

# Exploratory Analysis of Rainfall Data in India for Agriculture

---

## Project Documentation

**Team ID:** LTVIP2026TMIDS73971

**Team Leader:** Budigi Narasimhulu

**Team Member:** Kala Sai Pradeep

**Team Member:** Chirrareddy Tejeswar

**Team Member:** Jangam Chandu Reddy

**Team Member:** Gundabathina Bensun

## 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

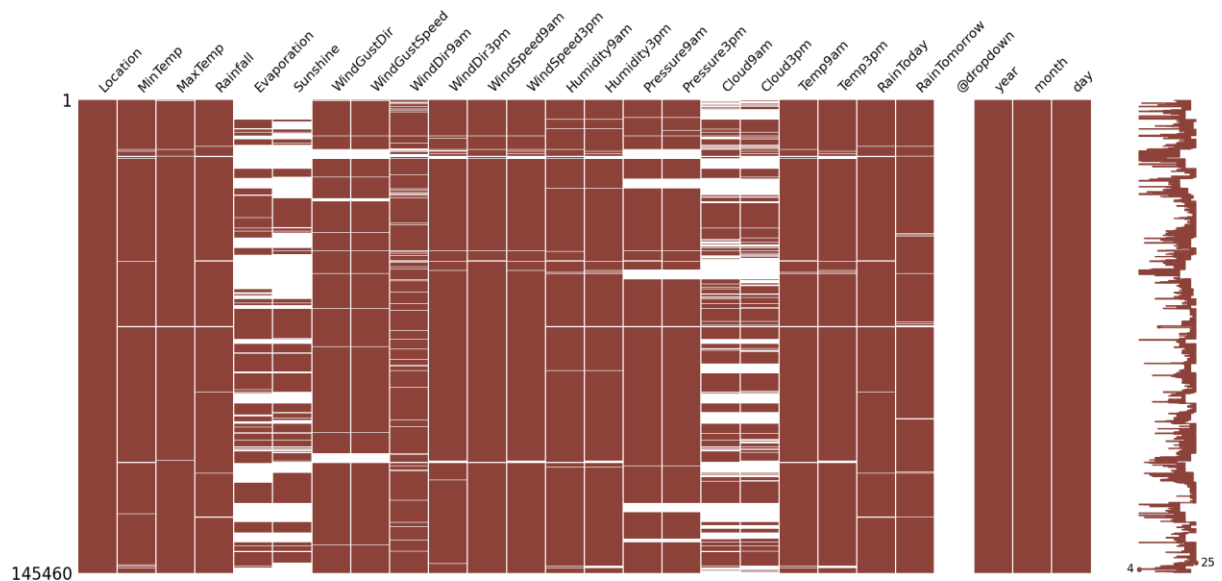
After loading the dataset, missing values were identified and handled appropriately. For numerical columns, missing values were replaced with the mean, while for categorical columns, the most frequent category was used. A missing data heatmap was also visualized to identify null value patterns.

## 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.



## Encoding and Scaling

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.

