# Importing Libraries

In [1]:

```python
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

# Importing Datasets

In [2]:

```
df=pd.read_csv(r"C:\Users\user\Downloads\C10_air\csvs_per_year\csvs(Dataset)\madrid_2011.
df
```

Out[2]:

| | date | BEN | CO | EBE | NMHC | NO | NO_2 | O_3 | PM10 | PM25 | SO_2 | TCH | TOI |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2011-11-01 01:00:00 | NaN | 1.0 | NaN | NaN | 154.0 | 84.0 | NaN | NaN | NaN | 6.0 | NaN | NaN |
| 1 | 2011-11-01 01:00:00 | 2.5 | 0.4 | 3.5 | 0.26 | 68.0 | 92.0 | 3.0 | 40.0 | 24.0 | 9.0 | 1.54 | 8.7 |
| 2 | 2011-11-01 01:00:00 | 2.9 | NaN | 3.8 | NaN | 96.0 | 99.0 | NaN | NaN | NaN | NaN | NaN | 7.2 |
| 3 | 2011-11-01 01:00:00 | NaN | 0.6 | NaN | NaN | 60.0 | 83.0 | 2.0 | NaN | NaN | NaN | NaN | NaN |
| 4 | 2011-11-01 01:00:00 | NaN | NaN | NaN | NaN | 44.0 | 62.0 | 3.0 | NaN | NaN | 3.0 | NaN | NaN |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | . |
| 209923 | 2011-09-01 00:00:00 | NaN | 0.2 | NaN | NaN | 5.0 | 19.0 | 44.0 | NaN | NaN | NaN | NaN | NaN |
| 209924 | 2011-09-01 00:00:00 | NaN | 0.1 | NaN | NaN | 6.0 | 29.0 | NaN | 11.0 | NaN | 7.0 | NaN | NaN |
| 209925 | 2011-09-01 00:00:00 | NaN | NaN | NaN | 0.23 | 1.0 | 21.0 | 28.0 | NaN | NaN | NaN | 1.44 | NaN |
| 209926 | 2011-09-01 00:00:00 | NaN | NaN | NaN | NaN | 3.0 | 15.0 | 48.0 | NaN | NaN | NaN | NaN | NaN |
| 209927 | 2011-09-01 00:00:00 | NaN | NaN | NaN | NaN | 4.0 | 33.0 | 38.0 | 13.0 | NaN | NaN | NaN | NaN |

209928 rows × 14 columns

# Data Cleaning and Data Preprocessing

In [3]:

```
df=df.dropna()
```

In [4]:

```python
df.columns
```

Out[4]:

```
Index(['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'P
M25',
       'SO_2', 'TCH', 'TOL', 'station'],
      dtype='object')
```

In [5]:

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 16460 entries, 1 to 209910
Data columns (total 14 columns):
 #   Column   Non-Null Count   Dtype
---  ------   --------------   -----
 0   date     16460 non-null   object
 1   BEN      16460 non-null   float64
 2   CO       16460 non-null   float64
 3   EBE      16460 non-null   float64
 4   NMHC     16460 non-null   float64
 5   NO       16460 non-null   float64
 6   NO_2     16460 non-null   float64
 7   O_3      16460 non-null   float64
 8   PM10     16460 non-null   float64
 9   PM25     16460 non-null   float64
 10  SO_2     16460 non-null   float64
 11  TCH      16460 non-null   float64
 12  TOL      16460 non-null   float64
 13  station  16460 non-null   int64
dtypes: float64(12), int64(1), object(1)
memory usage: 1.9+ MB
```

In [6]:

```python
data=df[['BEN', 'TOL', 'TCH']]
data
```

Out[6]:

|        | BEN | TOL | TCH  |
|--------|-----|-----|------|
| 1      | 2.5 | 8.7 | 1.54 |
| 6      | 0.7 | 1.7 | 1.36 |
| 25     | 1.8 | 7.4 | 1.71 |
| 30     | 1.0 | 2.9 | 1.40 |
| 49     | 1.3 | 6.2 | 1.75 |
| ...    | ... | ... | ...  |
| 209862 | 0.4 | 0.7 | 1.26 |
| 209881 | 0.9 | 4.9 | 1.34 |
| 209886 | 0.6 | 0.9 | 1.26 |
| 209905 | 0.6 | 3.8 | 1.32 |
| 209910 | 0.7 | 0.9 | 1.25 |

