

20104016

DEENA

Importing Libraries

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
```

Importing Datasets

```
In [2]: df=pd.read_csv("madrid_2016.csv")
```

Out[2]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL
0	2016-11-01 01:00:00	NaN	0.7	NaN	NaN	153.0	77.0	NaN	NaN	NaN	7.0	NaN	NaN
1	2016-11-01 01:00:00	3.1	1.1	2.0	0.53	260.0	144.0	4.0	46.0	24.0	18.0	2.44	14.4
2	2016-11-01 01:00:00	5.9	NaN	7.5	NaN	297.0	139.0	NaN	NaN	NaN	NaN	NaN	26.0
3	2016-11-01 01:00:00	NaN	1.0	NaN	NaN	154.0	113.0	2.0	NaN	NaN	NaN	NaN	NaN
4	2016-11-01 01:00:00	NaN	NaN	NaN	NaN	275.0	127.0	2.0	NaN	NaN	18.0	NaN	NaN
...
209491	2016-07-01 00:00:00	NaN	0.2	NaN	NaN	2.0	29.0	73.0	NaN	NaN	NaN	NaN	NaN
209492	2016-07-01 00:00:00	NaN	0.3	NaN	NaN	1.0	29.0	NaN	36.0	NaN	5.0	NaN	NaN
209493	2016-07-01 00:00:00	NaN	NaN	NaN	NaN	1.0	19.0	71.0	NaN	NaN	NaN	NaN	NaN
209494	2016-07-01 00:00:00	NaN	NaN	NaN	NaN	6.0	17.0	85.0	NaN	NaN	NaN	NaN	NaN
209495	2016-07-01 00:00:00	NaN	NaN	NaN	NaN	2.0	46.0	61.0	34.0	NaN	NaN	NaN	NaN

209496 rows × 14 columns

Data Cleaning and Data Preprocessing

In [3]:

In [4]:

Out[4]: Index(['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',
'SO_2', 'TCH', 'TOL', 'station'],
dtype='object')

In [5]:

```
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 16932 entries, 1 to 209478  
Data columns (total 14 columns):  
#   Column      Non-Null Count  Dtype  
---  -  
0   date        16932 non-null  object  
1   BEN         16932 non-null  float64  
2   CO          16932 non-null  float64  
3   EBE         16932 non-null  float64  
4   NMHC        16932 non-null  float64  
5   NO          16932 non-null  float64  
6   NO_2        16932 non-null  float64  
7   O_3         16932 non-null  float64  
8   PM10        16932 non-null  float64  
9   PM25        16932 non-null  float64  
10  SO_2        16932 non-null  float64  
11  TCH         16932 non-null  float64  
12  TOL         16932 non-null  float64  
13  station     16932 non-null  int64  
dtypes: float64(12), int64(1), object(1)  
memory usage: 1.9+ MB
```

```
In [6]: data=df[['CO' , 'station']]
```

```
Out[6]:
```

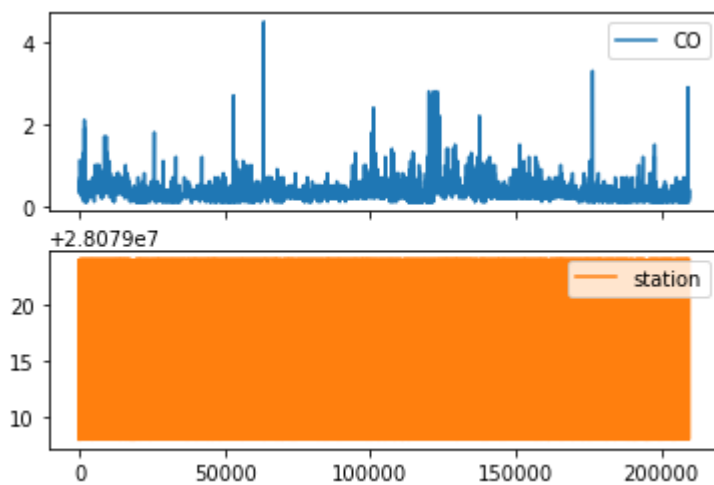
	CO	station
1	1.1	28079008
6	0.8	28079024
25	1.0	28079008
30	0.7	28079024
49	0.8	28079008
...
209430	0.2	28079024
209449	0.4	28079008
209454	0.2	28079024
209473	0.4	28079008
209478	0.2	28079024

16932 rows × 2 columns

Line chart