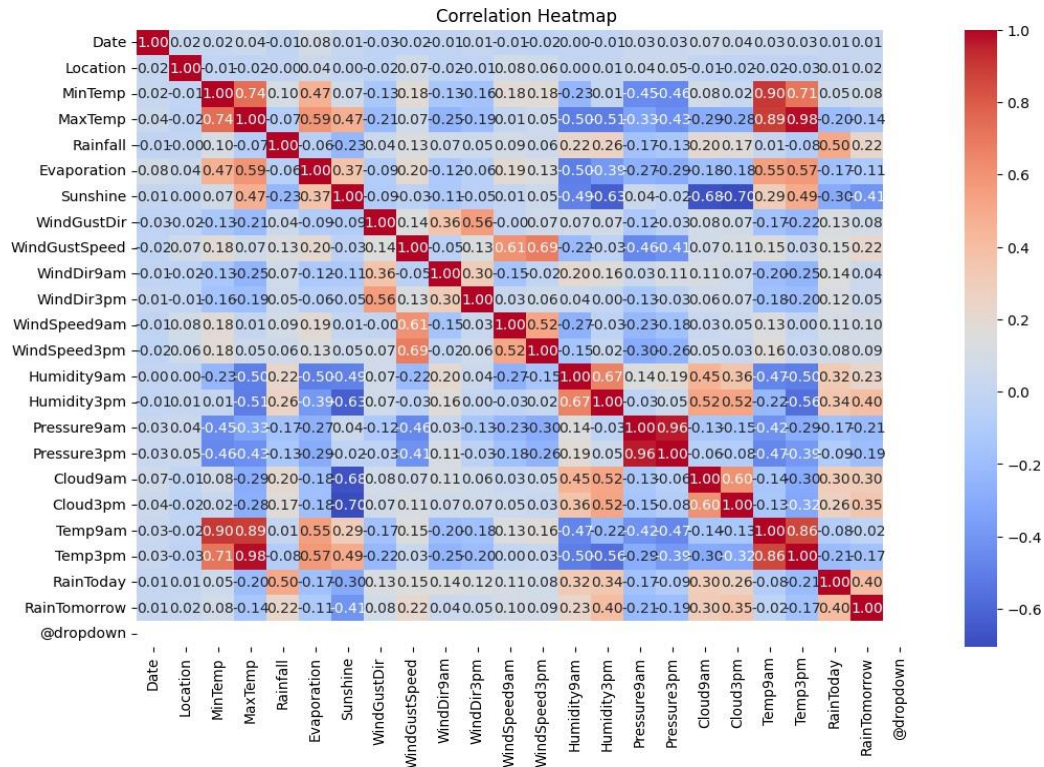
## Data Visualization

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

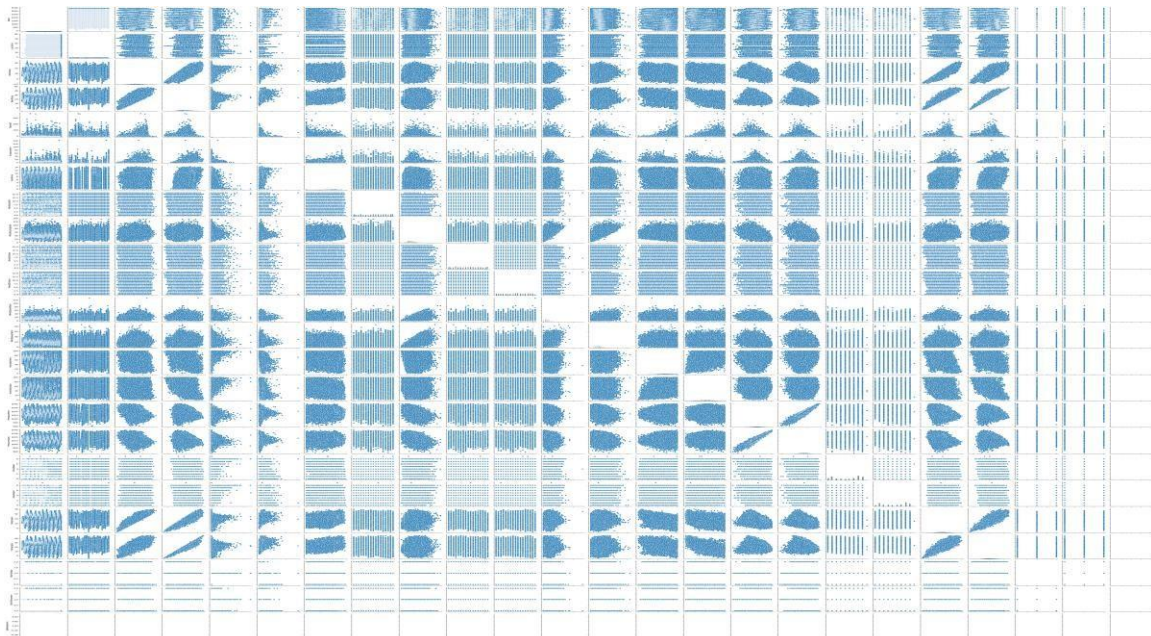
## Correlation Heatmap

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.



