



VIT[®]

Vellore Institute of Technology

(Deemed to be University under section 3 of UGC Act, 1956)

CHENNAI

RAINFALL PREDICTION

Name: Kush Desai, Raginee Titar, Raviraj Patil

Reg. no: 20MCB1002, 20MCB1004, 20MCB1016

Subject: Domain Specific Predictive Analytics(CSE6021)

OVERVIEW

- Objective
- Supervised learning
- Binary classification
- Methodology used for binary classification:
- Dataset description:
- Sample Dataset
- Procedure to predict the rain tomorrow
- Data preparation visualization
- Linear Regression to Rainfall Prediction
- Decision tree Regression to Rainfall Prediction
- Random Forest Regression to Rainfall Prediction
- Conclusion

OBJECTIVE:

Rainfall has a significant impact on society. Festivals, all sorts of activities, and sports matches are examples of events that could be heavily influenced by rainfall. Therefore, this topic is studied heavily. The goal of this analysis is to predict if it is going to rain the next day, based on the weather measures of the day before. This creates a binary classification problem (1 = Rain and 0 = No Rain). The dataset consists of weather data in Australia over 10 years (1 Nov 2007 to 25 Jun 2017).

SUPERVISED LEARNING

In supervised learning, algorithms learn from labeled data. After understanding the data, the algorithm determines which label should be given to new data by associating patterns to the unlabeled new data. Supervised learning can be divided into two categories: classification and regression.

Supervised learning can be divided into two categories: classification and regression.

Classification is a technique for determining which class the dependent belongs to based on one or more independent variables.

Logistic regression is kind of like linear regression, but is used when the dependent variable is not a number but something else (e.g., a "yes/no" response). It's called regression but performs classification based on the regression and it classifies the dependent variable into either of the classes.

BINARY CLASSIFICATION:

- Binary classification is the task of classifying the elements of a set into two groups on the basis of a classification rule. Accurate Sales Forecasts enable companies to make informed business decisions and predict short-term and long-term performance.
- Binary classification is dichotomization applied to a practical situation. In many practical binary classification problems, the two groups are not symmetric, and rather than overall accuracy, the relative proportion of different types of errors is of interest. For example, in medical testing, detecting a disease when it is not present (a false positive) is considered differently from not detecting a disease when it is present (a false negative).

METHODOLOGY USED FOR BINARY CLASSIFICATION:

The following are the algorithms that is be used binary classification;

- Decision Tree
- Logistic Regression
- Random Forest

DATASET DESCRIPTION:

Name: Rain in Australia

Size: 13.44 MB

Source: Kaggle

DATASET NAME:

weatherAUS.csv

ATTRIBUTES:

- **DATE-** The date represents each day Rainfall

- **Location**-The common name of the location of the weather station
- **MinTemp**- The minimum temperature in degrees celsius
- **MaxTemp**-The maximum temperature in degrees celsius
- **Rainfall**-The amount of rainfall recorded for the day in mm
- **Evaporation**-The so-called Class A pan evaporation (mm) in the 24 hours to 9am
- **Sunshine**-The number of hours of bright sunshine in the day
- **WindGustDir**-The direction of the strongest wind gust in the 24 hours to midnight
- **WindGustSpeed**-The speed (km/h) of the strongest wind gust in the 24 hours to midnight
- **WindDir9am**-Direction of the wind at 9am
- **WindDir3pm**-Direction of the wind at 3pm
- **WindSpeed9am**-Wind speed (km/hr) averaged over 10 minutes prior to 9am
- **WindSpeed3pm**-Wind speed (km/hr) averaged over 10 minutes prior to 3pm
- **Humidity9am**-Humidity (percent) at 9am
- **Humidity3pm**-Humidity (percent) at 3pm
- **Pressure9am**-Atmospheric pressure (hpa) reduced to mean sea level at 9am
- **Pressure3pm**-Atmospheric pressure (hpa) reduced to mean sea level at 3pm

- **Cloud9am**- Fraction of sky obscured by cloud at 9am. This is measured in "oktas", which are a unit of eighths. It records how many eighths of the sky are obscured by cloud. A 0 measure indicates completely clear sky whilst an 8 indicates that it is completely overcast.
- **Cloud3pm**-Fraction of sky obscured by cloud (in "oktas": eighths) at 3pm. See Cloud9am for a description of the values
- **Temp9am**-Temperature (degrees C) at 9am
- **Temp3pm**-Temperature (degrees C) at 3pm
- **RainToday**-Boolean: 1 if precipitation (mm) in the 24 hours to 9am exceeds 1mm, otherwise 0
- **RainTomorrow**-The amount of next day rain in mm. Used to create response variable RainTomorrow. A kind of measure of the "risk".