```
In [7]:
```

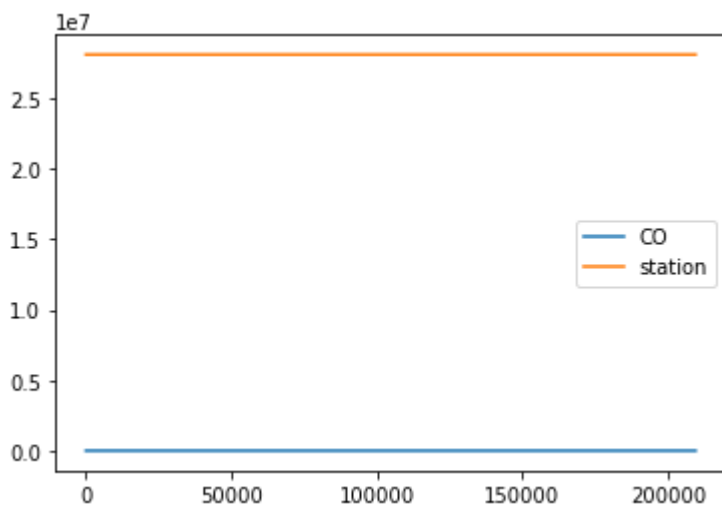
```
Out[7]: array([<AxesSubplot:>, <AxesSubplot:>], dtype=object)
```



Line chart

In [8]:

Out[8]: <AxesSubplot:>

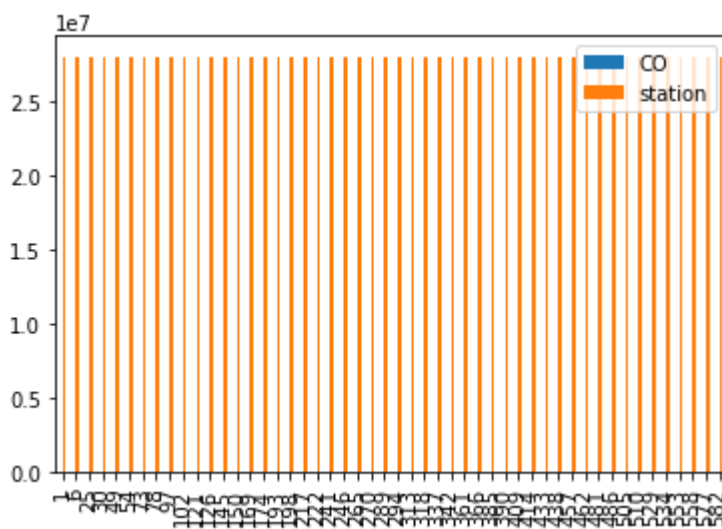


Bar chart

In [9]:

In [10]:

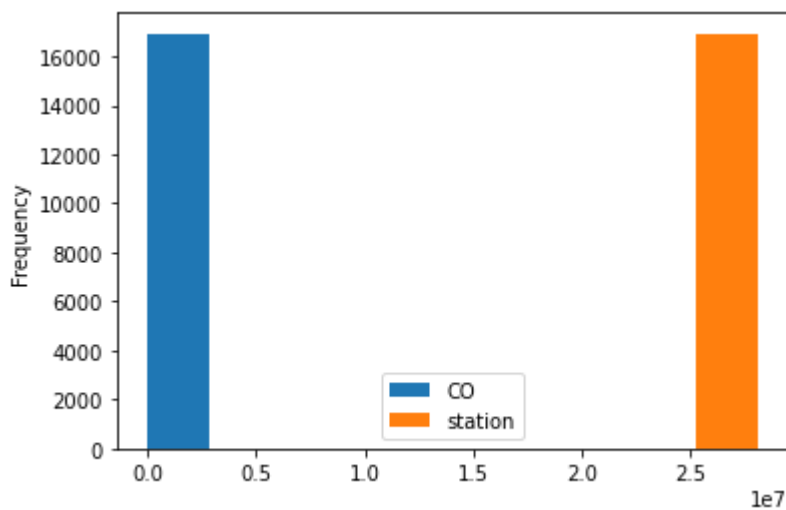
Out[10]: <AxesSubplot:>



Histogram

In [11]:

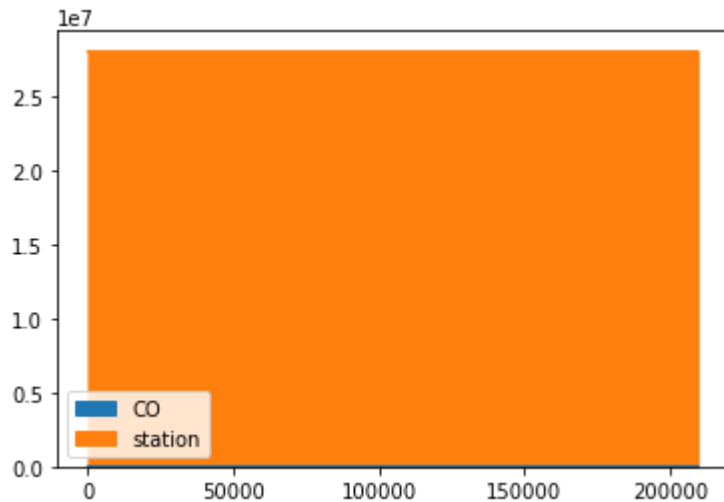
Out[11]: <AxesSubplot:ylabel='Frequency'>



Area chart

In [12]:

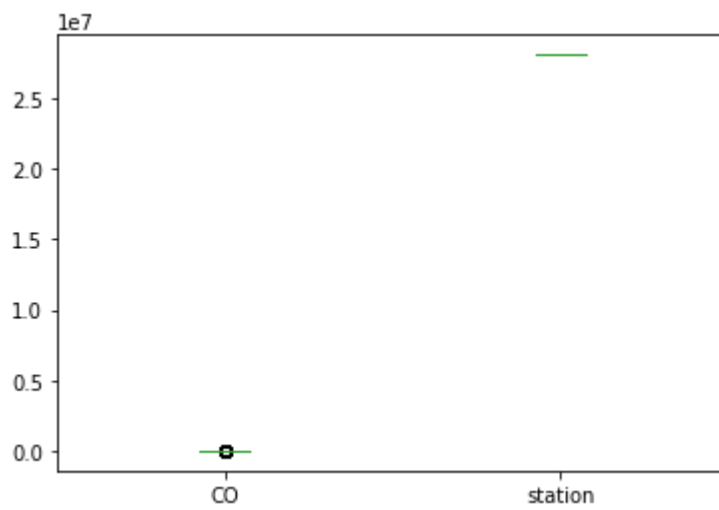
Out[12]: <AxesSubplot:>



Box chart

In [13]:

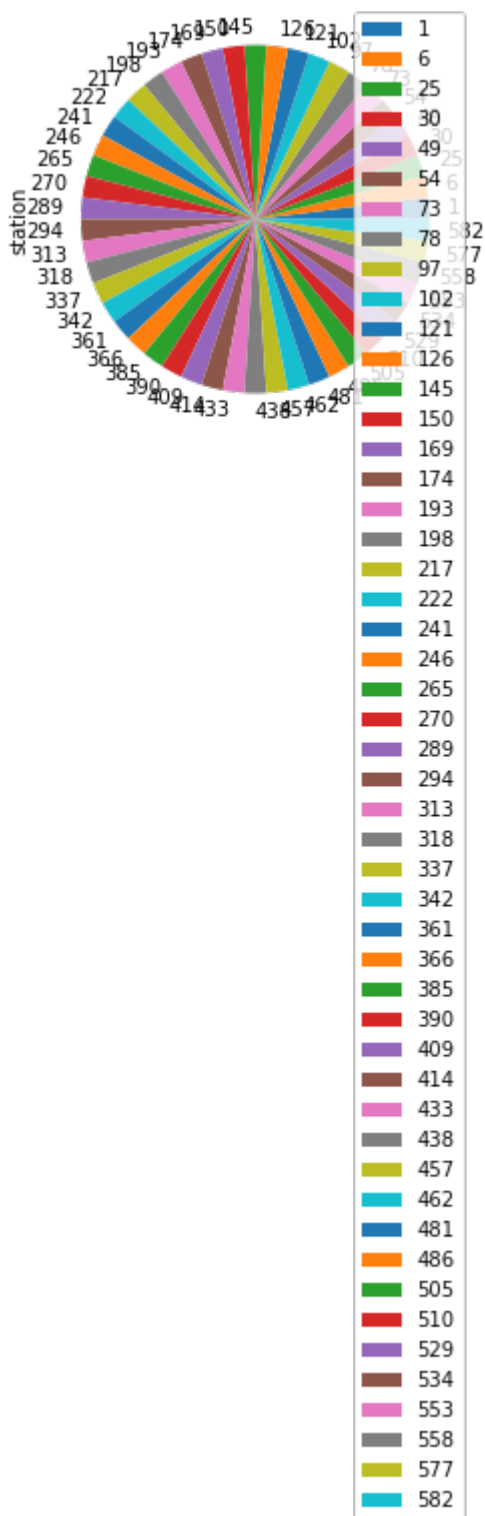
Out[13]: <AxesSubplot:>



Pie chart

In [14]:

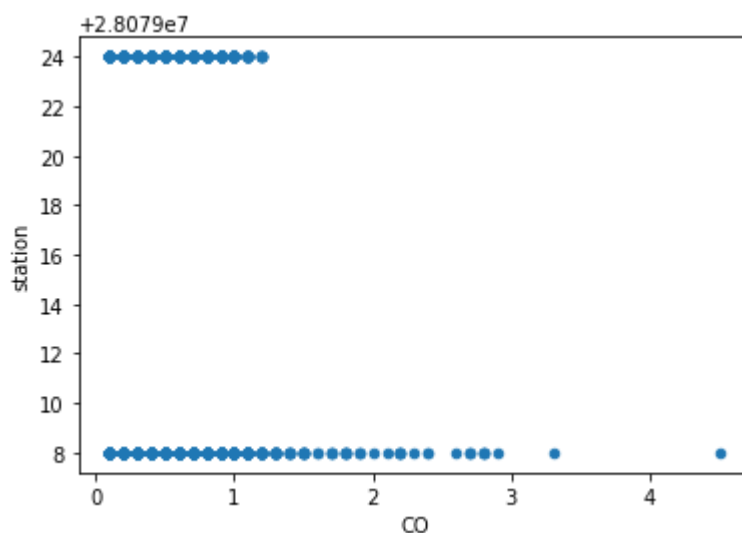
Out[14]: <AxesSubplot:ylabel='station'>



Scatter chart

In [15]:

Out[15]: <AxesSubplot:xlabel='CO', ylabel='station'>



In [16]:

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 16932 entries, 1 to 209478
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        16932 non-null  object
1   BEN         16932 non-null  float64
2   CO          16932 non-null  float64
3   EBE         16932 non-null  float64
4   NMHC        16932 non-null  float64
5   NO          16932 non-null  float64
6   NO_2        16932 non-null  float64
7   O_3         16932 non-null  float64
8   PM10        16932 non-null  float64
9   PM25        16932 non-null  float64
10  SO_2        16932 non-null  float64
11  TCH         16932 non-null  float64
12  TOL         16932 non-null  float64
13  station     16932 non-null  int64
```


