

LP1 - Data Analytics Mini project

Problem Statement:-

To classify air quality of a particular region based on the given parameters.

Team Members :

41250 - Sahil Patil

41254 - Pralhad Kulkarni

41258 - Srusti Raybhoge

Abstract:-

The main motive of the project is to classify the air quality of the input parameters by creating a model with the help of a given dataset . The attributes of the dataset include the amount of various elements present in the air at that very moment. We have used different ML algorithms to solve the same problem and create a model

H/W & S/W requirements:-

Python, Python libraries, Jupyter Notebook/Google Colab/ Python IDE.

Objective:-

1. To understand classification models
2. To predict air quality on the given inputs
3. Understand data preprocessing

Scope:-

1. To predict the air quality based on the given factors, place and time.
2. Study two classifiers

Introduction:-

Classification: refers to a predictive modeling problem where a class label is predicted for a given example of input data.

Types of Classification:

1. On the basis of Labels:
 - a. Binary labelled
 - b. Multi labelled

2. Linear Models:

- a. Logistic Regression
- b. Support Vector Machines

3. Nonlinear Models

- a. K-nearest Neighbors (KNN)
- b. Kernel Support Vector Machines (SVM)
- c. Naïve Bayes
- d. Decision Tree Classification
- e. Random Forest Classification

Data Preprocessing Steps:

- 1. Remove empty columns
- 2. Fill in empty cells with appropriate values (here median of the values in that particular column)
- 3. Convert the string type data into equivalent numeric data.

SMOT (Synthetic Minority Oversampling Technique):

Sometimes the dataset may be imbalanced i.e. some classes have a majority number of samples than other classes. This may affect the model and we get poor performance with minority class.

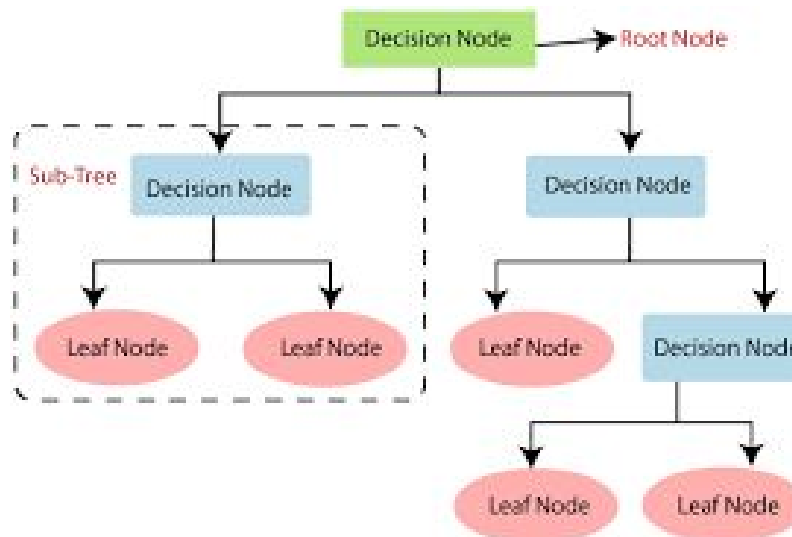
To avoid this we can do oversampling or undersampling of data. An improvement on duplicating examples from the minority class is to synthesize new examples from the minority class. This is a type of data augmentation for tabular data and can be very effective.

We have used random forest and decision tree classification to classify the dataset into appropriate air quality class and predict the class for the input.

Decision Tree:

Decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.

