

# Weather Prediction Model and MLOps

## Group 1 Team Members

**Aniket Shukla**

**Abhishek Yadav**

**Sagar Kamble**

**Sandeep Mamoriya**

**Sherine Martina**

**Shravya Pendyala**

**Sridhar Mulumoodi**

A dramatic landscape photograph showing a vast field of tall grass in the foreground. The sky is filled with dark, heavy clouds, with a bright, jagged lightning bolt striking down on the right side. The horizon is visible in the distance, and the overall scene is lit with a mix of warm, golden light from the left and cooler, blue light from the right.

# Problem Description

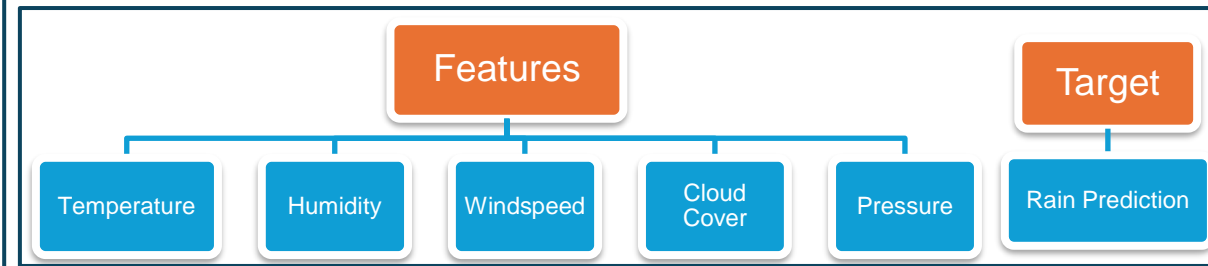
- The objective of this project is to develop a ML model that can forecast rainfall using a variety of input variables, such as temperature, humidity, wind speed, atmospheric pressure & cloud cover.
- Classification Problem

# Dataset info & Cleaning

Dataset consists of 2,500 weather observations, it's a simple yet practical dataset for learning how to **predict rainfall** based on various weather conditions.

Ref: [Weather Forecast Dataset \(kaggle.com\)](https://www.kaggle.com/datasets/ashishpatel26/weather-forecast-dataset)

A	B	C	D	E	F
Temperature	Humidity	Wind_Speed	Cloud_Cover	Pressure	Rain
23.72034	89.59264	7.33560	50.50169	1032.37876	rain
27.87973	46.48970	5.95248	4.99005	992.61419	no rain
25.06908	83.07284	1.37199	14.85578	1007.23162	no rain
23.62208	74.36776	7.05055	67.25528	982.63201	rain
20.59137	96.85882	4.64392	47.67644	980.82514	no rain
26.14735	48.21726	15.25855	59.76628	1049.73875	no rain
20.93968	40.79944	2.23257	45.82751	1014.17377	no rain
32.29433	51.84847	2.87362	92.55150	1006.04173	no rain
34.09157	48.05711	5.57021	82.52487	993.73205	no rain
19.58604	82.97829	5.76054	98.01445	1036.50346	rain
29.79313	81.31765	16.92610	93.92329	1029.40269	no rain



## Data importing

```
url = "https://raw.githubusercontent.com/shravyapendyala/CCE_Assignment_1/refs/heads/main/weather_forecast_data.csv"
dataset=pd.read_csv(url)
data=pd.read_csv(url)
dataset.head(10)
```



	Temperature	Humidity	Wind_Speed	Cloud_Cover	Pressure	Rain
0	23.720338	89.592641	7.335604	50.501694	1032.378759	rain
1	27.879734	46.489704	5.952484	4.990053	992.614190	no rain
2	25.069084	83.072843	1.371992	14.855784	1007.231620	no rain
3	23.622080	74.367758	7.050551	67.255282	982.632013	rain
4	20.591370	96.858822	4.643921	47.676444	980.825142	no rain
5	26.147353	48.217260	15.258547	59.766279	1049.738751	no rain
6	20.939680	40.799444	2.232566	45.827508	1014.173766	no rain
7	32.294325	51.848471	2.873621	92.551497	1006.041733	no rain
8	34.091569	48.057114	5.570206	82.524873	993.732047	no rain
9	19.586038	82.978293	5.760537	98.014450	1036.503457	rain

