

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_2005.csv")
```

Out[2]:

	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3	PM1
0	2005-11-01 01:00:00	NaN	0.77	NaN	NaN	NaN	57.130001	128.699997	NaN	14.720000	14.9
1	2005-11-01 01:00:00	1.52	0.65	1.49	4.57	0.25	86.559998	181.699997	1.27	11.680000	30.9
2	2005-11-01 01:00:00	NaN	0.40	NaN	NaN	NaN	46.119999	53.000000	NaN	30.469999	14.6
3	2005-11-01 01:00:00	NaN	0.42	NaN	NaN	NaN	37.220001	52.009998	NaN	21.379999	15.1
4	2005-11-01 01:00:00	NaN	0.57	NaN	NaN	NaN	32.160000	36.680000	NaN	33.410000	5.0
...
236995	2006-01-01 00:00:00	1.08	0.36	1.01	NaN	0.11	21.990000	23.610001	NaN	43.349998	5.0
236996	2006-01-01 00:00:00	0.39	0.54	1.00	1.00	0.11	2.200000	4.220000	1.00	69.639999	4.9
236997	2006-01-01 00:00:00	0.19	NaN	0.26	NaN	0.08	26.730000	30.809999	NaN	43.840000	4.3
236998	2006-01-01 00:00:00	0.14	NaN	1.00	NaN	0.06	13.770000	17.770000	NaN	NaN	5.0
236999	2006-01-01 00:00:00	0.50	0.40	0.73	1.84	0.13	20.940001	26.950001	1.49	48.259998	5.6

237000 rows × 17 columns

Data Cleaning and Data Preprocessing

In [3]:

In [4]:

Out[4]: Index(['date', 'BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
'PM10', 'PM25', 'PXY', 'SO_2', 'TCH', 'TOL', 'station'],
dtype='object')

In [5]:

```
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 20070 entries, 5 to 236999  
Data columns (total 17 columns):  
#   Column      Non-Null Count  Dtype  
---  -  
0   date        20070 non-null  object  
1   BEN         20070 non-null  float64  
2   CO          20070 non-null  float64  
3   EBE         20070 non-null  float64  
4   MXY         20070 non-null  float64  
5   NMHC        20070 non-null  float64  
6   NO_2        20070 non-null  float64  
7   NOx         20070 non-null  float64  
8   OXY         20070 non-null  float64  
9   O_3         20070 non-null  float64  
10  PM10        20070 non-null  float64  
11  PM25        20070 non-null  float64  
12  PXY         20070 non-null  float64  
13  SO_2        20070 non-null  float64  
14  TCH         20070 non-null  float64  
15  TOL         20070 non-null  float64  
16  station     20070 non-null  int64  
dtypes: float64(15), int64(1), object(1)  
memory usage: 2.8+ MB
```

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

```
Out[6]:
```

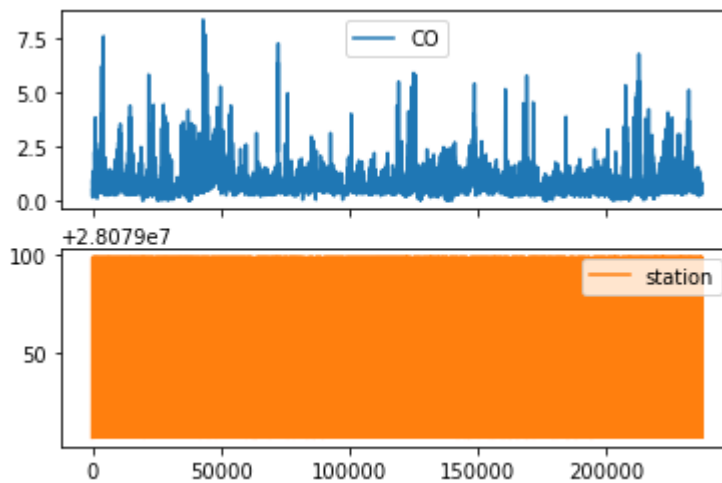
	CO	station
5	0.88	28079006
22	0.22	28079024
25	0.49	28079099
31	0.84	28079006
48	0.20	28079024
...