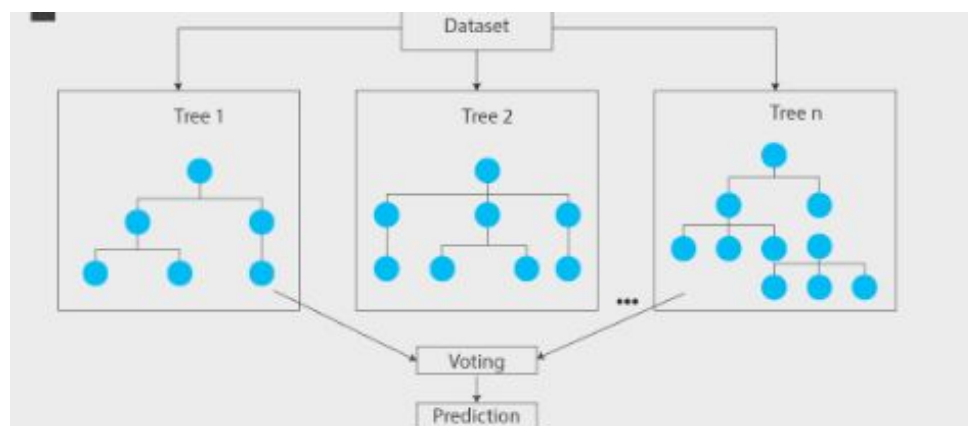
In a Decision tree, there are two nodes, which are the Decision Node and Leaf Node . Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.



Random Forest:

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.



Features from the dataset:

City, Date, PM2.5, PM10, NO, NO2, NOx, NH3, CO, SO2, O3, Benzene, Toluene, Xylene, AQI, Air_quality

Results:-

Dataset overview:

colab.research.google.com/drive/1clrZ66kL17qs5sdJMXOWJAEvgxZCPvNY#scrollTo=zmS21kpA_vTI

AIR_QUALITY.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

```
dataset = pd.read_csv("/content/City_day.csv");
dataset.dropna(axis=0, subset = ["Air_quality", "Xylene", "AQI", "Toluene", "Benzene", "O3", "SO2", "CO", "NH3", "NOx", "NO2", "PM10", "PM2.5", "NO"], how = 'all')
dataset.dropna(subset = ["Air_quality"], inplace=True)
x = dataset.iloc[:, :-1].values
y = dataset.iloc[:, 15].values
print(dataset.head())
```

	City	Date	PM2.5	PM10	...	Toluene	Xylene	AQI	Air_quality
28	Ahmedabad	1/29/2015	83.13	NaN	...	0.00	3.14	209.0	Poor
29	Ahmedabad	1/30/2015	79.84	NaN	...	0.00	4.81	328.0	Very Poor
30	Ahmedabad	1/31/2015	94.52	NaN	...	0.01	7.67	514.0	Severe
31	Ahmedabad	2/1/2015	135.99	NaN	...	0.04	25.87	782.0	Severe
32	Ahmedabad	2/2/2015	178.33	NaN	...	0.06	35.61	914.0	Severe

[5 rows x 16 columns]

```
# Check for possible null values in the dataset as missing values potentially screw the ml models
dataset.isnull().sum()
```

City	0
Date	0
PM2.5	635
PM10	6975
NO	308
NO2	305
NOx	1785
NH3	6323
CO	416
SO2	554
O3	709
Benzene	3270
Toluene	5435
Xylene	14208
AQI	0
Air_quality	0
dtype:	int64

colab.research.google.com/drive/1clrZ66kL17qs5sdJMXOWJAEvgxZCPvNY#scrollTo=HOr4oojb_vTI

