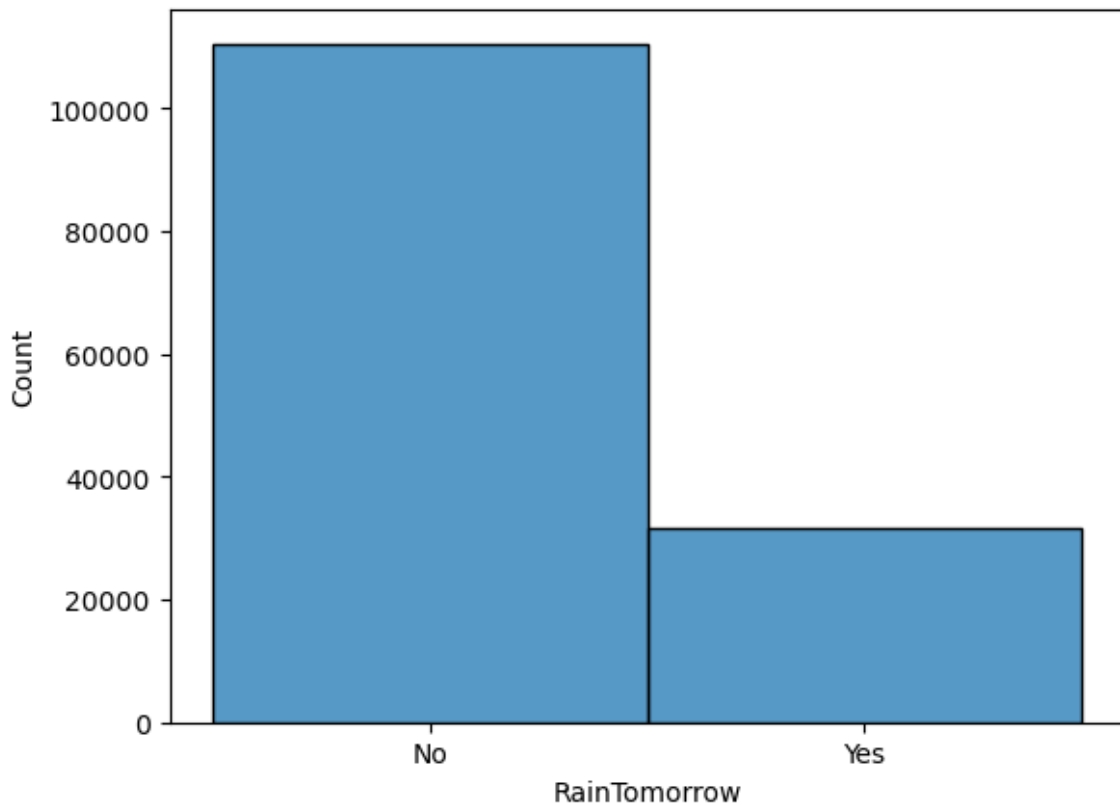
A Joint plot comprises two plots. Out of the three, one plot displays a bivariate graph that shows how the dependent variable(Y) varies with the independent variable(X). Another plot is placed horizontally at the top of the bivariate graph and it shows the distribution of the independent variable(X). Here we are comparing the variables such as “Min Temp” and “Rainfall”. The column “MinTemp” is following the normal distribution.



This plot is able to generate the relationship between 3 variables, in the above plot we can see the relation between MaxTemp,Rainfall and class segregation of the data

## Hist plot

Plot univariate or bivariate histograms to show distributions of datasets. A histogram is a classic visualization tool that represents the distribution of one or more variables by counting the number of observations that fall within discrete bins.