## Data pre-processing

Using `LabelEncoder()` from `sklearn`, converted the categorical data into numerical data so that it can be used for fitting the model

	Temperature	Humidity	Wind_Speed	Cloud_Cover	Pressure	Rain
0	23.720338	89.592641	7.335604	50.501694	1032.378759	rain
1	27.879734	46.489704	5.952484	4.990053	992.614190	no rain
2	25.069084	83.072843	1.371992	14.855784	1007.231620	no rain
3	23.622080	74.367758	7.050551	67.255282	982.632013	rain
4	20.591370	96.858822	4.643921	47.676444	980.825142	no rain

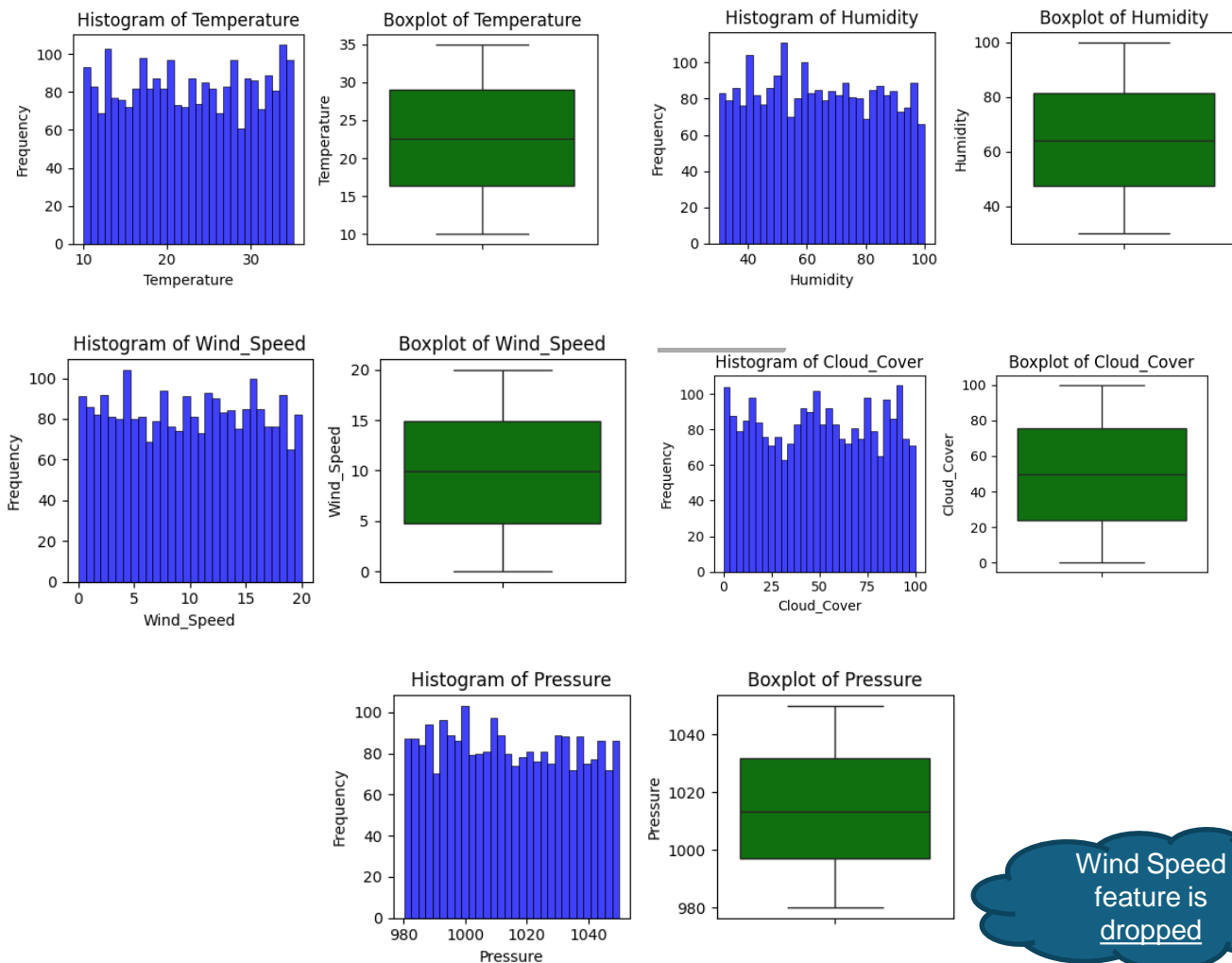


	Temperature	Humidity	Wind_Speed	Cloud_Cover	Pressure	Rain
0	23.720338	89.592641	7.335604	50.501694	1032.378759	1
1	27.879734	46.489704	5.952484	4.990053	992.614190	0
2	25.069084	83.072843	1.371992	14.855784	1007.231620	0
3	23.622080	74.367758	7.050551	67.255282	982.632013	1
4	20.591370	96.858822	4.643921	47.676444	980.825142	0



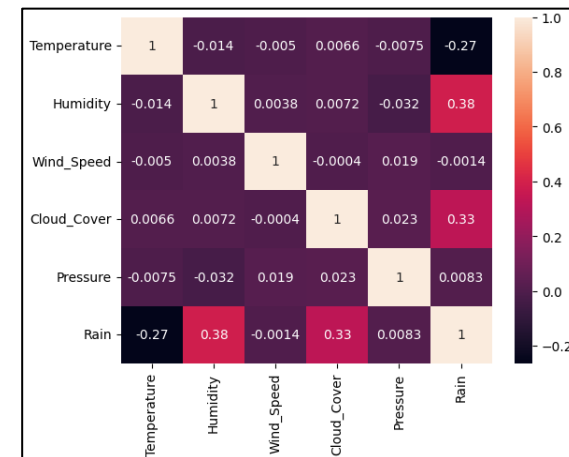
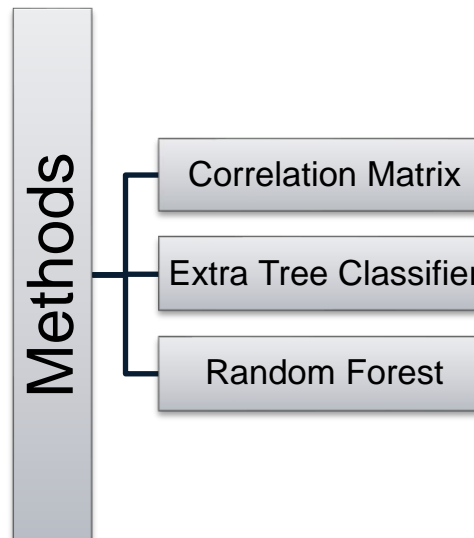
# ML Methodologies