	BEN	CO	EBE	NMHC	NO	NO_2
count	16932.000000	16932.000000	16932.000000	16932.000000	16932.000000	16932.000000
mean	0.537970	0.349941	0.298955	0.099913	20.815734	39.373376
std	0.599479	0.203807	0.450204	0.079850	40.986063	31.170307
min	0.100000	0.100000	0.100000	0.000000	1.000000	1.000000
25%	0.200000	0.200000	0.100000	0.050000	1.000000	14.000000
50%	0.400000	0.300000	0.200000	0.090000	7.000000	34.000000
75%	0.700000	0.400000	0.300000	0.120000	23.000000	58.000000
max	12.300000	4.500000	13.500000	2.210000	829.000000	319.000000

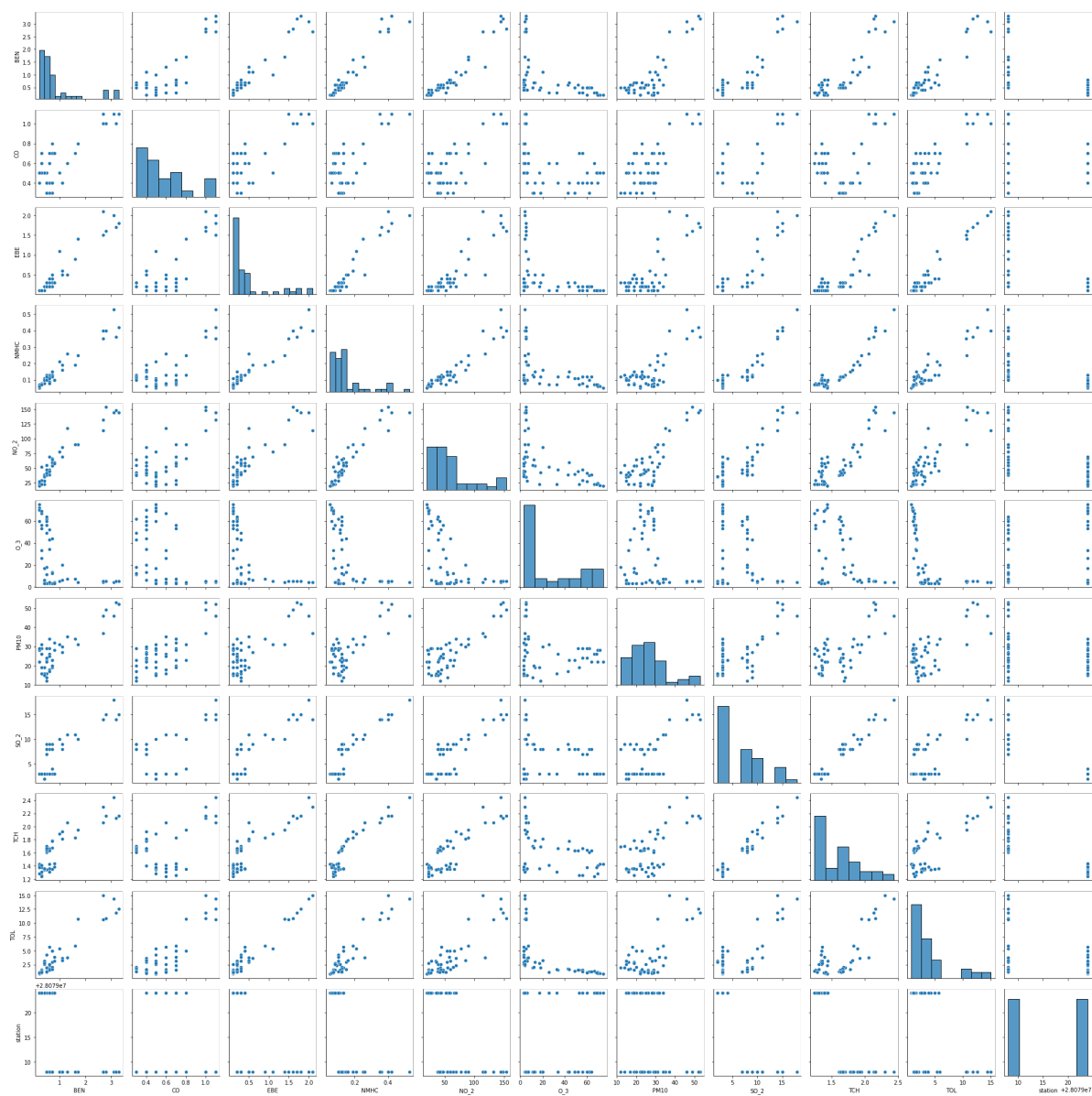
```
df1=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',  
        'PM10', 'SO_2', 'TCH', 'TCH', 'TCH', 'TCH']]
```