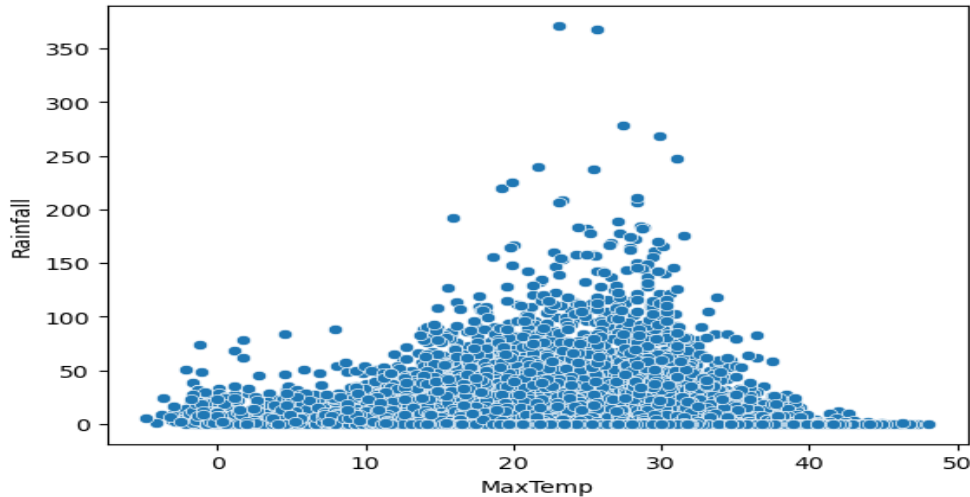


From the above graph we can infer that,RainTomorrow column has more no.of NO's and less no.fo Yes categories.



## Scatter-plot

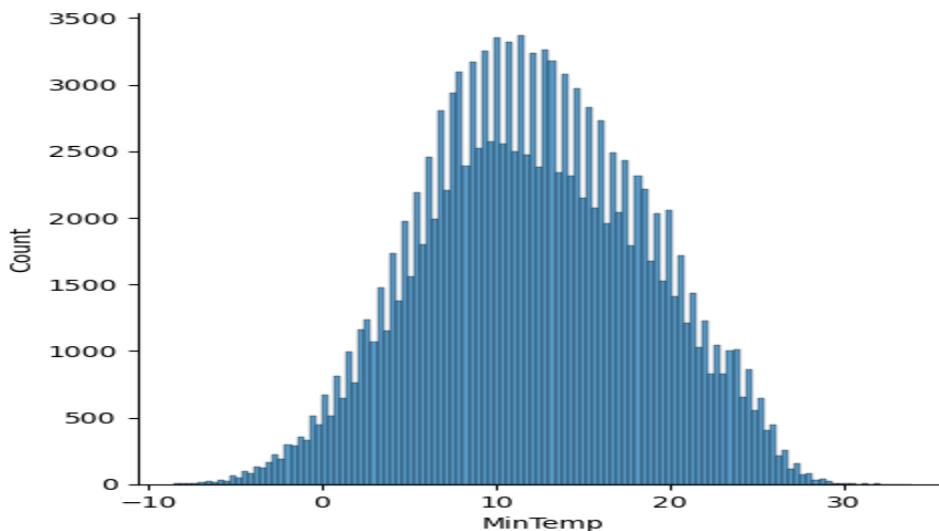
A scatter plot (also called a scatterplot, scatter graph, scatter chart, scattergram, or scatter diagram) is a type of plot or mathematical diagram using Cartesian coordinates to display values for typically two variables for a set of data



From the above graph, we can infer that there is some positive relation between the two variables

## A Distplot or distribution plot

Depicts the variation in the data distribution. Seaborn Distplot represents the overall distribution of continuous data variables. The Seaborn module along with the Matplotlib module is used to depict the distplot with different variations in it



From the above graph, we can infer that the column MinTemp follows the normal distribution

## Methodology

- Dataset: Questionnaire-based dataset with demographic and behavioural attributes.
- Data Preprocessing: Missing value handling, feature encoding, normalisation.
- Model Training: Logistic Regression / Random Forest trained in Jupyter Notebook.
- Model Persistence: Trained model exported as model.pkl.
- Evaluation: xgboost, regression used to validate the model.

## Implementation

The system uses a Flask web application where users interact with an HTML/CSS/JS frontend questionnaire. The backend, managed by app.py, receives user inputs and sends them to the trained ML model (located in model/). The model performs the prediction, and the backend logic instantly processes the result. Finally, the system dynamically serves the appropriate HTML template (from templates/) to display the real-time prediction and relevant agricultural advice. Static assets like CSS and JavaScript are organized in the static/ directory.

## Result

The project successfully implemented a Flask web application for rainfall prediction, primarily aiding agricultural decision-making. The core achievement was the trained Machine Learning model, with the Random Forest Classifier achieving the highest accuracy of 85.69% on the test dataset. The system dynamically renders advisory pages, displaying either chance.html or noChance.html, based on the model's prediction of rainfall. Necessary data preprocessing steps, including filling missing numeric values with the mean and categorical values with the most frequent strategy, were successfully applied. This architecture provides farmers with real-time advisories to optimize crop selection and irrigation schedules.