## Histogram & Boxplot of Features Distribution

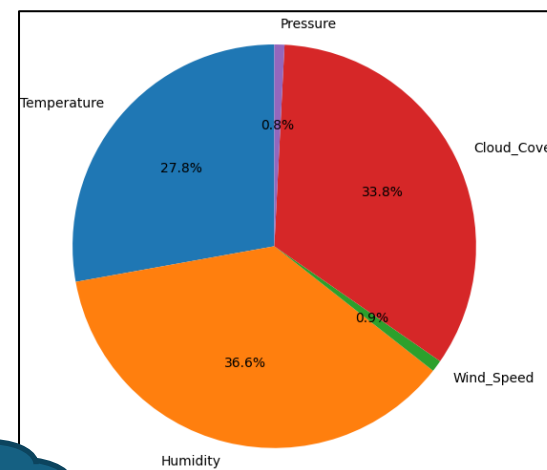


## Feature Selection

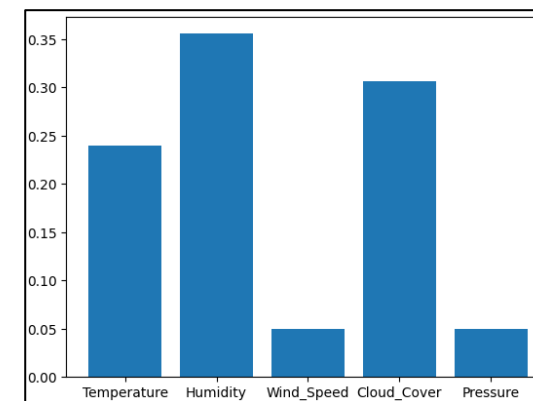
To improve model performance, reduce over-fitting, and simplify the model