EDA AND VISUALIZATION

In [19]:

```
sns.pairplot(df)
```

Out[19]: <seaborn.axisgrid.PairGrid at 0x1ac855aaac0>

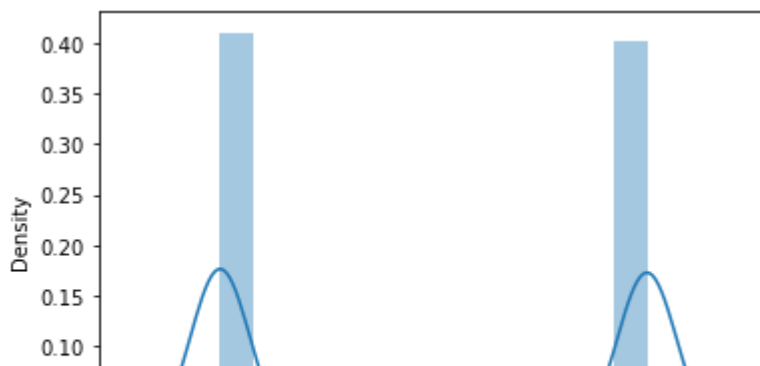


In [20]:

```
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 function for histograms).
```

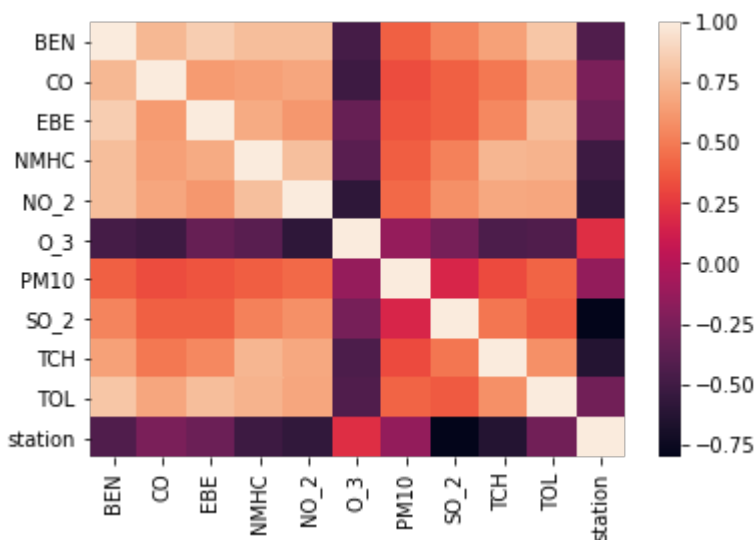
```
warnings.warn(msg, FutureWarning)
```

Out[20]: <AxesSubplot:xlabel='station', ylabel='Density'>



In [21]:

Out[21]: <AxesSubplot:>



TO TRAIN THE MODEL AND MODEL BUILDING

```
In [22]: x=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',
               'PM10', 'SO_2', 'TCH', 'TOL']]
```

```
In [23]: 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]: from sklearn.linear_model import LinearRegression  
lr=LinearRegression()
```

```
Out[24]: LinearRegression()
```

```
In [25]: lr.intercept_
```

```
Out[25]: 28079040.16666302
```

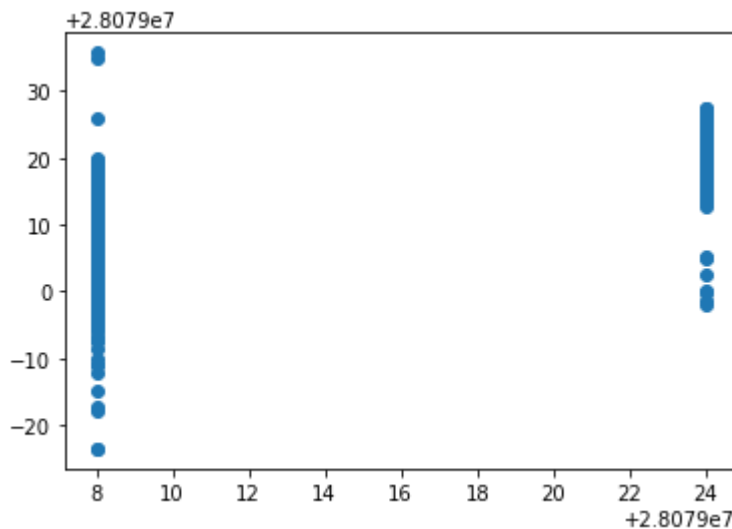
```
In [26]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
```

```
Out[26]:
```