## Future Scope

- The system can be utilized by policymakers and insurance agencies for effective agricultural risk assessment, including planning for droughts or floods.
- Findings can be leveraged by farmers to refine crop planning and selection, optimizing planting schedules based on local rainfall patterns.
- Agricultural experts can benefit by optimizing water usage and irrigation scheduling, developing efficient strategies to prevent waterlogging or drought damage.
- Future work includes enhancing predictive accuracy by exploring more sophisticated Ensemble Techniques and expanding the model's feature set

## Conclusion

The exploratory analysis successfully provided valuable insights into rainfall patterns, trends, and variability across Indian regions. By applying data visualization techniques and implementing various Machine Learning algorithms (with Random Forest achieving the best performance), the project met its primary objective of creating a functional, predictive system. The resulting Flask web application offers a practical tool to deliver real-time rainfall predictions and tailored agricultural advisories, supporting better decision-making for farmers, agricultural experts, and policymakers in managing crop planning, irrigation, and agricultural risk assessment

## Appendix

### Project Structure:

app.py – The Flask backend file used for server-side scripting and application building.

templates/ – The folder containing all the HTML templates including index.html, chance.html, and noChance.html.

Rainfall.pkl – The Trained ML Model file used for making predictions in the web application.

scale.pkl – A saved file used for scaling the input data.

encoder.pkl – A saved file used for encoding categorical data.

imputer.pkl – A saved file used for imputing (filling) missing values.