ABOUT DATASET:

This dataset contains about 10 years of daily weather observations from many locations across Australia.

RainTomorrow is the target variable to predict. It means -- did it rain the next day, Yes or No? This column is Yes if the rain for that day was 1mm or more.

SAMPLE DATASET:

A1																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																					
----	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

PROCEDURE TO PREDICT THE RAIN TOMORROW:

STEP 1: We imported the packages namely NumPy, pandas, matplotlib lib, seaborn, sklearn.

STEP 2: Data Preparation

In Data Preparation we undergo following steps namely,

2.1 Load dataset

2.2 Dealing with missing value

2.3 Count missing values in each column

2.4 Replacing the missing values in features with mean, median and drop null value.

STEP 3: Feature Selection

Select Feature Importance using Filter Method (Chi-Square)

STEP 4: Exploratory Data Analytics

5.1 Correlation Heatmap

STEP 6: Rain tomorrow Prediction

Here we apply various classification techniques to predict the sales.

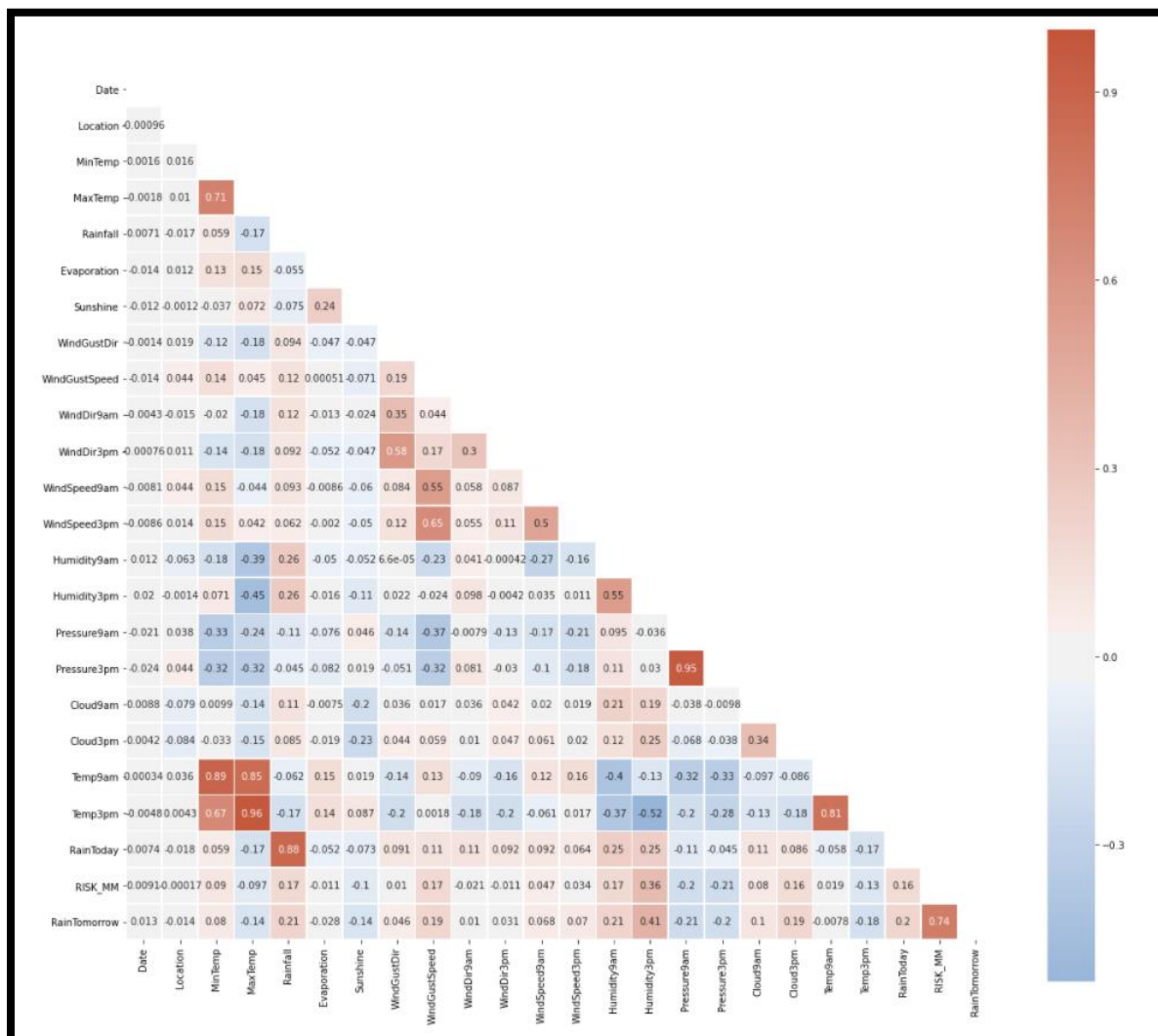
6.1 Logistic Regression

6.2 Decision Tree

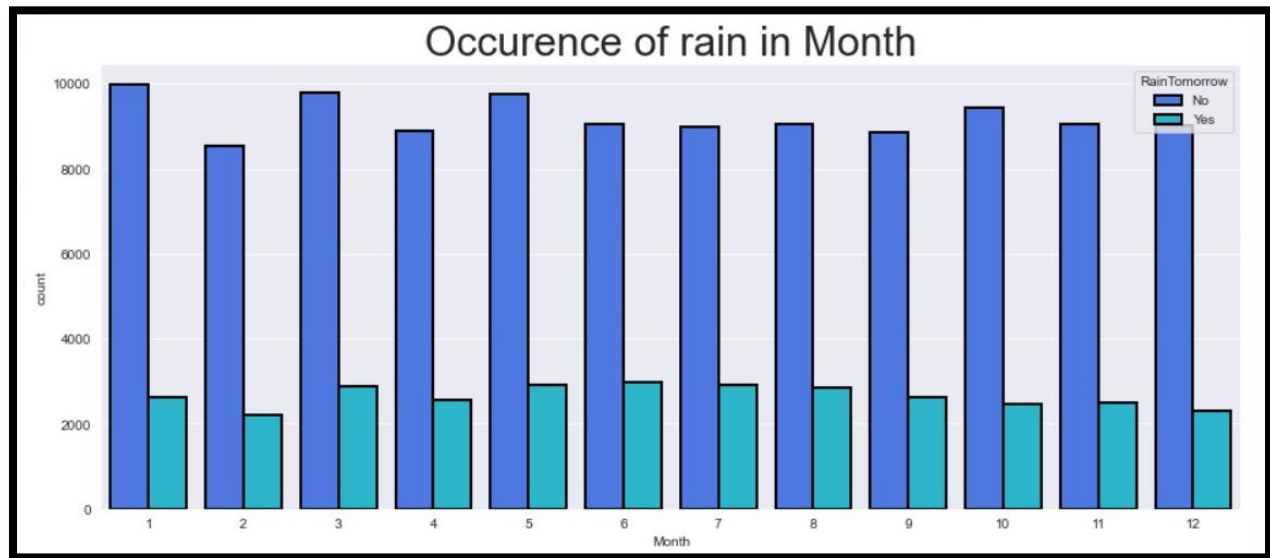
6.3 Random Forest

DATA PREPERATION VISUALIZATION

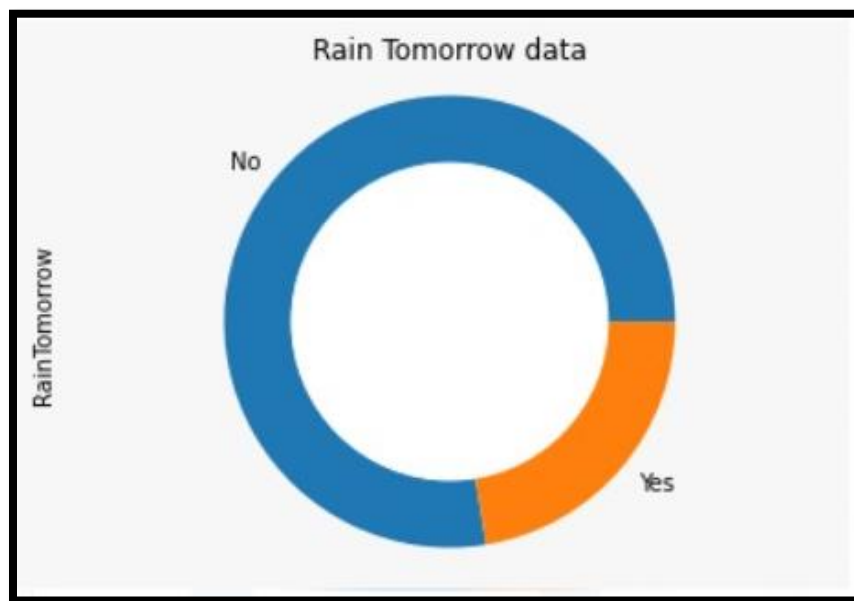
CORRELATION HEATMAP



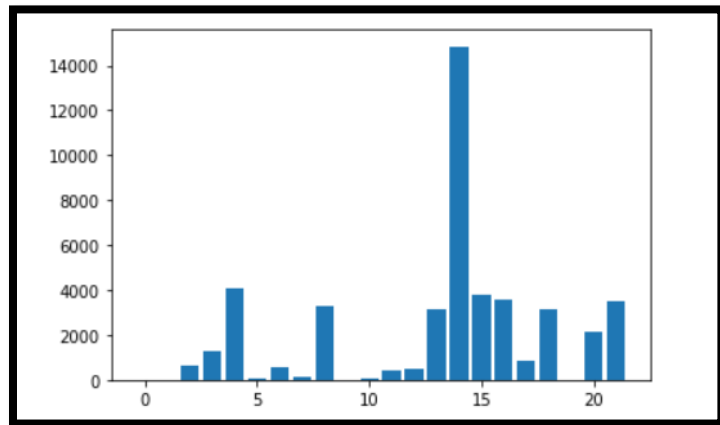
OCCURRENCE OF RAIN IN MONTH



CLASS DISTRIBUTION



FEATURE SELCION USING Chi-Square METHOD



LOGISTIC REGRESSION TO RAINFALL PREDICTION:

- To perform logistic regression, the sigmoid function. In a binary logistic regression model, the dependent variable has two levels (IN this case yes or no). Outputs with more than two values are modeled by multinomial logistic regression and, if the multiple categories are ordered, by ordinal logistic regression.