	Co-efficient
BEN	1.317056
CO	6.426048
EBE	0.511071
NMHC	6.174209
NO_2	-0.069588
O_3	-0.028794
PM10	0.022648
SO_2	-0.839052
TCH	-13.088157
TOL	0.212470

```
In [27]: prediction =lr.predict(x_test)
```

```
Out[27]: <matplotlib.collections.PathCollection at 0x1ac8eacff40>
```



ACCURACY

In [28]:

Out[28]: 0.7977550894516356

In [29]:

Out[29]: 0.7932129126373898

Ridge and Lasso

In [30]:

In [31]: `rr=Ridge(alpha=10)`

Out[31]: Ridge(alpha=10)

Accuracy(Ridge)

In [32]:

Out[32]: 0.7979040290991042

In [33]:

Out[33]: 0.7929894292760988

In [34]: `la=Lasso(alpha=10)`

Out[34]: Lasso(alpha=10)

In [35]:

Out[35]: 0.6185114894305468

Accuracy(Lasso)

In [36]:

Out[36]: 0.6236707279152129

In [37]: `from sklearn.linear_model import ElasticNet
en=ElasticNet()
en.fit(x_train,y_train)`

Out[37]: ElasticNet()

In [38]:

```
Out[38]: array([ 0.          ,  0.          ,  0.          , -0.          , -0.08205541,
        -0.02674217,  0.02212625, -0.85133871, -0.          ,  0.25003088])
```

In [39]:

```
Out[39]: 28079025.946373824
```

In [40]: `prediction=en.predict(x_test)`

In [41]:

```
Out[41]: 0.6828599253726944
```

Evaluation Metrics

```
In [42]: from sklearn import metrics
print(metrics.mean_absolute_error(y_test,prediction))
print(metrics.mean_squared_error(y_test,prediction))
```

```
3.494153625629549
```

```
20.29530052319224
```

```
4.505030579606784
```

Logistic Regression

In [43]:

```
In [44]: feature_matrix=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',
        'PM10', 'SO_2', 'TCH', 'TOL']]
```

In [45]:

```
Out[45]: (16932, 10)
```

In [46]:

```
Out[46]: (16932,)
```

In [47]:

In [48]:

```
In [49]: logr=LogisticRegression(max_iter=10000)
```

```
Out[49]: LogisticRegression(max_iter=10000)
```

```
In [50]:
```

```
In [51]: prediction=logr.predict(observation)
[28079008]
```

```
In [52]:
Out[52]: array([28079008, 28079024], dtype=int64)
```

```
In [53]:
Out[53]: 0.9923812898653437
```

```
In [54]:
Out[54]: 1.0
```

```
In [55]:
Out[55]: array([[1.0000000e+00, 1.6336121e-46]])
```

Random Forest

```
In [56]:
```

```
In [57]: rfc=RandomForestClassifier()
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="ac
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]:
Out[60]: 0.9946000674991562
```

```
In [61]:
```

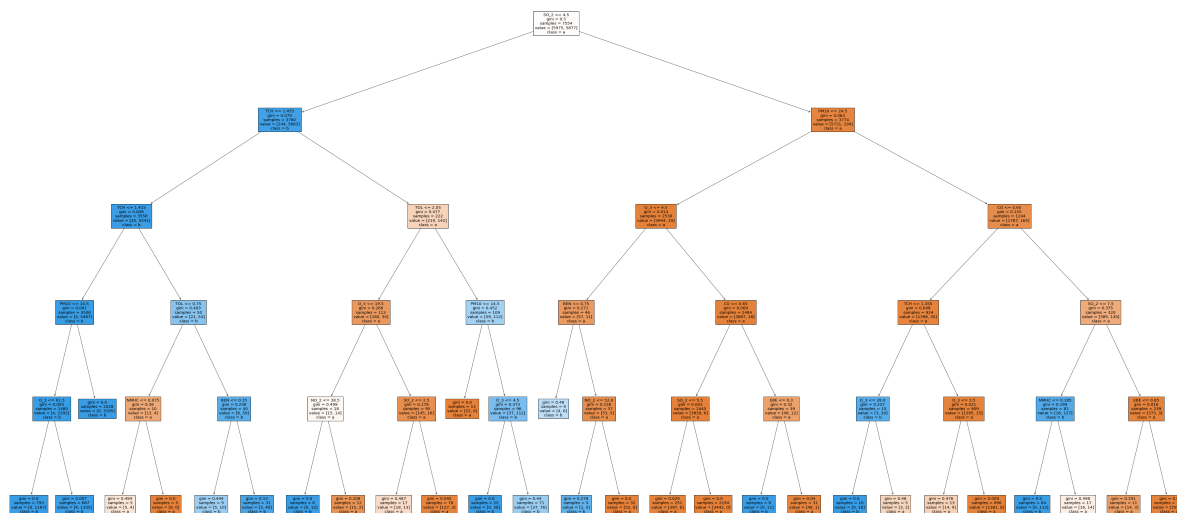
In [62]: `from sklearn.tree import plot_tree`

```
plt.figure(figsize=(80,40))
```