Correlation matrix



Random Forest

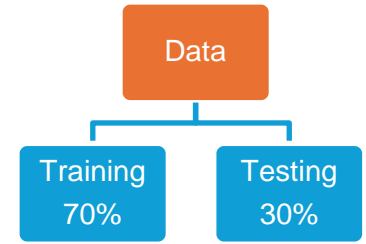


Extra Tree Classifier

Wind Speed feature is dropped

# ML Methodologies

## Training and Testing the model

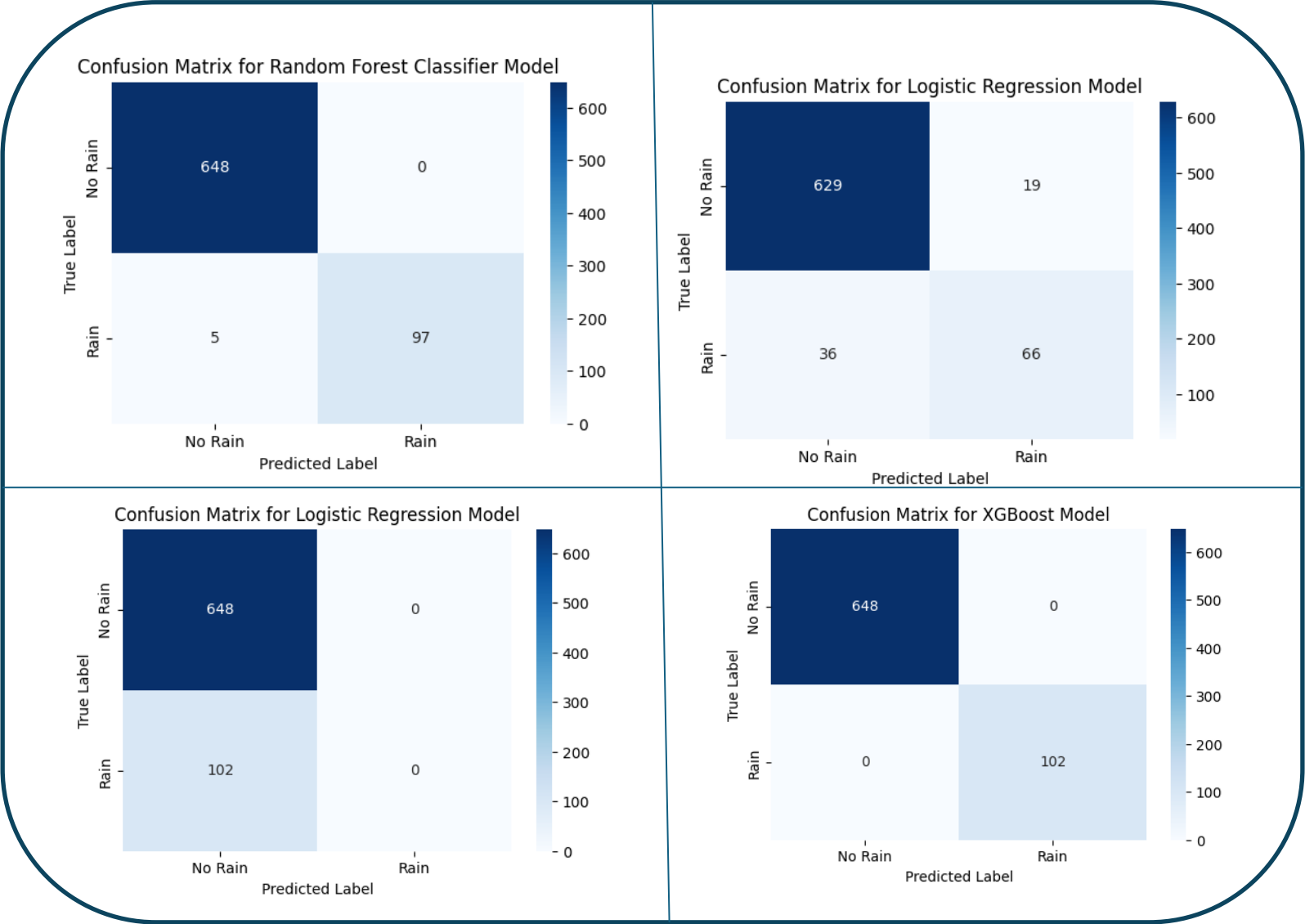


	Random Forest Classifier	Logistic Regression	SVM	XGBoost
Pros	<ul style="list-style-type: none"><li>• High accuracy</li><li>• Versatile in dealing with complex datasets containing outliers</li><li>• Effective classification and regression</li></ul>	<ul style="list-style-type: none"><li>• Good accuracy for simply datasets</li><li>• Easy to implement and interpret</li><li>• Interprets model coefficients as indicators of feature importance</li></ul>	<ul style="list-style-type: none"><li>• Models high-dimensional data</li><li>• Possess good generalization performance (classifies new and unseen data)</li></ul>	<ul style="list-style-type: none"><li>• High accuracy and high precision</li><li>• Regularization techniques avoids over-fitting</li></ul>
Cons	<ul style="list-style-type: none"><li>• Handling huge datasets can be time consuming</li><li>• Not preferred when the model must be highly interpretable.</li></ul>	<ul style="list-style-type: none"><li>• Non-linear problems can't be dealt with this model</li><li>• Not preferred when the number of features is more than the observations</li></ul>	<ul style="list-style-type: none"><li>• Handling huge datasets is not possible</li><li>• Limited to two-class problems (multi-class is dealt using other strategies)</li></ul>	<ul style="list-style-type: none"><li>• Tuning the hyper-parameters of this algorithm can be time consuming</li><li>• “Black Box” algorithm as it is difficult to interpret and understand the predictions.</li></ul>

# Results

## Accuracy and Confusion Matrix

Training Model	Accuracy
Random Forest Classifier	99.3
Logistic Regression	93.2
SVM	86.4
XGBoost	100



# MLOps

Train and store model in MLFlow and later fetch it to deploy on end user side

mlflow2.18.0

Experiments

Models

Search Experiments

☐ Default

☒ MLflow WeatherForecast-1

MLflow WeatherForecast-1

Provide Feedback

Add Description

Share

Runs

Evaluation

Experimental

Traces

Experimental

metrics.rmse < 1 and params.model = "tree"