- The logistic regression model itself simply models probability of output in terms of input and does not perform statistical classification, though it can be used to make a classifier, for instance by choosing a cutoff value and classifying inputs with probability greater than the cutoff as one class, below the cutoff as the other; this is a common way to make a binary classifier.

DECISION TREE REGRESSION TO RAINFALL PREDICTION:

- Decision tree builds classification models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes.
- A decision node has two or more branches. Leaf node represents a classification or decision. The topmost decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data.

RANDOM FOREST REGRESSION TO RAINFALL PREDICTION:

- The random forest is a model made up of many decision trees. Rather than just simply averaging the prediction of trees (which we could call a “forest”)
- The random forest combines hundreds or thousands of decision trees, trains each one on a slightly different set of the observations, splitting nodes in each tree considering a limited number of the features. The final predictions of the random forest are made by averaging the predictions of each individual tree.

SAMPLE CODE:

LOGISTICREGRESSION

```
from sklearn.linear_model import LogisticRegression
from sklearn import metrics
from sklearn.metrics import precision_recall_fscore_support
from sklearn import metrics
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

X_train, X_test, y_train, y_test =
train_test_split(features,target,test_size=0.1,random_state=1)
logreg = LogisticRegression()
logreg.fit(X_train, y_train)
y_pred= logreg.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: %.2f%%" % (accuracy * 100.0))
print("Precision Score:-", metrics.precision_score(y_test, y_pred))
print("Recall Score:-", metrics.recall_score(y_test, y_pred))
print("F1 Score:-", metrics.f1_score(y_test, y_pred))
#print("Average Precision Score:-",
metrics.average_precision_score(y_test, predictions))
#print("Log Loss:-", metrics.log_loss(y_test, predictions))
print("ROC-AUC Score:-", metrics.roc_auc_score(y_test, y_pred))
#print(precision_recall_fscore_support(y_test, predictions,
average='binary'))
```

Decision Tree:

```
from sklearn import metrics
from sklearn.metrics import precision_recall_fscore_support
from sklearn import metrics
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.tree import DecisionTreeClassifier
X_train, X_test, y_train, y_test =
train_test_split(features,target,test_size=0.1,random_state=1)

dtree=DecisionTreeClassifier()
dtree.fit(X_train,y_train)
y_pred = dtree.predict(X_test)

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: %.2f%%" % (accuracy * 100.0))
print("Precision Score:-", metrics.precision_score(y_test, y_pred))
print("Recall Score:-", metrics.recall_score(y_test, y_pred))
print("F1 Score:-", metrics.f1_score(y_test, y_pred))
#print("Average Precision Score:-",
metrics.average_precision_score(y_test, predictions))
#print("Log Loss:-", metrics.log_loss(y_test, predictions))
print("ROC-AUC Score:-", metrics.roc_auc_score(y_test, y_pred))
```

```
#print(precision_recall_fscore_support(y_test, predictions,
average='binary'))
```