Out[62]: [Text(2071.0384615384614, 1993.2, 'SO₂ <= 4.5\ngini = 0.5\nsamples = 7554\nvalue = [5975, 5877]\nclass = a'),
Text(1030.1538461538462, 1630.8000000000002, 'TCH <= 1.455\ngini = 0.079\nsamples = 3780\nvalue = [244, 5683]\nclass = b'),
Text(472.15384615384613, 1268.4, 'TCH <= 1.415\ngini = 0.009\nsamples = 3558\nvalue = [25, 5541]\nclass = b'),
Text(257.53846153846155, 906.0, 'PM10 <= 10.5\ngini = 0.001\nsamples = 3508\nvalue = [4, 5487]\nclass = b'),
Text(171.69230769230768, 543.5999999999999, 'O₃ <= 61.5\ngini = 0.003\nsamples = 1480\nvalue = [4, 2292]\nclass = b'),
Text(85.84615384615384, 181.19999999999982, 'gini = 0.0\nsamples = 793\nvalue = [0, 1187]\nclass = b'),
Text(257.53846153846155, 181.19999999999982, 'gini = 0.007\nsamples = 687\nvalue = [4, 1105]\nclass = b'),
Text(343.38461538461536, 543.5999999999999, 'gini = 0.0\nsamples = 2028\nvalue = [0, 3195]\nclass = b'),
Text(686.7692307692307, 906.0, 'TOL <= 0.75\ngini = 0.403\nsamples = 50\nvalue = [21, 54]\nclass = b'),
Text(515.0769230769231, 543.5999999999999, 'NMHC <= 0.075\ngini = 0.36\nsamples = 10\nvalue = [13, 4]\nclass = a'),
Text(429.23076923076917, 181.19999999999982, 'gini = 0.494\nsamples = 5\nvalue = [5, 4]\nclass = a'),
Text(600.9230769230769, 181.19999999999982, 'gini = 0.0\nsamples = 5\nvalue = [8, 0]\nclass = a'),
Text(858.4615384615383, 543.5999999999999, 'BEN <= 0.35\ngini = 0.238\nsamples = 40\nvalue = [8, 50]\nclass = b'),
Text(772.6153846153845, 181.19999999999982, 'gini = 0.444\nsamples = 9\nvalue = [5, 10]\nclass = b'),
Text(944.3076923076923, 181.19999999999982, 'gini = 0.13\nsamples = 31\nvalue = [3, 40]\nclass = b'),
Text(1588.153846153846, 1268.4, 'TOL <= 2.05\ngini = 0.477\nsamples = 222\nvalue = [219, 142]\nclass = a'),
Text(1373.5384615384614, 906.0, 'O₃ <= 19.5\ngini = 0.266\nsamples = 113\nvalue = [160, 30]\nclass = a'),
Text(1201.8461538461538, 543.5999999999999, 'NO₂ <= 38.5\ngini = 0.499\nsamples = 18\nvalue = [15, 14]\nclass = a'),
Text(1116.0, 181.19999999999982, 'gini = 0.0\nsamples = 6\nvalue = [0, 12]\nclass = b'),
Text(1287.6923076923076, 181.19999999999982, 'gini = 0.208\nsamples = 12\nvalue = [15, 2]\nclass = a'),
Text(1545.230769230769, 543.5999999999999, 'SO₂ <= 2.5\ngini = 0.179\nsamples = 95\nvalue = [145, 16]\nclass = a'),
Text(1459.3846153846152, 181.19999999999982, 'gini = 0.487\nsamples = 17\nvalue = [18, 13]\nclass = a'),
Text(1631.0769230769229, 181.19999999999982, 'gini = 0.045\nsamples = 78\nvalue = [127, 3]\nclass = a'),
Text(1802.7692307692307, 906.0, 'PM10 <= 14.5\ngini = 0.452\nsamples = 109\nvalue = [59, 112]\nclass = b'),
Text(1716.9230769230767, 543.5999999999999, 'gini = 0.0\nsamples = 13\nvalue = [22, 0]\nclass = a'),
Text(1888.6153846153845, 543.5999999999999, 'O₃ <= 4.5\ngini = 0.373\nsamples = 96\nvalue = [37, 112]\nclass = b'),


```
Text(1802.7692307692307, 181.19999999999982, 'gini = 0.0\nsamples = 25\nvalue = [0, 36]\nclass = b'),
Text(1974.4615384615383, 181.19999999999982, 'gini = 0.44\nsamples = 71\nvalue = [37, 76]\nclass = b'),
Text(3111.9230769230767, 1630.8000000000002, 'PM10 <= 24.5\ngini = 0.063\nsamples = 3774\nvalue = [5731, 194]\nclass = a'),
Text(2446.6153846153843, 1268.4, 'O_3 <= 4.5\ngini = 0.014\nsamples = 2530\nvalue = [3944, 29]\nclass = a'),