16460 rows × 3 columns

# Line chart

In [7]:

```python
data.plot.line(subplots=True)
```

Out[7]:

```
array([<AxesSubplot:>, <AxesSubplot:>, <AxesSubplot:>], dtype=object)
```



# Line chart

In [8]:

```
data.plot.line()
```

Out[8]:

```
<AxesSubplot:>
```



## Bar chart

In [9]:

```
b=data[0:50]
```
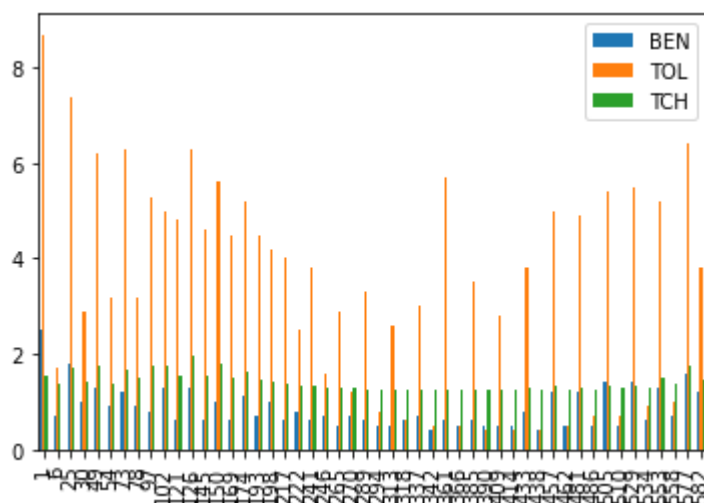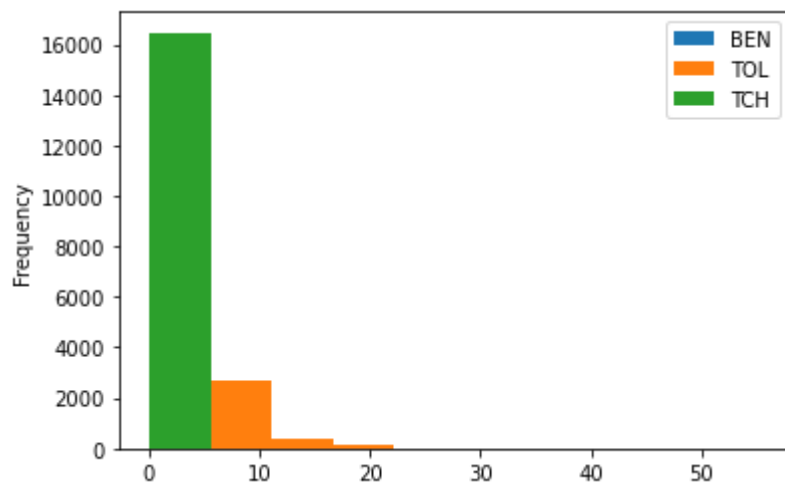
In [10]:

```
b.plot.bar()
```

Out[10]:

```
<AxesSubplot:>
```



## Histogram

In [11]:

```
data.plot.hist()
```
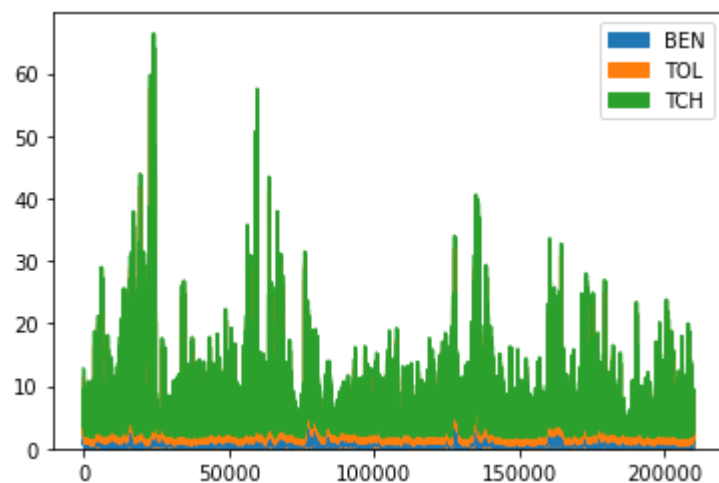
Out[11]:

```
<AxesSubplot:ylabel='Frequency'>
```