AIR_QUALITY.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

```
[6]
```

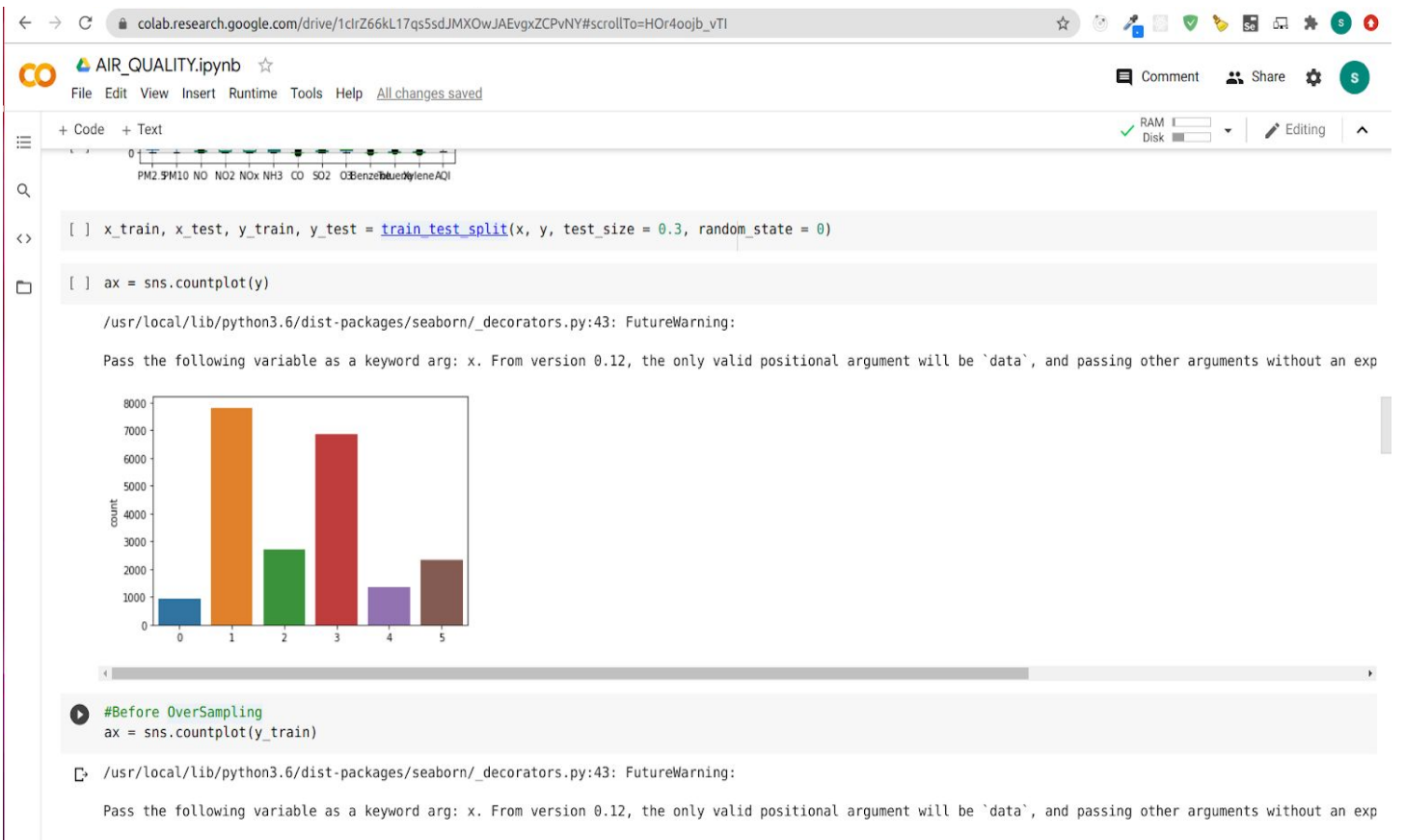
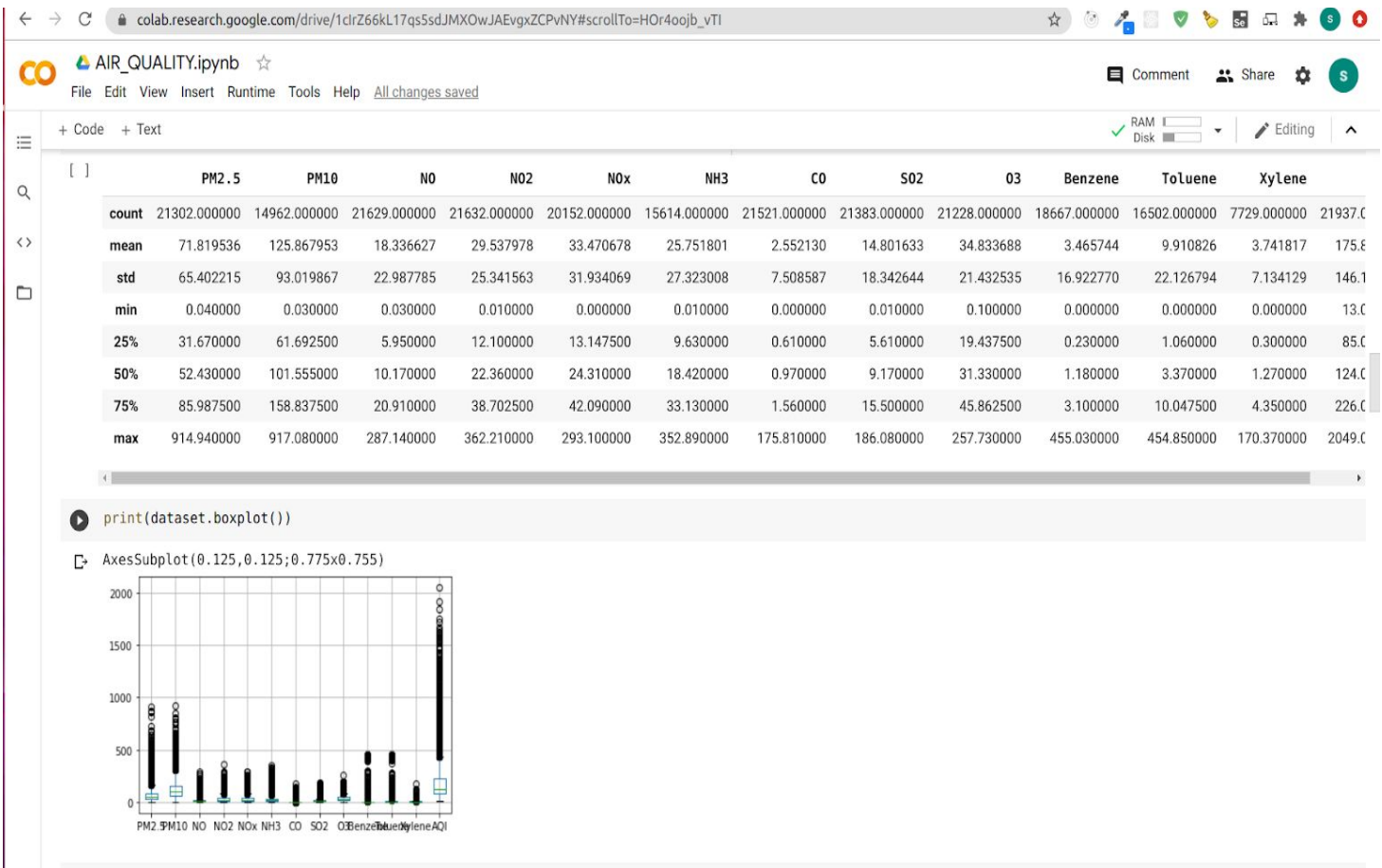
```
le_X_city = LabelEncoder()
le_X_date = LabelEncoder()
le_Y = LabelEncoder()
y = le_Y.fit_transform(y)

x[:,0] = le_X_city.fit_transform(x[:,0])
x[:,1] = le_X_date.fit_transform(x[:,1])
```

```
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 21937 entries, 28 to 26218
Data columns (total 16 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   City            21937 non-null  object
 1   Date            21937 non-null  object
 2   PM2.5           21302 non-null  float64
 3   PM10            14962 non-null  float64
 4   NO              21629 non-null  float64
 5   NO2             21632 non-null  float64
 6   NOx             20152 non-null  float64
 7   NH3             15614 non-null  float64
 8   CO              21521 non-null  float64
 9   SO2             21383 non-null  float64
10   O3              21228 non-null  float64
11   Benzene         18667 non-null  float64
12   Toluene         16502 non-null  float64
13   Xylene          7729 non-null   float64
14   AQI             21937 non-null  float64
15   Air_quality     21937 non-null  object
dtypes: float64(13), object(3)
memory usage: 2.8+ MB
```