Text(2146.153846153846, 906.0, 'BEN <= 0.75\ngini = 0.271\nsamples = 46\nvalue = [57, 11]\nclass = a'),
Text(2060.3076923076924, 543.5999999999999, 'gini = 0.48\nsamples = 9\nvalue = [4, 6]\nclass = b'),
Text(2232.0, 543.5999999999999, 'NO_2 <= 52.0\ngini = 0.158\nsamples = 37\nvalue = [53, 5]\nclass = a'),
Text(2146.153846153846, 181.19999999999982, 'gini = 0.278\nsamples = 5\nvalue = [1, 5]\nclass = b'),
Text(2317.846153846154, 181.19999999999982, 'gini = 0.0\nsamples = 32\nvalue = [52, 0]\nclass = a'),
Text(2747.076923076923, 906.0, 'CO <= 0.65\ngini = 0.009\nsamples = 2484\nvalue = [3887, 18]\nclass = a'),
Text(2575.3846153846152, 543.5999999999999, 'SO_2 <= 5.5\ngini = 0.003\nsamples = 2445\nvalue = [3839, 6]\nclass = a'),
Text(2489.5384615384614, 181.19999999999982, 'gini = 0.029\nsamples = 251\nvalue = [397, 6]\nclass = a'),
Text(2661.230769230769, 181.19999999999982, 'gini = 0.0\nsamples = 2194\nvalue = [3442, 0]\nclass = a'),
Text(2918.7692307692305, 543.5999999999999, 'EBE <= 0.3\ngini = 0.32\nsamples = 39\nvalue = [48, 12]\nclass = a'),
Text(2832.9230769230767, 181.19999999999982, 'gini = 0.0\nsamples = 8\nvalue = [0, 11]\nclass = b'),
Text(3004.6153846153843, 181.19999999999982, 'gini = 0.04\nsamples = 31\nvalue = [48, 1]\nclass = a'),
Text(3777.230769230769, 1268.4, 'CO <= 0.65\ngini = 0.155\nsamples = 1244\nvalue = [1787, 165]\nclass = a'),
Text(3433.8461538461534, 906.0, 'TCH <= 1.355\ngini = 0.048\nsamples = 924\nvalue = [1398, 35]\nclass = a'),
Text(3262.1538461538457, 543.5999999999999, 'O_3 <= 28.0\ngini = 0.227\nsamples = 15\nvalue = [3, 20]\nclass = b'),
Text(3176.307692307692, 181.19999999999982, 'gini = 0.0\nsamples = 10\nvalue = [0, 18]\nclass = b'),
Text(3347.9999999999995, 181.19999999999982, 'gini = 0.48\nsamples = 5\nvalue = [3, 2]\nclass = a'),
Text(3605.5384615384614, 543.5999999999999, 'O_3 <= 3.5\ngini = 0.021\nsamples = 909\nvalue = [1395, 15]\nclass = a'),
Text(3519.6923076923076, 181.19999999999982, 'gini = 0.476\nsamples = 13\nvalue = [14, 9]\nclass = a'),
Text(3691.3846153846152, 181.19999999999982, 'gini = 0.009\nsamples = 896\nvalue = [1381, 6]\nclass = a'),
Text(4120.615384615385, 906.0, 'SO_2 <= 7.5\ngini = 0.375\nsamples = 320\nvalue = [389, 130]\nclass = a'),
Text(3948.9230769230767, 543.5999999999999, 'NMHC <= 0.185\ngini = 0.199\nsamples = 81\nvalue = [16, 127]\nclass = b'),
Text(3863.076923076923, 181.19999999999982, 'gini = 0.0\nsamples = 64\nvalue = [0, 113]\nclass = b'),
Text(4034.7692307692305, 181.19999999999982, 'gini = 0.498\nsamples = 17\nvalue = [16, 14]\nclass = a'),
```

```
Text(4292.307692307692, 543.5999999999999, 'EBE <= 0.65\ngini = 0.016\nsamples = 239\nvalue = [373, 3]\nnclass = a'),
Text(4206.461538461538, 181.19999999999982, 'gini = 0.291\nsamples = 11\nvalue = [14, 3]\nnclass = a'),
Text(4378.153846153846, 181.19999999999982, 'gini = 0.0\nsamples = 228\nvalue = [359, 0]\nnclass = a')]
```



Conclusion

Accuracy

Linear Regression :0.7932129126373898

Ridge Regression :0.6185114894305468

Lasso Regression :0.6236707279152129

ElasticNet Regression : 0.6828599253726944

Logistic Regression : 0.9923812898653437

Random Forest :0.9946000674991562

Random Forest is suitable for this dataset