# Area chart

In [12]:

```
data.plot.area()
```

Out[12]:

```
<AxesSubplot:>
```

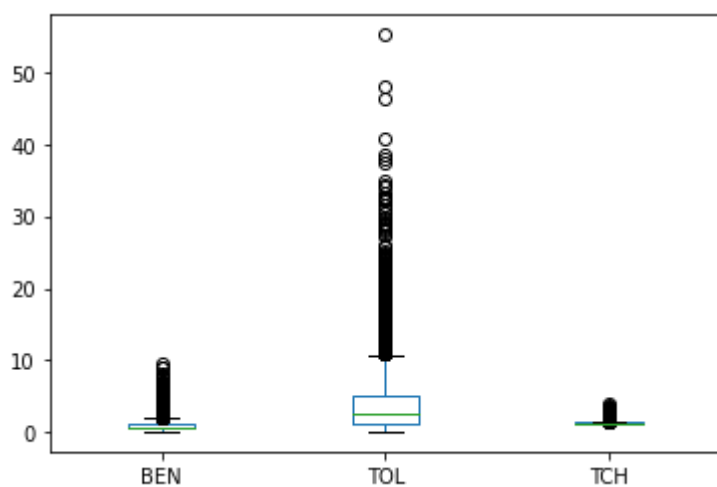

# Box chart

In [13]:

```python
data.plot.box()
```

Out[13]:

```
<AxesSubplot:>
```



# Pie chart

In [14]:
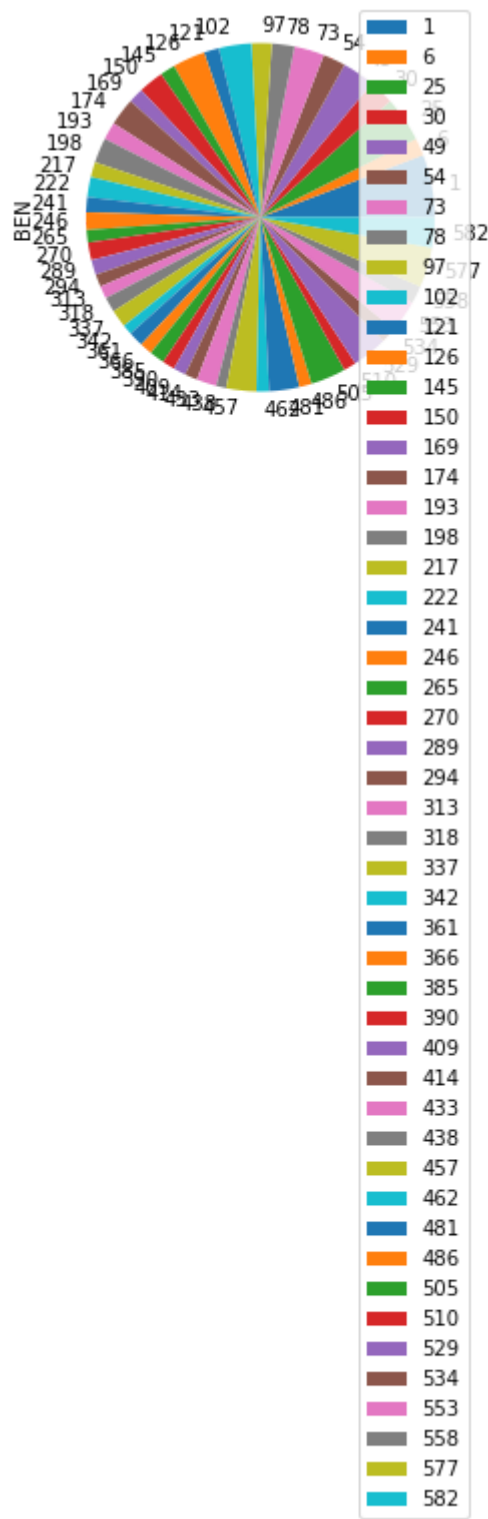
```
b.plot.pie(y='BEN' )
```

Out[14]:

<AxesSubplot:ylabel='BEN'>



# Scatter chart

In [15]:

```python
data.plot.scatter(x='BEN' ,y='TOL')
```
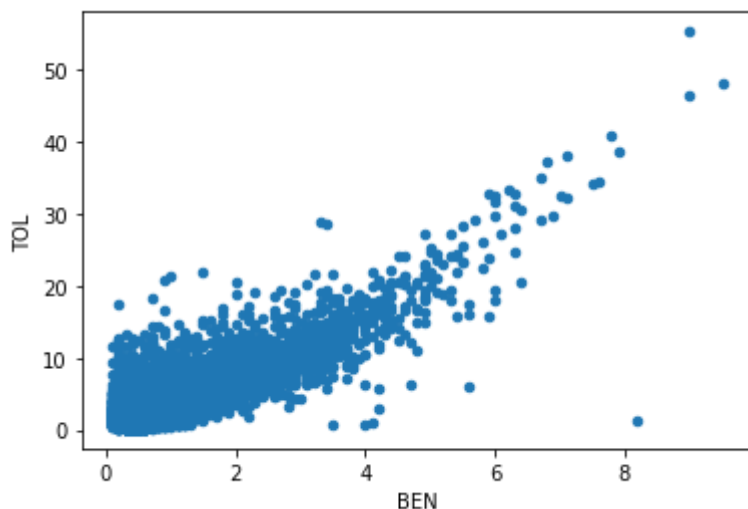
Out[15]:

```
<AxesSubplot:xlabel='BEN', ylabel='TOL'>
```



In [16]:

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 16460 entries, 1 to 209910
Data columns (total 14 columns):
 #   Column   Non-Null Count   Dtype
---  ------   --------------   -----
 0   date     16460 non-null   object
 1   BEN      16460 non-null   float64
 2   CO       16460 non-null   float64
 3   EBE      16460 non-null   float64
 4   NMHC     16460 non-null   float64
 5   NO       16460 non-null   float64
 6   NO_2     16460 non-null   float64
 7   O_3      16460 non-null   float64
 8   PM10     16460 non-null   float64
 9   PM25     16460 non-null   float64
 10  SO_2     16460 non-null   float64
 11  TCH      16460 non-null   float64
 12  TOL      16460 non-null   float64
 13  station  16460 non-null   int64
dtypes: float64(12), int64(1), object(1)
memory usage: 1.9+ MB
```

In [17]:

```python
df.describe()
```

Out[17]:

|        | BEN          | CO           | EBE          | NMHC         | NO           | NO_2         |
|--------|--------------|--------------|--------------|--------------|--------------|--------------|
| count  | 16460.000000 | 16460.000000 | 16460.000000 | 16460.000000 | 16460.000000 | 16460.000000 |
| mean   | 0.900680     | 0.277758     | 1.471871     | 0.167043     | 23.671810    | 44.583961    |
| std    | 0.768892     | 0.206143     | 1.051004     | 0.075068     | 44.362859    | 31.569185    |
| min    | 0.100000     | 0.100000     | 0.200000     | 0.010000     | 1.000000     | 1.000000     |
| 25%    | 0.500000     | 0.200000     | 0.800000     | 0.120000     | 2.000000     | 19.000000    |
| 50%    | 0.700000     | 0.200000     | 1.200000     | 0.160000     | 7.000000     | 40.000000    |
| 75%    | 1.100000     | 0.300000     | 1.700000     | 0.200000     | 25.000000    | 63.000000    |
| max    | 9.500000     | 3.200000     | 12.800000    | 0.840000     | 615.000000   | 289.000000   |

In [18]:

```python
df1=df[['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',
        'SO_2', 'TCH', 'TOL', 'station']]
```
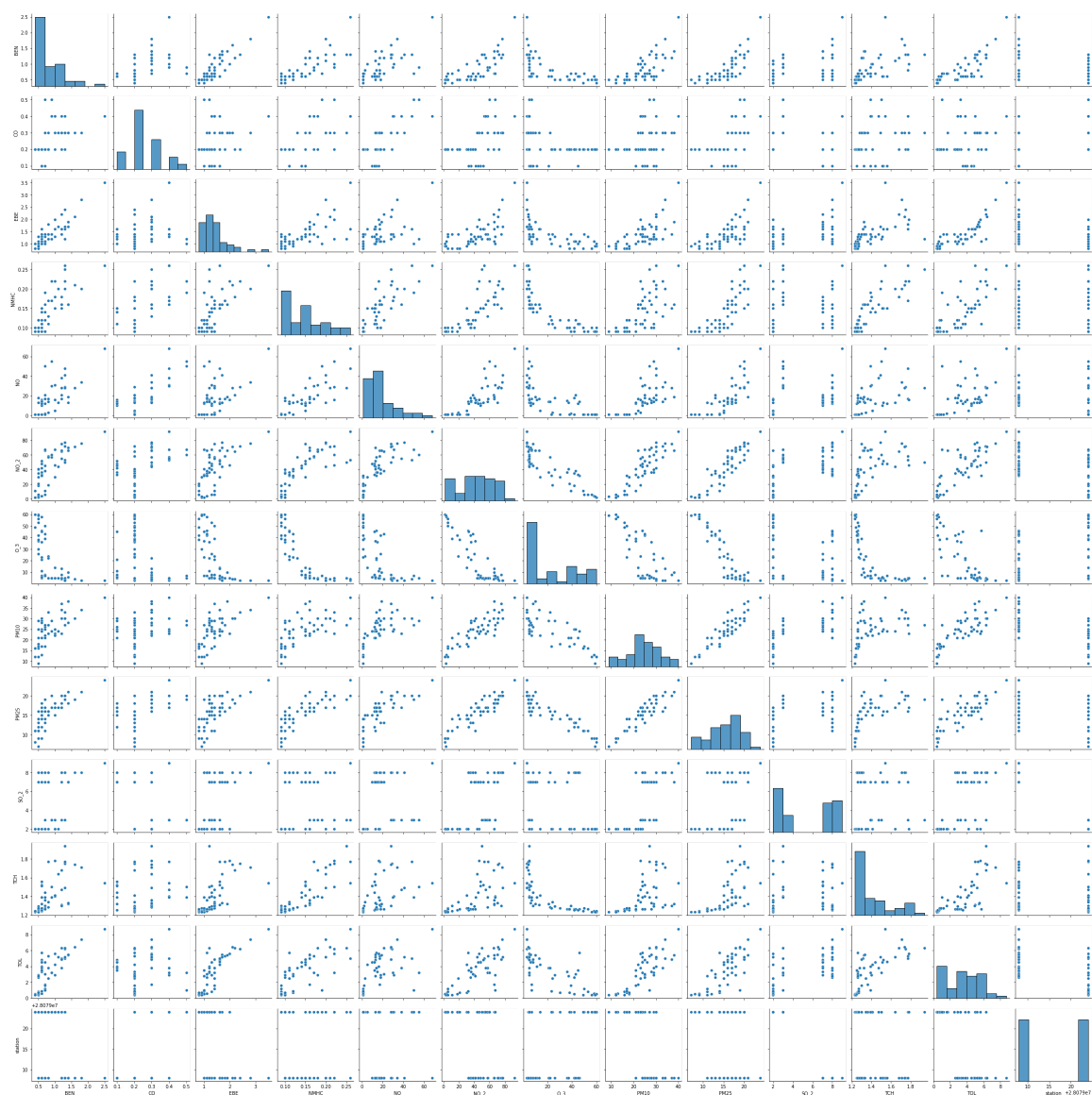
# EDA AND VISUALIZATION

In [19]:

```
sns.pairplot(df1[0:50])
```

Out[19]:

```
<seaborn.axisgrid.PairGrid at 0x1bb2a70d7f0>
```
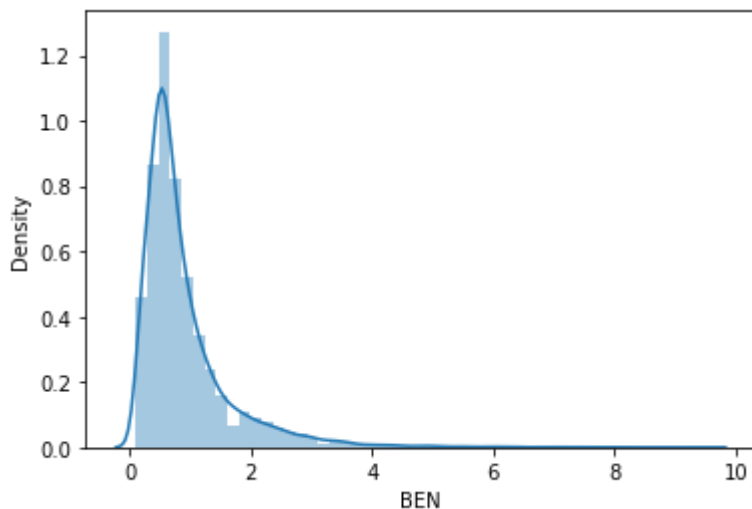
In [20]:

```python
sns.distplot(df1['BEN'])
```

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557:
FutureWarning: `distplot` is a deprecated function and will be removed in
a future version. Please adapt your code to use either `displot` (a figure
-level function with similar flexibility) or `histplot` (an axes-level fun
ction for histograms).
  warnings.warn(msg, FutureWarning)

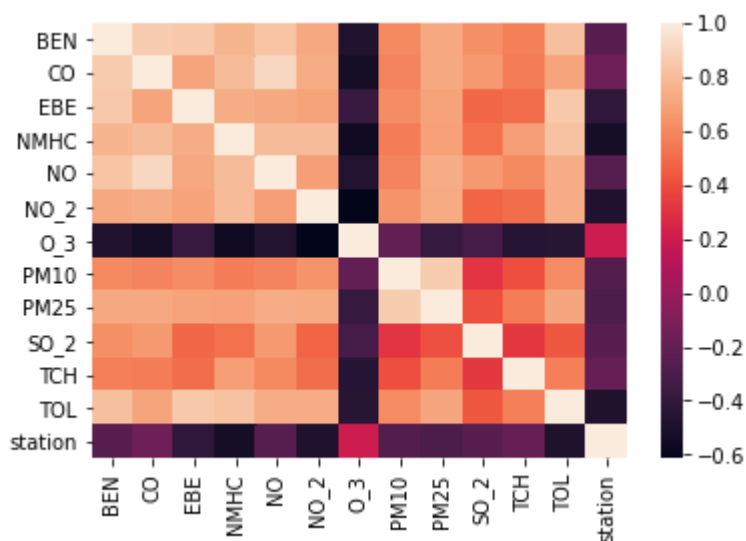Out[20]:

```
<AxesSubplot:xlabel='BEN', ylabel='Density'>
```



In [21]:

```python
sns.heatmap(df1.corr())
```

Out[21]:

```
<AxesSubplot:>
```



# TO TRAIN THE MODEL AND MODEL BULDING

In [22]:

```python
x=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',
       'SO_2', 'TCH', 'TOL']]
y=df['station']
```

In [23]:

```python
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

# Linear Regression

In [24]:

```python
from sklearn.linear_model import LinearRegression
lr=LinearRegression()
lr.fit(x_train,y_train)
```

Out[24]:

LinearRegression()

In [25]:

```python
lr.intercept_
```

Out[25]:

28079015.223056305

In [26]:

```python
coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```
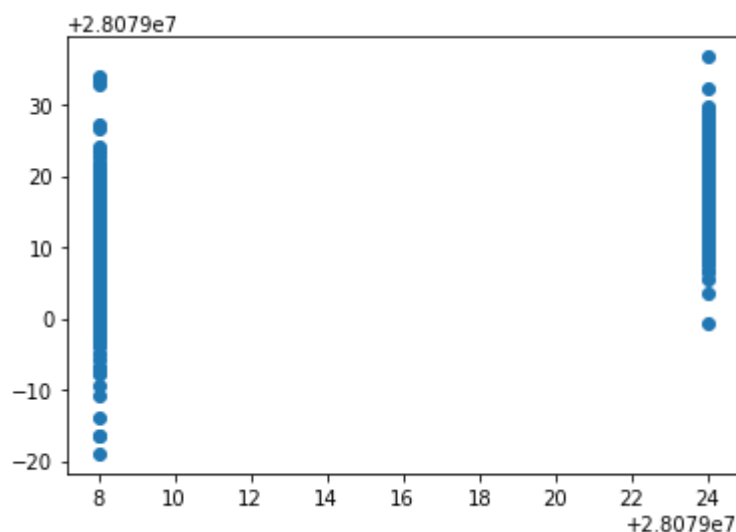
Out[26]:

|      | Co-efficient |
|------|--------------|
| BEN  | 3.805170 |
| CO   | 37.603729 |
| EBE  | -1.858668 |
| NMHC | -93.094409 |
| NO   | -0.036306 |
| NO_2 | -0.089968 |
| O_3  | -0.015230 |
| PM10 | 0.014540 |
| PM25 | -0.035677 |
| SO_2 | -0.474032 |
| TCH  | 11.074567 |
| TOL  | -0.415248 |

In [27]:

```
prediction =lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[27]:

```
<matplotlib.collections.PathCollection at 0x1bb36eccac0>
```



# ACCURACY

In [28]:

```
lr.score(x_test,y_test)
```

Out[28]:

```
0.6176242622614179
```

In [29]:

```
lr.score(x_train,y_train)
```

Out[29]:

```
0.6312376538837546
```

# Ridge and Lasso

In [30]:

```
from sklearn.linear_model import Ridge,Lasso
```

In [31]:

```
rr=Ridge(alpha=10)
rr.fit(x_train,y_train)
```

Out[31]:

```
Ridge(alpha=10)
```

# Accuracy(Ridge)

In [32]:

```
rr.score(x_test,y_test)
```

Out[32]:

0.580321441574621

In [33]:

```
rr.score(x_train,y_train)
```

Out[33]:

0.5979243288685172

In [34]:

```
la=Lasso(alpha=10)
la.fit(x_train,y_train)
```

Out[34]:

Lasso(alpha=10)

In [35]:

```
la.score(x_test,y_test)
```

Out[35]:

0.24603868232736525

# Accuracy(Lasso)

In [36]:

```
la.score(x_train,y_train)
```

Out[36]:

0.23241032291592945

# Accuracy(Elastic Net)

In [37]:

```
from sklearn.linear_model import ElasticNet
en=ElasticNet()
en.fit(x_train,y_train)
```

Out[37]:

ElasticNet()

In [38]:

```
en.coef_
```

Out[38]:

```
array([ 0.30437437,  0.        , -0.        , -0.        ,  0.05117634,
       -0.13281545, -0.04260736,  0.02915675,  0.09360197, -0.1808079 ,
        0.        , -1.00493552])
```

In [39]:

```
en.intercept_
```

Out[39]:

```
28079025.017157324
```

In [40]:

```
prediction=en.predict(x_test)
```

In [41]:

```
en.score(x_test,y_test)
```

Out[41]:

```
0.3502416888161278
```

# Evaluation Metrics

In [42]:

```
from sklearn import metrics
print(metrics.mean_absolute_error(y_test,prediction))
print(metrics.mean_squared_error(y_test,prediction))
print(np.sqrt(metrics.mean_squared_error(y_test,prediction)))
```

```
5.604138864118284
41.58278557254639
6.448471568716605
```

# Logistic Regression

In [43]:

```
from sklearn.linear_model import LogisticRegression
```

In [44]:

```
feature_matrix=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',
      'PM10', 'SO_2', 'TCH', 'TOL']]
target_vector=df['station']
```

In [45]:

```python
feature_matrix.shape
```

Out[45]:

```
(16460, 10)
```

In [46]:

```python
target_vector.shape
```

Out[46]:

```
(16460,)
```

In [47]:

```python
from sklearn.preprocessing import StandardScaler
```

In [48]:

```python
fs=StandardScaler().fit_transform(feature_matrix)
```

In [49]:

```python
logr=LogisticRegression(max_iter=10000)
logr.fit(fs,target_vector)
```

Out[49]:

```
LogisticRegression(max_iter=10000)
```

In [50]:

```python
observation=[[1,2,3,4,5,6,7,8,9,10]]
```

In [51]:

```python
prediction=logr.predict(observation)
print(prediction)
```

```
[28079008]
```