236970	0.39	28079024
236973	0.45	28079099
236979	0.38	28079006
236996	0.54	28079024
236999	0.40	28079099

20070 rows × 2 columns

Line chart

```
In [7]:
```

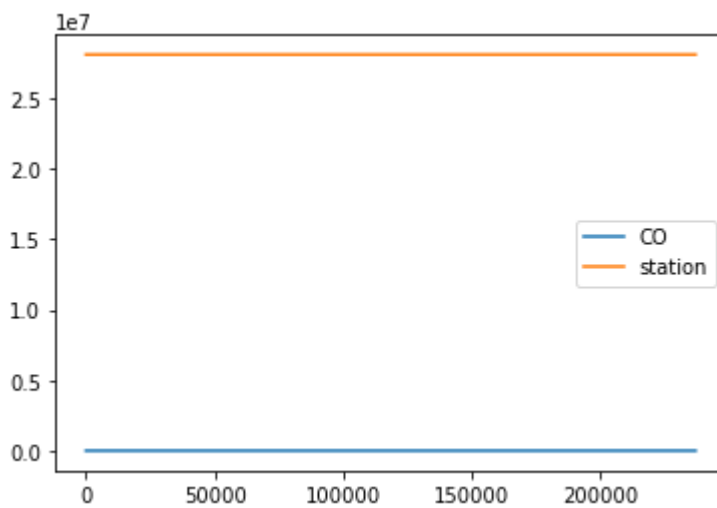
```
Out[7]: array([<AxesSubplot:>, <AxesSubplot:>], dtype=object)
```



Line chart

In [8]:

Out[8]: <AxesSubplot:>

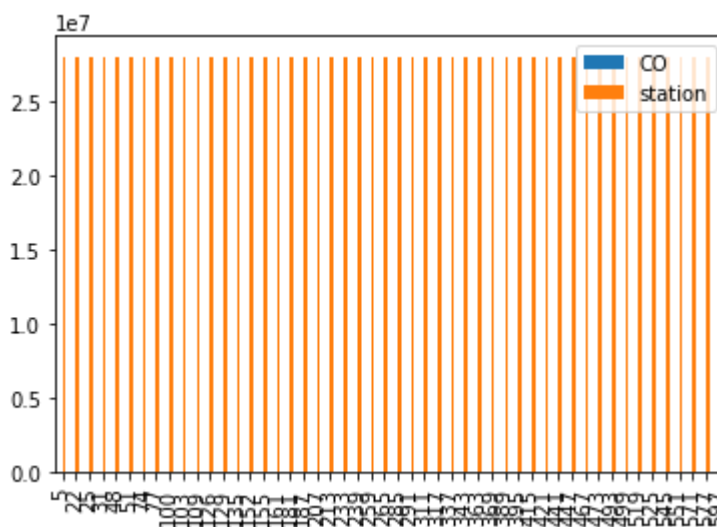


Bar chart

In [9]:

In [10]:

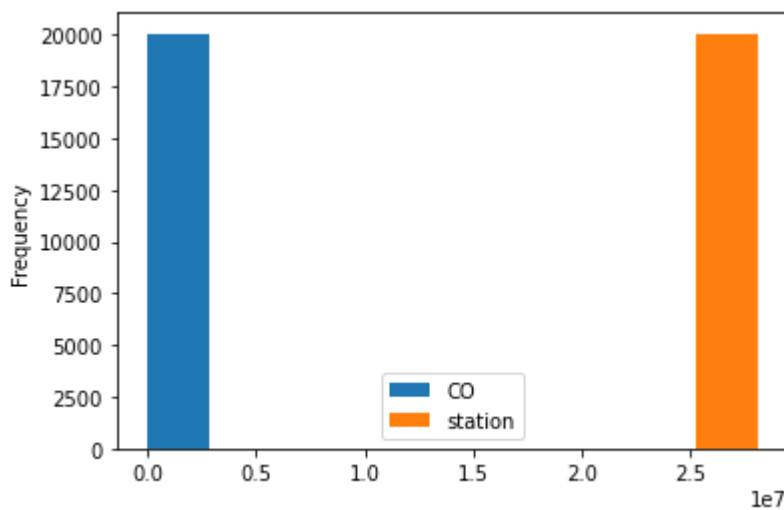
Out[10]: <AxesSubplot:>



Histogram

In [11]:

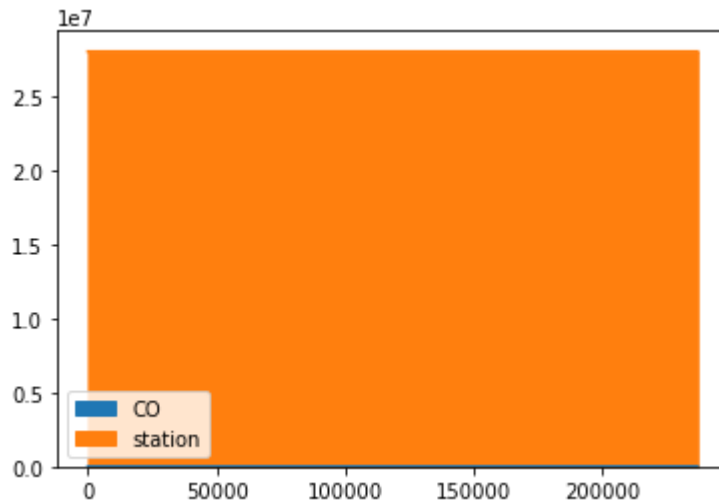
Out[11]: <AxesSubplot:ylabel='Frequency'>



Area chart

In [12]:

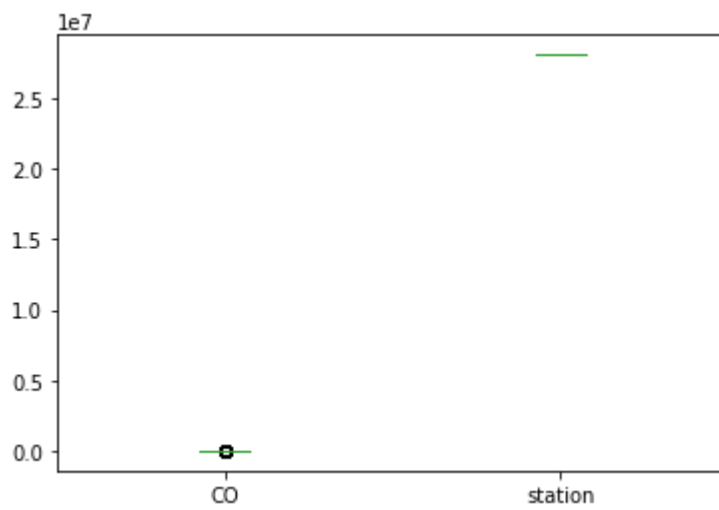
Out[12]: <AxesSubplot:>



Box chart

In [13]:

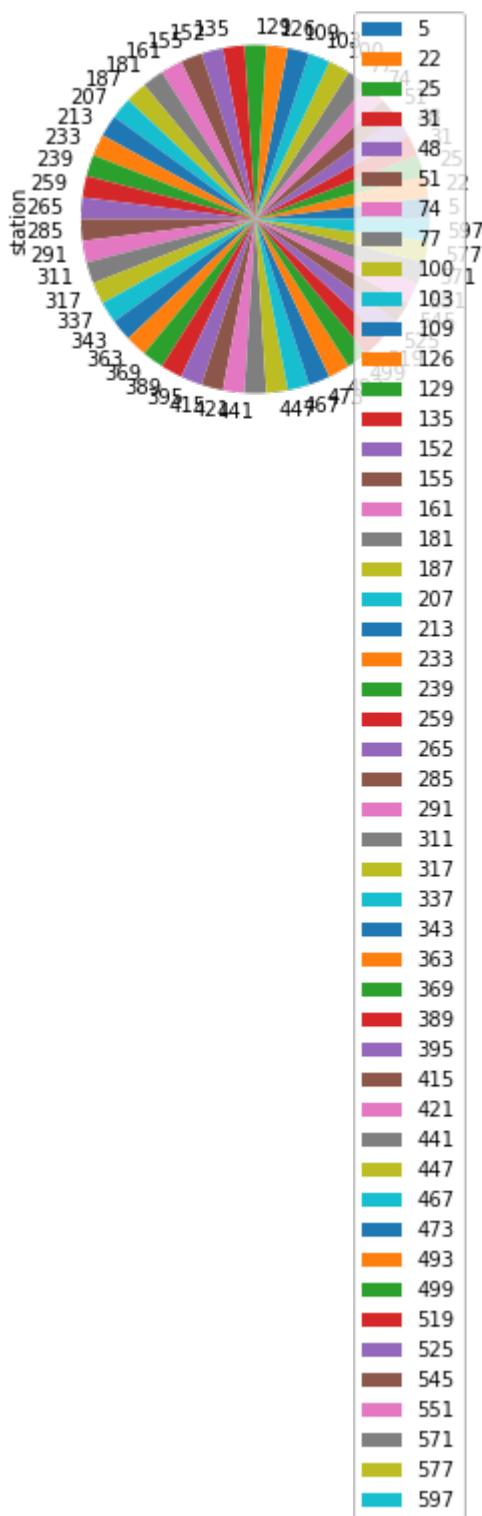
Out[13]: <AxesSubplot:>



Pie chart