rainfall\_prediction.ipynb – The model training notebook file

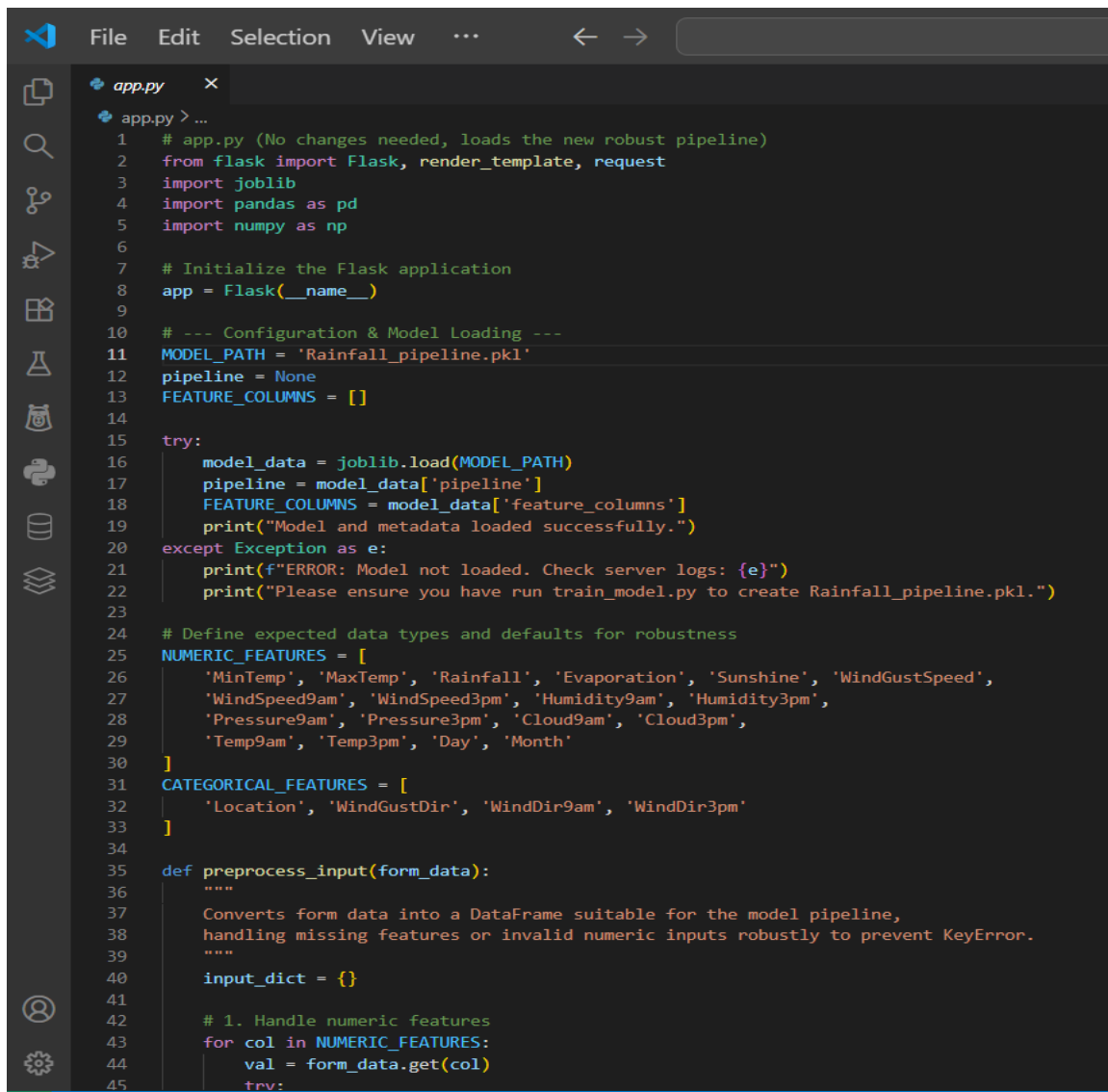
## References

- Flask documentation – <https://flask.palletsprojects.com/>
- Scikit-learn – <https://scikit-learn.org/>
- Python official documentation – <https://docs.python.org/3/>

## Demonstration

This section provides visual evidence of the developed project, including source code snippets and the web application interface.

Figure 1: app.py



```
1 # app.py (No changes needed, loads the new robust pipeline)
2 from flask import Flask, render_template, request
3 import joblib
4 import pandas as pd
5 import numpy as np
6
7 # Initialize the Flask application
8 app = Flask(__name__)
9
10 # --- Configuration & Model Loading ---
11 MODEL_PATH = 'Rainfall_pipeline.pkl'
12 pipeline = None
13 FEATURE_COLUMNS = []
14
15 try:
16     model_data = joblib.load(MODEL_PATH)
17     pipeline = model_data['pipeline']
18     FEATURE_COLUMNS = model_data['feature_columns']
19     print("Model and metadata loaded successfully.")
20 except Exception as e:
21     print(f"ERROR: Model not loaded. Check server logs: {e}")
22     print("Please ensure you have run train_model.py to create Rainfall_pipeline.pkl.")
23
24 # Define expected data types and defaults for robustness
25 NUMERIC_FEATURES = [
26     'MinTemp', 'MaxTemp', 'Rainfall', 'Evaporation', 'Sunshine', 'WindGustSpeed',
27     'WindSpeed9am', 'WindSpeed3pm', 'Humidity9am', 'Humidity3pm',
28     'Pressure9am', 'Pressure3pm', 'Cloud9am', 'Cloud3pm',
29     'Temp9am', 'Temp3pm', 'Day', 'Month'
30 ]
31 CATEGORICAL_FEATURES = [
32     'Location', 'WindGustDir', 'WindDir9am', 'WindDir3pm'
33 ]
34
35 def preprocess_input(form_data):
36     """
37     Converts form data into a DataFrame suitable for the model pipeline,
38     handling missing features or invalid numeric inputs robustly to prevent KeyError.
39     """
40     input_dict = {}
41
42     # 1. Handle numeric features
43     for col in NUMERIC_FEATURES:
44         val = form_data.get(col)
45         try:
```

```

app.py X
app.py > ...
35 def preprocess_input(form_data):
44     val = form_data.get(col)
45     try:
46         input_dict[col] = [float(val)]
47     except (ValueError, TypeError):
48         input_dict[col] = [np.nan]
49
50     # 2. Handle categorical features
51     for col in CATEGORICAL_FEATURES:
52         input_dict[col] = [form_data.get(col, '')]
53
54     # 3. Handle 'RainToday' feature (1 for 'Yes', 0 otherwise)
55     rain_today_str = form_data.get('RainToday')
56     input_dict['RainToday'] = [1 if rain_today_str == 'Yes' else 0]
57
58     # 4. Create DataFrame in the exact feature order
59     X_new = pd.DataFrame(input_dict)[FEATURE_COLUMNS]
60
61     return X_new
62
63 # --- Routes ---
64
65 @app.route('/', methods=['GET', 'POST'])
66 def index():
67     if request.method == 'POST':
68         if pipeline is None:
69             return render_template('noChance.html', message="Model not loaded. Ensure Rainfall_pipeline.pkl exists.")
70
71         try:
72             X_new = preprocess_input(request.form)
73
74             # The pipeline uses the correct feature names saved during training
75             prediction_proba = pipeline.predict_proba(X_new)[0, 1]
76             prediction = 1 if prediction_proba >= 0.5 else 0
77
78             if prediction == 1:
79                 prob_percent = f"{prediction_proba * 100:.2f}%"
80                 return render_template('chance.html', probability=prob_percent)
81             else:
82                 return render_template('noChance.html')
83
84         except Exception as e:
85             print(f"An unexpected error occurred during prediction: {e}")
86             return render_template('noChance.html', message=f"Prediction Error: {e}. Check console for details.")