In [52]:

```python
logr.classes_
```

Out[52]:

```
array([28079008, 28079024], dtype=int64)
```

In [53]:

```python
logr.score(fs,target_vector)
```

Out[53]:

```
0.9237545565006076
```

In [54]:

```
logr.predict_proba(observation)[0][0]
```

Out[54]:

```
0.9999999999999966
```

In [55]:

```
logr.predict_proba(observation)
```

Out[55]:

```
array([[1.00000000e+00, 3.47334507e-15]])
```

# Random Forest

In [56]:

```
from sklearn.ensemble import RandomForestClassifier
```

In [57]:

```
rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

Out[57]:

```
RandomForestClassifier()
```

In [58]:

```
parameters={'max_depth':[1,2,3,4,5],
            'min_samples_leaf':[5,10,15,20,25],
            'n_estimators':[10,20,30,40,50]
}
```

In [59]:

```
from sklearn.model_selection import GridSearchCV
grid_search =GridSearchCV(estimator=rfc,param_grid=parameters,cv=2,scoring="accuracy")
grid_search.fit(x_train,y_train)
```

Out[59]:

```
GridSearchCV(cv=2, estimator=RandomForestClassifier(),
             param_grid={'max_depth': [1, 2, 3, 4, 5],
                         'min_samples_leaf': [5, 10, 15, 20, 25],
                         'n_estimators': [10, 20, 30, 40, 50]},
             scoring='accuracy')
```

In [60]:

```
grid_search.best_score_
```
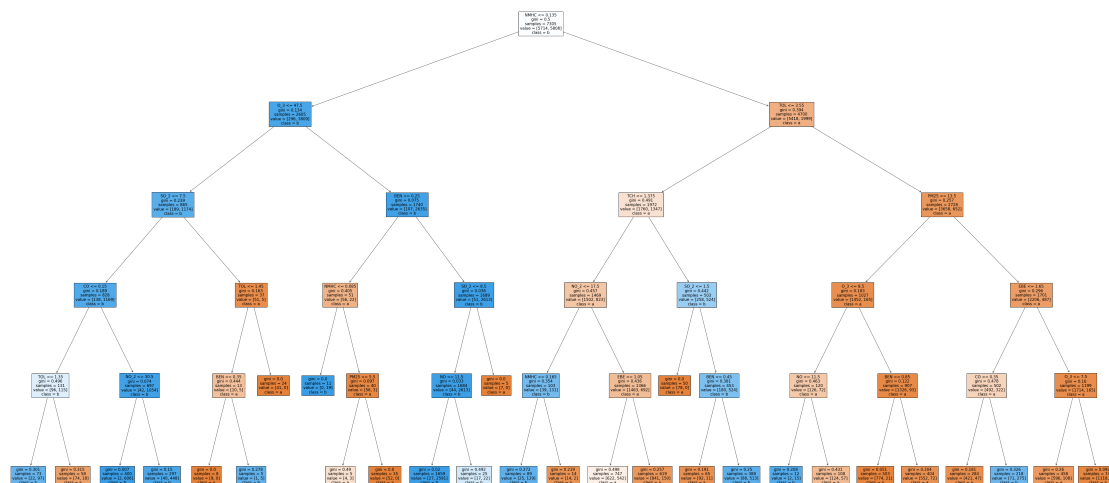
Out[60]:

```
0.9395938205172714
```

In [61]:

```python
rfc_best=grid_search.best_estimator_
```

In [62]:

```python
from sklearn.tree import plot_tree

plt.figure(figsize=(80,40))
plot_tree(rfc_best.estimators_[5],feature_names=x.columns,class_names=['a','b','c','d'],f
```

```
[596, 108]\nclass = a'),
 Text(4374.72, 181.19999999999982, 'gini = 0.092\nsamples = 741\nvalue =
[1118, 57]\nclass = a')]
```



# Conclusion

## Accuracy

*Linear Regression:0.6312376538837546*

*Ridge Regression:0.5979243288685172*

*Lasso Regression:0.23241032291592945*

*ElasticNet Regression:0.3502416888161278*

*Logistic Regression:0.9999999999999966*

*Random Forest:0.9395938205172714*

### Logistic Regression is suitable for this dataset