Time created

State: Active

Datasets

+ New run

Sort: Created

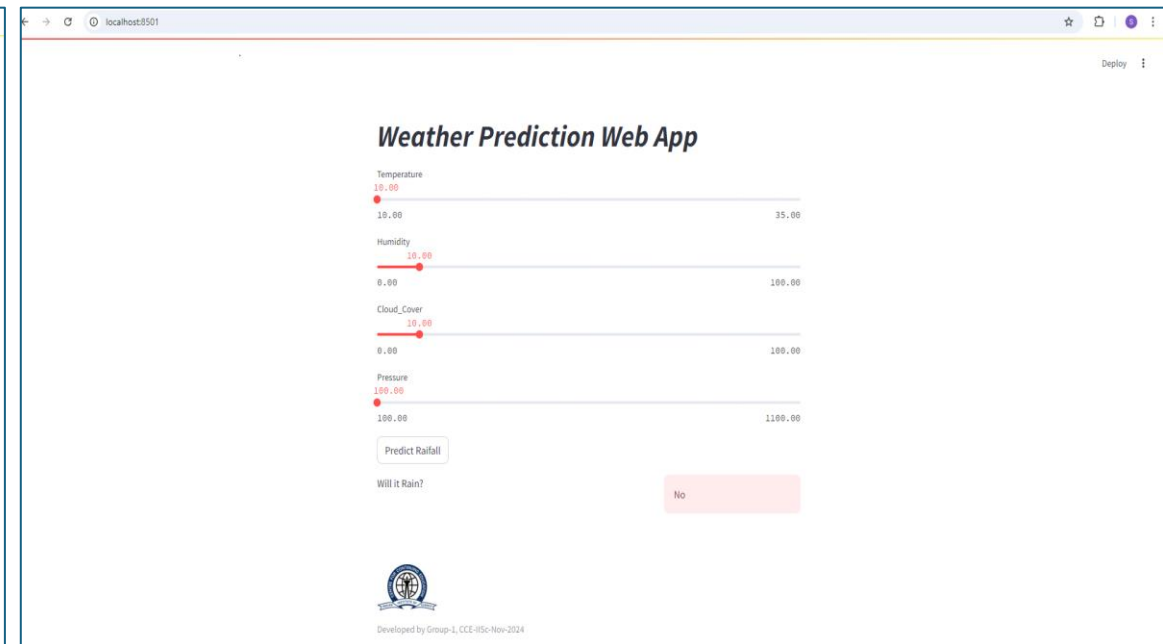
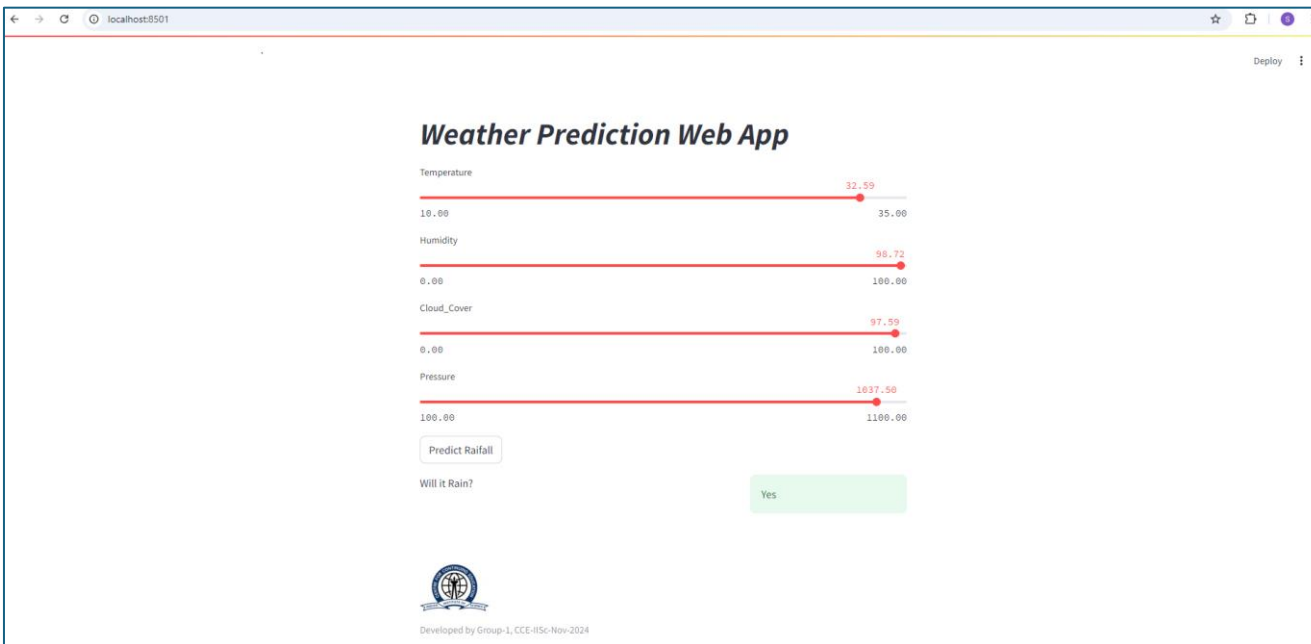
Columns