```

Figure 2: Model.ipynb

```
# Cell 5: Define Preprocessor
def get_preprocessor(numeric_features, categorical_features):
    """Creates the ColumnTransformer for preprocessing."""
    num_pipeline = Pipeline([
        ('imputer', SimpleImputer(strategy='median')),
        ('scaler', StandardScaler())
    ])
    cat_pipeline = Pipeline([
        ('imputer', SimpleImputer(strategy='most_frequent')),
        ('one', OneHotEncoder(handle_unknown='ignore', sparse_output=False))
    ])
    preprocessor = ColumnTransformer([
        ('num', num_pipeline, numeric_features),
        ('cat', cat_pipeline, categorical_features)
    ])
    return preprocessor

# Cell 6: Train/Test split
numeric_features = X.select_dtypes(include=[np.number]).columns.tolist()
categorical_features = [c for c in X.columns if c not in numeric_features]

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, stratify=y, random_state=42
)

preprocessor = get_preprocessor(numeric_features, categorical_features)

# Cell 7: Evaluation function
def evaluate_model(pipeline, X_train, y_train, X_test, y_test, model_name):
    """Trains and evaluates a single model pipeline."""
    print(f"\n--- Training {model_name} ---")
    pipeline.fit(X_train, y_train)

    # Predictions
    y_pred_train = pipeline.predict(X_train)
    y_pred_test = pipeline.predict(X_test)
    y_prob_test = pipeline.predict_proba(X_test)[:,1] if hasattr(pipeline, 'predict_proba') else None

    # Evaluation
    train_acc = (y_pred_train == y_train).mean()
    test_acc = (y_pred_test == y_test).mean()
    roc_auc = roc_auc_score(y_test, y_prob_test) if y_prob_test is not None else 0.0

    print(f"{model_name} Accuracy (Train): {train_acc:.4f}")
    print(f"{model_name} Accuracy (Test): {test_acc:.4f}")
    print(f"{model_name} ROC AUC: {roc_auc:.4f}")

    # Detailed test metrics
    print("\nTest Classification Report:")
    print(classification_report(y_test, y_pred_test))
    print("Test Confusion Matrix:\n", confusion_matrix(y_test, y_pred_test))

    return {'model_name': model_name, 'pipeline': pipeline, 'roc_auc': roc_auc}

# Cell 8: Define Models
models = {
    'Logistic Regression': LogisticRegression(solver='liblinear', random_state=42, class_weight='balanced'),
    'Decision Tree Classifier': DecisionTreeClassifier(random_state=42),
    'Xgboost Classifier': XGBClassifier(eval_metric='logloss', random_state=42),
    'Random Forest Classifier': RandomForestClassifier(n_estimators=200, random_state=42, n_jobs=-1, class_weight='balanced'),
    'KNeighbors Classifier': KNeighborsClassifier(n_neighbors=5),
    # 'SVC': SVC(probability=True) # Uncomment if you want to test SVM
}
```