Mean, Max, Min, Standard deviation, Count of each column



colab.research.google.com/drive/1c1rZ66kL17qs5sdJMXOwJAEvgxZCPvNY#scrollTo=HOr4oojb_vTI

AIR_QUALITY.ipynb

File Edit View Insert Runtime Tools Help All changes saved

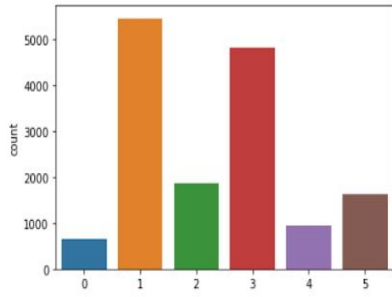
+ Code + Text

RAM Disk Editing

```
[ ] #Before OverSampling
ax = sns.countplot(y_train)
```

/usr/local/lib/python3.6/dist-packages/seaborn/_decorators.py:43: FutureWarning:

Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an exp



```
print('Classes and number of values in trainset',Counter(y_train))
from imblearn.over_sampling import SMOTE
oversample = SMOTE()
x_train, y_train = oversample.fit_resample(x_train,y_train)
print('Classes and number of values in trainset after SMOTE:',Counter(y_train))
med=np.median(x_train,axis=0)
```

Classes and number of values in trainset Counter({1: 5450, 3: 4808, 2: 1867, 5: 1618, 4: 951, 0: 661})

/usr/local/lib/python3.6/dist-packages/sklearn/externals/six.py:31: FutureWarning:

colab.research.google.com/drive/1c1rZ66kL17qs5sdJMXOwJAEvgxZCPvNY#scrollTo=HOr4oojb_vTI

AIR_QUALITY.ipynb

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

RAM Disk Editing

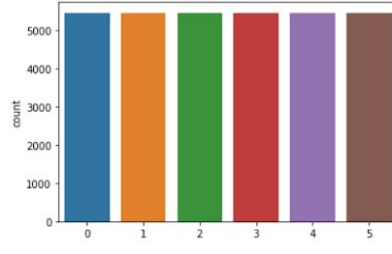
```
[ ] Function safe_indexing is deprecated; safe_indexing is deprecated in version 0.22 and will be removed in version 0.24.
Classes and number of values in trainset after SMOTE: Counter({3: 5450, 1: 5450, 2: 5450, 5: 5450, 4: 5450, 0: 5450})
/usr/local/lib/python3.6/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning:
Function safe_indexing is deprecated; safe_indexing is deprecated in version 0.22 and will be removed in version 0.24.
/usr/local/lib/python3.6/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning:
Function safe_indexing is deprecated; safe_indexing is deprecated in version 0.22 and will be removed in version 0.24.
```