RANDOM FOREST

```
from sklearn import metrics
from sklearn.metrics import precision_recall_fscore_support
from sklearn.ensemble import RandomForestClassifier
from sklearn import metrics
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score

x_train, x_test, y_train, y_test= train_test_split(features, target,
test_size= 0.1, random_state=1)
#class_weight={0:1,1:3}
#feature Scaling
st_x= StandardScaler()
x_train= st_x.fit_transform(x_train)
x_test= st_x.transform(x_test)

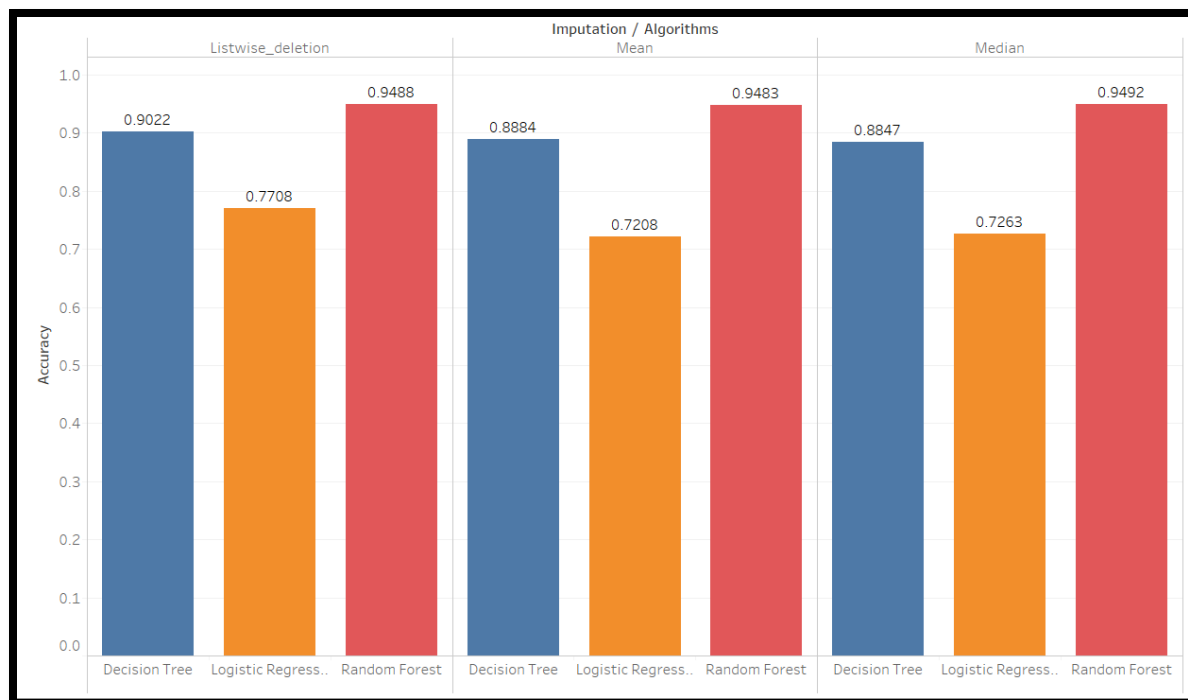
classifier= RandomForestClassifier(n_estimators= 100,
criterion="entropy",class_weight='balanced')
classifier.fit(x_train, y_train)
y_pred= classifier.predict(x_test)

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: %.2f%%" % (accuracy * 100.0))
print("Precision Score:-", metrics.precision_score(y_test, y_pred))
print("Recall Score:-", metrics.recall_score(y_test, y_pred))
print("F1 Score:-", metrics.f1_score(y_test, y_pred))
#print("Average Precision Score:-",
metrics.average_precision_score(y_test, predictions))
#print("Log Loss:-", metrics.log_loss(y_test, predictions))
print("ROC-AUC Score:-", metrics.roc_auc_score(y_test, y_pred))
#print(precision_recall_fscore_support(y_test, predictions,
average='binary'))
```