## 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.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.

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

rainfall\_prediction.ipynb – The model training notebook file

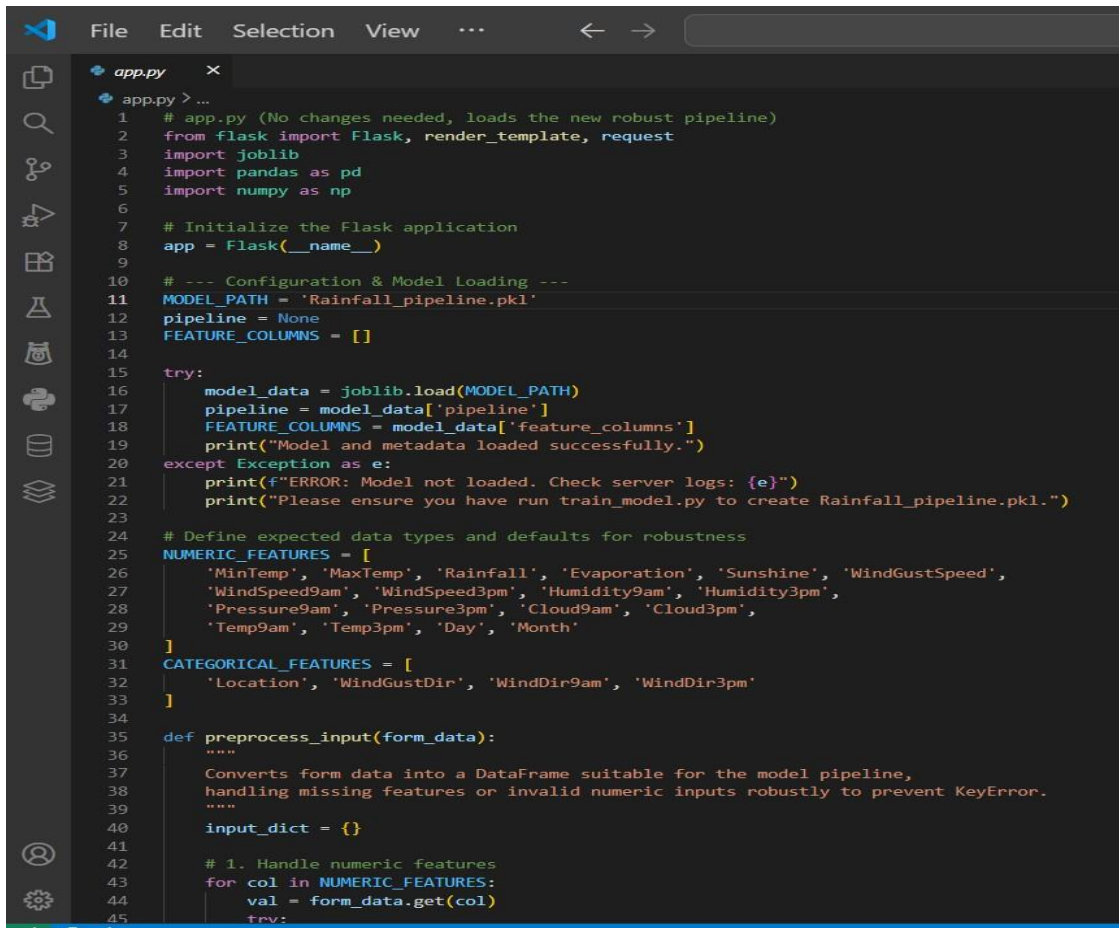
## References

- 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.