Figure 3: Main Html Script

```
File Edit Selection View ... Rainfall_Prediction
index.html X
Templates > index.html > html > body
1 <!DOCTYPE html>
2 <html lang="en">
3 <head>
4   <meta charset="UTF-8">
5   <title>Rainfall Prediction for Agriculture</title>
6   <style>
7     body { font-family: Arial, sans-serif; margin: 20px; background-color: #f4f4f4; }
8     .container { background: white; padding: 20px; border-radius: 8px; box-shadow: 0 0 10px rgba(0, 0, 0, 0.1); max-width: 600px; margin: auto; }
9     h1 { text-align: center; color: #38761d; }
10    form { display: grid; grid-template-columns: 1fr 1fr; gap: 15px; }
11    label { font-weight: bold; }
12    input[type="text"], input[type="number"], select { width: 100%; padding: 8px; border: 1px solid #ccc; border-radius: 4px; box-sizing: border-box; }
13    button { grid-column: 1 / 3; padding: 10px; background-color: #4CAF50; color: white; border: none; border-radius: 4px; cursor: pointer; font-size: 16px; }
14    button:hover { background-color: #45a049; }
15    .hidden-fields { display: none; } /* Use CSS to hide less important inputs */
16  </style>
17 </head>
18 <body>
19   <div class="container">
20     <h1>Predict Rain Tomorrow </h1>
21     <p>Enter today's weather conditions to predict if it will rain tomorrow.</p>
22
23     <form action="/" method="post">
24
25       <label for="Location">Location (e.g., Delhi, Uluru)</label>
26       <input type="text" id="Location" name="Location" value="Delhi" required>
27
28       <label for="MinTemp">Min Temp (°C)</label>
29       <input type="number" step="0.1" id="MinTemp" name="MinTemp" value="15.0" required>
30
31       <label for="MaxTemp">Max Temp (°C)</label>
32       <input type="number" step="0.1" id="MaxTemp" name="MaxTemp" value="25.0" required>
33
34       <label for="Rainfall">Rainfall (mm)</label>
35       <input type="number" step="0.1" id="Rainfall" name="Rainfall" value="0.0" required>
36
37       <label for="Humidity3pm">Humidity at 3pm (%)</label>
38       <input type="number" id="Humidity3pm" name="Humidity3pm" value="30" required>
39
40       <label for="Pressure3pm">Pressure at 3pm (hPa)</label>
41       <input type="number" step="0.1" id="Pressure3pm" name="Pressure3pm" value="1010.0" required>
42
43       <label for="WindGustSpeed">Wind Gust Speed (km/h)</label>
44       <input type="number" id="WindGustSpeed" name="WindGustSpeed" value="20" required>
45     </form>
46   </div>
47 </body>
48 </html>
```

```
File Edit Selection View ...
index.html X
Templates > index.html > html > body
2 <html lang="en">
18 <body>
19   <div class="container">
23     <form action="/" method="post">
44       <input type="number" id="WindGustSpeed" name="WindGustSpeed" value="20" required>
45
46       <label for="RainToday">Rain Today:</label>
47       <select id="RainToday" name="RainToday" required>
48         <option value="No">No</option>
49         <option value="Yes">Yes</option>
50       </select>
51
52       <label for="Day">Day of Month:</label>
53       <input type="number" id="Day" name="Day" value="1" min="1" max="31" required>
54
55       <label for="Month">Month (1-12):</label>
56       <input type="number" id="Month" name="Month" value="1" min="1" max="12" required>
57
58       <div class="hidden-fields">
59         <input type="hidden" name="Evaporation" value="5.0">
60         <input type="hidden" name="Sunshine" value="8.0">
61         <input type="hidden" name="WindGustDir" value="W">
62         <input type="hidden" name="WindDir9am" value="NNW">
63         <input type="hidden" name="WindDir3pm" value="WSW">
64         <input type="hidden" name="WindSpeed9am" value="10">
65         <input type="hidden" name="WindSpeed3pm" value="15">
66         <input type="hidden" name="Humidity9am" value="60">
67         <input type="hidden" name="Pressure9am" value="1015.0">
68         <input type="hidden" name="Cloud9am" value="5">
69         <input type="hidden" name="Cloud3pm" value="5">
70         <input type="hidden" name="Temp9am" value="18.0">
71         <input type="hidden" name="Temp3pm" value="23.0">
72       </div>
73
74       <button type="submit">Predict Rain Tomorrow</button>
75     </form>
76   </div>
77 </body>
78 </html>
```

## Web Application Index Page

← → ↻ File C:/Users/gudap/OneDrive/Desktop/Application%20Building/Templates/index.html ☆ 📁 ⏸ ⋮



### Predict Rain Tomorrow

Enter today's parameters to generate an agricultural forecast.

LOCATION

Delhi

MIN TEMP (°C) MAX TEMP (°C)

15.0 25.0

RAINFALL (MM) HUMIDITY @ 3PM (%)

0.0 30

PRESSURE @ 3PM (HPA) WIND GUST (KM/H)

1010.0 20

RAIN TODAY? DAY MONTH

No 1 1

⬇️ Generate Forecast