OUTPUT

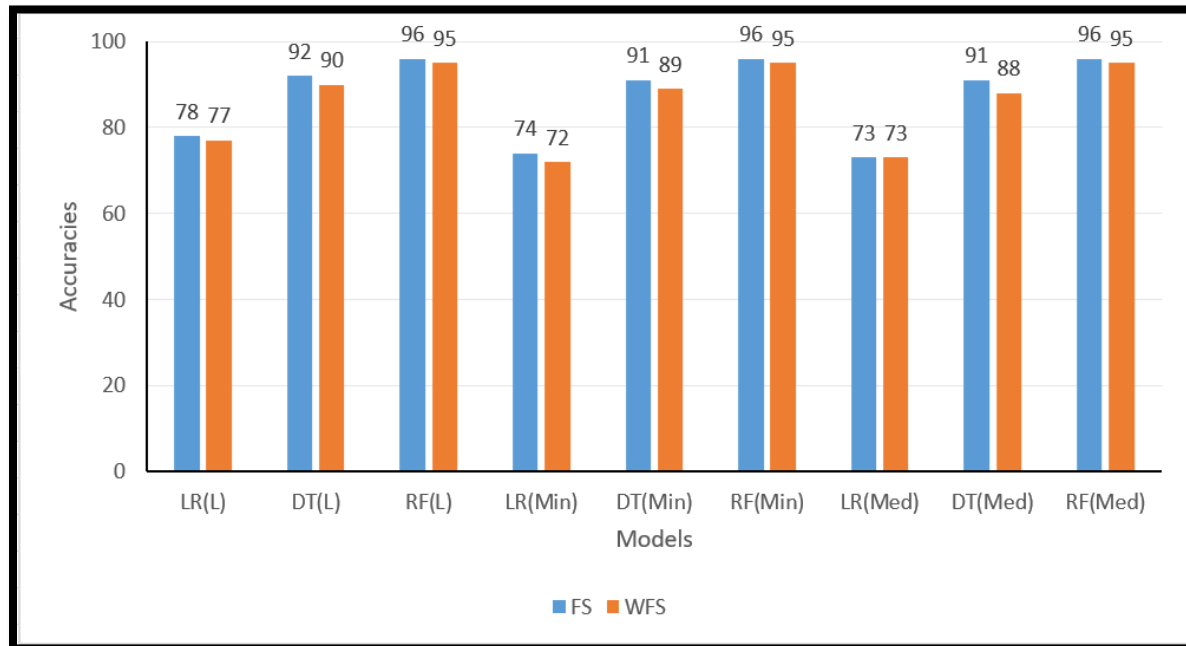
1) Without Feature Selection

Algorithms	Accuracy			Precision score			Recall score			F1 Score		
Model	Mean	Median	Drop null	Mean	Median	Drop null	Mean	Median	Drop null	Mean	Median	Drop null
Logistic Regression	72	73	77	0.70	0.73	0.77	0.55	0.55	0.64	0.62	0.62	0.98
Decision Tree	89	88	90	0.81	0.88	0.90	0.95	0.95	0.96	0.87	0.87	0.88
Random Forest	95	95	95	0.92	0.95	0.95	0.95	0.95	0.96	0.94	0.94	0.94



2) With Feature Selection

Algorithms	Accuracy			Precision score			Recall score			F1 Score		
Model	Mean	Median	Drop null	Mean	Median	Drop null	Mean	Median	Drop null	Mean	Median	Drop null
Logistic Regression	74	73	78	0.71	0.69	0.74	0.59	0.58	0.66	0.64	0.63	0.70
Decision Tree	91	91	92	0.83	0.83	0.84	0.97	0.97	0.98	0.90	0.89	0.90
Random Forest	96	96	96	0.94	0.93	0.93	0.97	0.97	0.98	0.95	0.95	0.95



LR: Logistic Regression

DT: Decision Trees

RF: Random Forest

(L): Missing Value imputation using Listwise deletion

(Min): Missing Value imputation using Mean

(Med): Missing Value imputation using Median

FD: With feature selection

WFS: Without Feature Selection

Inference: Random Forest gives best accuracy for all types of missing value imputation.

CONCLUSION:

In this project, handled missing values by using the mean, listwise deletion, and median imputation. Over-sampling method was used to create “synthetic” data to balance dataset. For feature selection Chi-Square method is used to select best top 15 columns to get better accuracy. We consider different binary classification machine learning approach for tomorrow rain fall prediction. From the all the three model Random Forest has got the best accuracy i.e., 96% accuracy.