### 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
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.")
```

## Model.ipynb

```
File Edit Selection View ...
train_model.ipynb
Cell 1: Import libraries
Generate Code Markdown Run All Clear All Outputs Outline ...
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

[9]

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)

[10]

Cell 7: Evaluation function
def evaluate_model(pipeline, X_train, y_train, X_test, y_test, model_name):
    """Train 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)
```

```
File Edit Selection View ... ← → Rainfall_Prediction

train_model.ipynb x
train_model.ipynb > # Cell 1: Import libraries
Generate + Code + Markdown | Run All | Clear All Outputs | Outline ...

# 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}

[11]

# 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
}

[12]
```

## Main Html Script

```
<!DOCTYPE html>
<html lang="en">
<head>
  <meta charset="UTF-8">
  <title>AgroRain - Smart Rain Prediction</title>
  <meta name="viewport" content="width=device-width, initial-scale=1.0">
  <script src="https://cdn.tailwindcss.com"></script>
</head>
<body class="bg-gradient-to-br from-blue-900 via-indigo-800 to-blue-600 min-h-screen flex items-center justify-center p-6">
  <div class="bg-white/10 backdrop-blur-lg shadow-2xl rounded-3xl p-10 w-full max-w-4xl border border-white/20">
    <h1 class="text-4xl font-extrabold text-white text-center mb-2">
      AgroRain Intelligence
    </h1>
    <p class="text-center text-blue-100 mb-8">
      AI-Powered Rainfall Prediction for Smart Agriculture
    </p>
    <form action="/" method="post" class="grid md:grid-cols-2 gap-6">
      <div>
        <label class="block text-white font-medium mb-1">Location</label>
        <input type="text" name="Location" required
          class="w-full p-3 rounded-xl bg-white/20 text-white placeholder-white/70 border border-white/30 focus:outline-none focus:ring-2 focus:ring-blue-300">
      </div>
      <div>
        <label class="block text-white font-medium mb-1">Min Temperature (°C)</label>
        <input type="number" step="0.1" name="MinTemp" required
          class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
      </div>
      <div>
        <label class="block text-white font-medium mb-1">Max Temperature (°C)</label>
        <input type="number" step="0.1" name="MaxTemp" required
          class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
      </div>
      <div>
        <label class="block text-white font-medium mb-1">Rainfall (mm)</label>
        <input type="number" step="0.1" name="Rainfall" required
          class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
      </div>
    </form>
  </div>
```

```

<div>
  <label class="block text-white font-medium mb-1">Humidity (3PM %)</label>
  <input type="number" name="Humidity3pm" required
    class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
</div>
<div>
  <label class="block text-white font-medium mb-1">Pressure (hPa)</label>
  <input type="number" step="0.1" name="Pressure3pm" required
    class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
</div>
<div>
  <label class="block text-white font-medium mb-1">Wind Gust Speed (km/h)</label>
  <input type="number" name="WindGustSpeed" required
    class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
</div>
<div>
  <label class="block text-white font-medium mb-1">Rain Today</label>
  <select name="RainToday"
    class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
    <option value="No">No</option>
    <option value="Yes">Yes</option>
  </select>
</div>
<div>
  <label class="block text-white font-medium mb-1">Day</label>
  <input type="number" name="Day" min="1" max="31" required
    class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
</div>
<div>
  <label class="block text-white font-medium mb-1">Month</label>
  <input type="number" name="Month" min="1" max="12" required
    class="w-full p-3 rounded-xl bg-white/20 text-white border border-white/30 focus:ring-2 focus:ring-blue-300">
</div>
</div>
<div class="md:col-span-2 mt-6">
  <button type="submit"
    class="w-full bg-white text-blue-900 font-bold py-4 rounded-xl hover:scale-105 transition duration-300 shadow-xl">
    Predict Rainfall
  </button>
</div>
</form>
</div>
</body>
</html>

```

**Output: -**

 Rainfall Expected Tomorrow

 No Rainfall Expected Tomorrow

GitHub Link: <https://github.com/Narasimhulubudigi27/exploratory-analysis-of-rain-fall-data-in-india-for-agriculture>