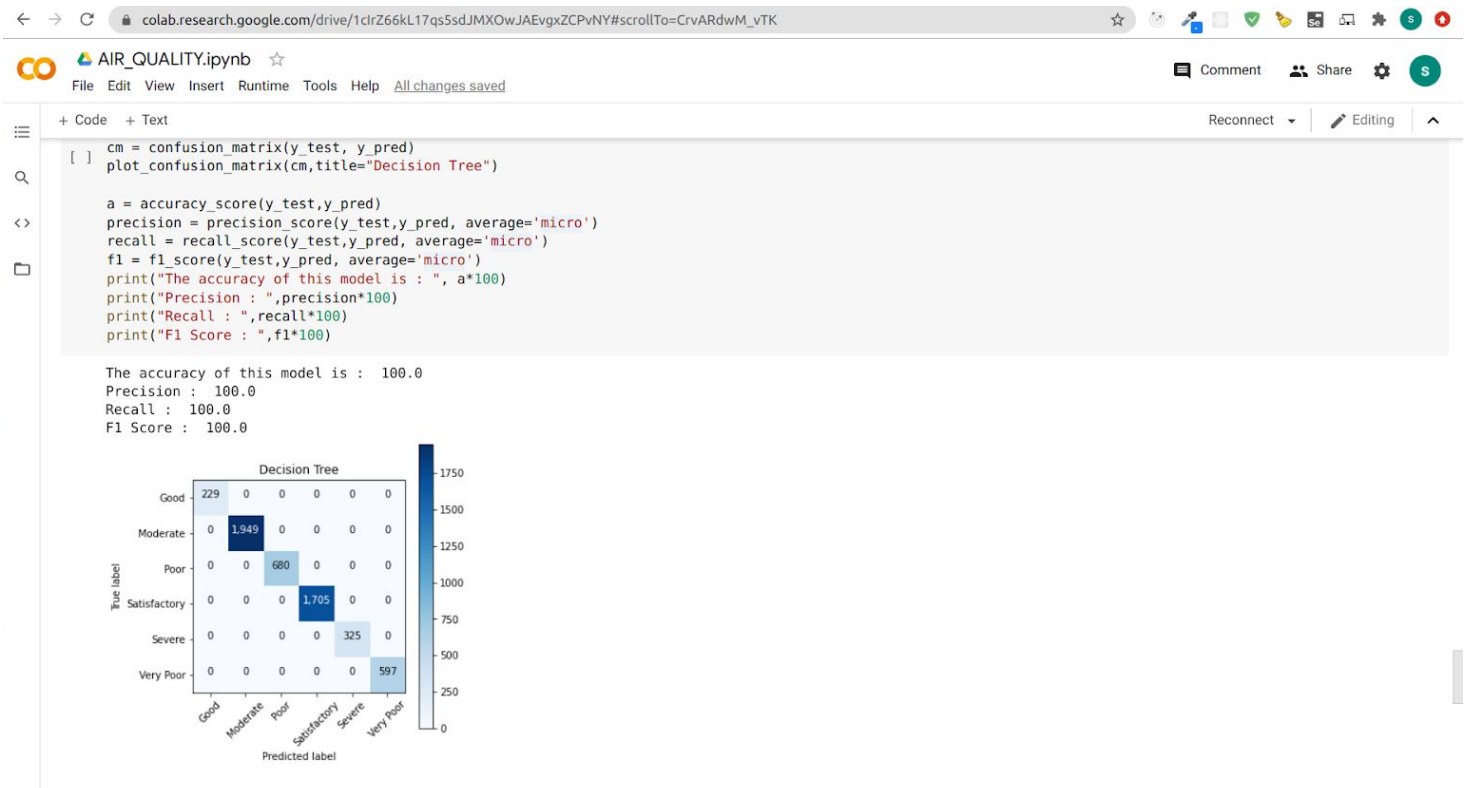
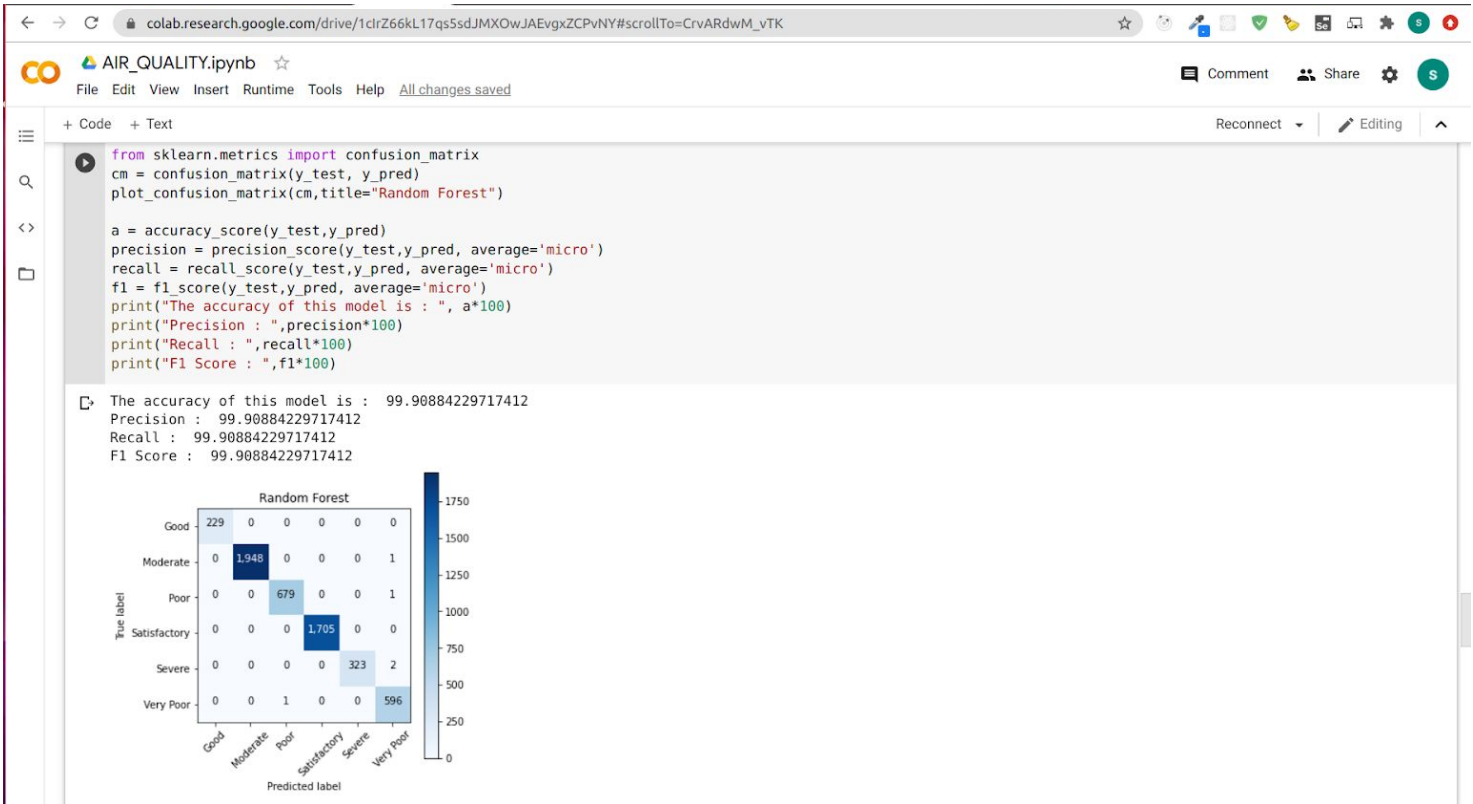
```
[ ] #After Oversampling
sns.countplot(y_train)
```