Group by

	Run Name	Created	Dataset	Duration	Source	Models
<input type="checkbox"/>	<div><div></div>serious-horse-401</div>	<div><div></div>25 minutes ago</div>	-	9.6s	<div>c:\Users\...</div>	<div>ML-model-LR v9</div>
<input type="checkbox"/>	<div><div></div>angry-steed-752</div>	<div><div></div>7 hours ago</div>	-	4.8s	<div>c:\Users\...</div>	<div>ML-model-LR v8</div>
<input type="checkbox"/>	<div><div></div>monumental-chimp-451</div>	<div><div></div>7 hours ago</div>	-	5.0s	<div>c:\Users\...</div>	<div>ML-model-LR v7</div>
<input type="checkbox"/>	<div><div></div>dashing-calf-550</div>	<div><div></div>7 hours ago</div>	-	7.7s	<div>c:\Users\...</div>	<div>ML-model-LR v6</div>
<input type="checkbox"/>	<div><div></div>peaceful-newt-731</div>	<div><div></div>19 hours ago</div>	-	8.7s	<div>c:\Users\...</div>	<div>ML-model-LR v5</div>
<input type="checkbox"/>	<div><div></div>puzzled-dove-875</div>	<div><div></div>2 days ago</div>	-	5.2s	<div>c:\Users\...</div>	<div>ML-model-LR v4</div>
<input type="checkbox"/>	<div><div></div>caring-sheep-520</div>	<div><div></div>2 days ago</div>	-	5.3s	<div>c:\Users\...</div>	<div>ML-model-LR v3</div>
<input type="checkbox"/>	<div><div></div>luxuriant-vole-856</div>	<div><div></div>2 days ago</div>	-	5.7s	<div>c:\Users\...</div>	<div>ML-model-LR v2</div>
<input type="checkbox"/>	<div><div></div>hilarious-bird-346</div>	<div><div></div>2 days ago</div>	-	4.5s	<div>c:\Users\...</div>	<div>ML-model-LR v1</div>

# MLOps

Webapp using Streamlit to do live prediction using user inputs by fetching model from MLFlow registry





# Contributions

Team Members	Data cleanup and acquisition	ML model selection / training	Hyper parameter tuning/MLops	Metrics	Presentation	Documentation
Aniket Shukla	5/ 100	5/ 100	5/ 100	10/ 100	5/ 100	5/ 100
Abhishek Yadav	30/ 100	30/ 100	5/ 100	20/ 100	10/ 100	30/ 100
Sagar Kamble	10/ 100	15/ 100	70/ 100	10/ 100	10/ 100	10/ 100
Sandeep Mamoriya	10/ 100	5/ 100	5/ 100	20/ 100	35/ 100	20/ 100
Sherine Martina	5/ 100	5/ 100	5/ 100	10/ 100	30/ 100	25/ 100
Shravya Pendyala	30/ 100	30/ 100	5/ 100	20/ 100	5/ 100	5/ 100
Sridhar Mulumoodi	10/ 100	10/ 100	5/ 100	10/ 100	5/ 100	5/ 100
Total	100 / 100	100 / 100	100 / 100	100 / 100	100 / 100	100 / 100

# Conclusions

- Using EDA, we can nicely analyse input data/features and based on such analysis, it helps to choose/drop any feature based on their relevance
- It was helpful to check available classification models and compare them for our problem statement to get best fit. The maximum accuracy attained is using the XGBoost model which is 100%. The model looks overfitting so we chose RF model to conclude
- Using MLOps, it was very helpful to use trained ML model in reality. Streamlit helped to create webapp for user inputs and MLFlow helped to store trained models. With this, we achieved MLOps in low cost manner by hosting it on local PC servers. This approach helps in organizations who wants to keep sensitive data and models internal. Otherwise, cloud options can be used for further hosting of models.
- As a futuristic scope, long-term predictions can be estimated leading to the forecasting of rainfall conditions in a particular region