/usr/local/lib/python3.6/dist-packages/seaborn/_decorators.py:43: FutureWarning:

Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an exp

<matplotlib.axes._subplots.AxesSubplot at 0x7fc687c26e10>





TEST CASES:

```
[16] res = []
city = le_X_city.fit_transform(["Pune"])
date = le_X_date.fit_transform(["12/12/2020"])
ls = [city[0], date[0], 83.13, 101.555, 6.93, 28.71, 33.72, 18.42, 6.93, 49.52, 59.76, 0.02, 0.0, 3.14, 209.0]
lst = []
lst.append(ls)
temp = le_Y.inverse_transform(clf.predict(lst))
temp = temp.tolist()
res.append(temp[0])

temp = le_Y.inverse_transform(dt.predict(lst))
temp = temp.tolist()
res.append(temp[0])

print("Random Forest : ", res[0])
print("Decision Tree : ", res[1])
```

Random Forest : Poor
Decision Tree : Poor

```
res = []
city = le_X_city.fit_transform(["Patna"])
date = le_X_date.fit_transform(["12/12/2020"])
ls = [city[0], date[0], 49.34, 0, 20.33, 26.7, 21.5, 0, 0.91, 49.93, 49.47, 1.43, 16.43, 5.49, 94]
lst = []
lst.append(ls)
temp = le_Y.inverse_transform(clf.predict(lst))
temp = temp.tolist()
res.append(temp[0])

temp = le_Y.inverse_transform(dt.predict(lst))
temp = temp.tolist()
res.append(temp[0])
```

```
res = []
city = le_X_city.fit_transform(["Patna"])
date = le_X_date.fit_transform(["12/12/2020"])
ls = [city[0], date[0], 49.34, 0, 20.33, 26.7, 21.5, 0, 0.91, 49.93, 49.47, 1.43, 16.43, 5.49, 94]
lst = []
lst.append(ls)
temp = le_Y.inverse_transform(clf.predict(lst))
temp = temp.tolist()
res.append(temp[0])

temp = le_Y.inverse_transform(dt.predict(lst))
temp = temp.tolist()
res.append(temp[0])

print("Random Forest : ", res[0])
print("Decision Tree : ", res[1])
```

Random Forest : Satisfactory
Decision Tree : Satisfactory

Conclusion:-

Thus we studied different classification algorithms and used two of them to predict the quality of air in a particular region at a particular time.