In [14]:

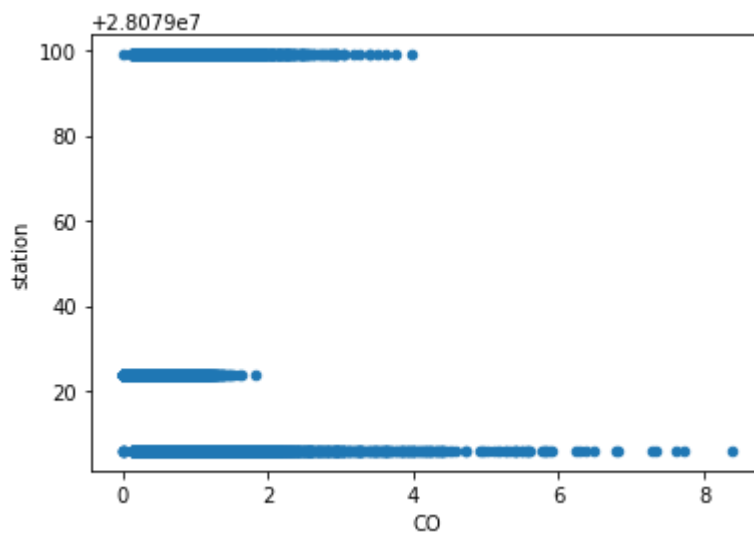
Out[14]: <AxesSubplot:ylabel='station'>



Scatter chart

In [15]:

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



In [16]:

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 20070 entries, 5 to 236999
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        20070 non-null  object
1   BEN         20070 non-null  float64
2   CO          20070 non-null  float64
3   EBE         20070 non-null  float64
4   MXY         20070 non-null  float64
5   NMHC        20070 non-null  float64
6   NO_2        20070 non-null  float64
7   NOx         20070 non-null  float64
8   OXY         20070 non-null  float64
9   O_3         20070 non-null  float64
10  PM10        20070 non-null  float64
11  PM25        20070 non-null  float64
12  PXY         20070 non-null  float64
13  SO_2        20070 non-null  float64
14  TCU         20070 non-null  float64
```


	BEN	CO	EBE	MXV	NMHC	NO_2	
count	20070.000000	20070.000000	20070.000000	20070.000000	20070.000000	20070.000000	20070.000000
mean	1.923656	0.720657	2.345423	5.457855	0.179282	66.226924	1.923656
std	2.019061	0.549723	2.379219	5.495147	0.152783	40.568197	2.019061
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.690000	0.400000	0.950000	1.930000	0.090000	36.602499	0.690000
50%	1.260000	0.580000	1.480000	3.800000	0.150000	60.525000	1.260000
75%	2.510000	0.880000	2.950000	7.210000	0.220000	89.317499	2.510000
max	26.570000	8.380000	29.870001	71.050003	1.880000	419.500000	26.570000

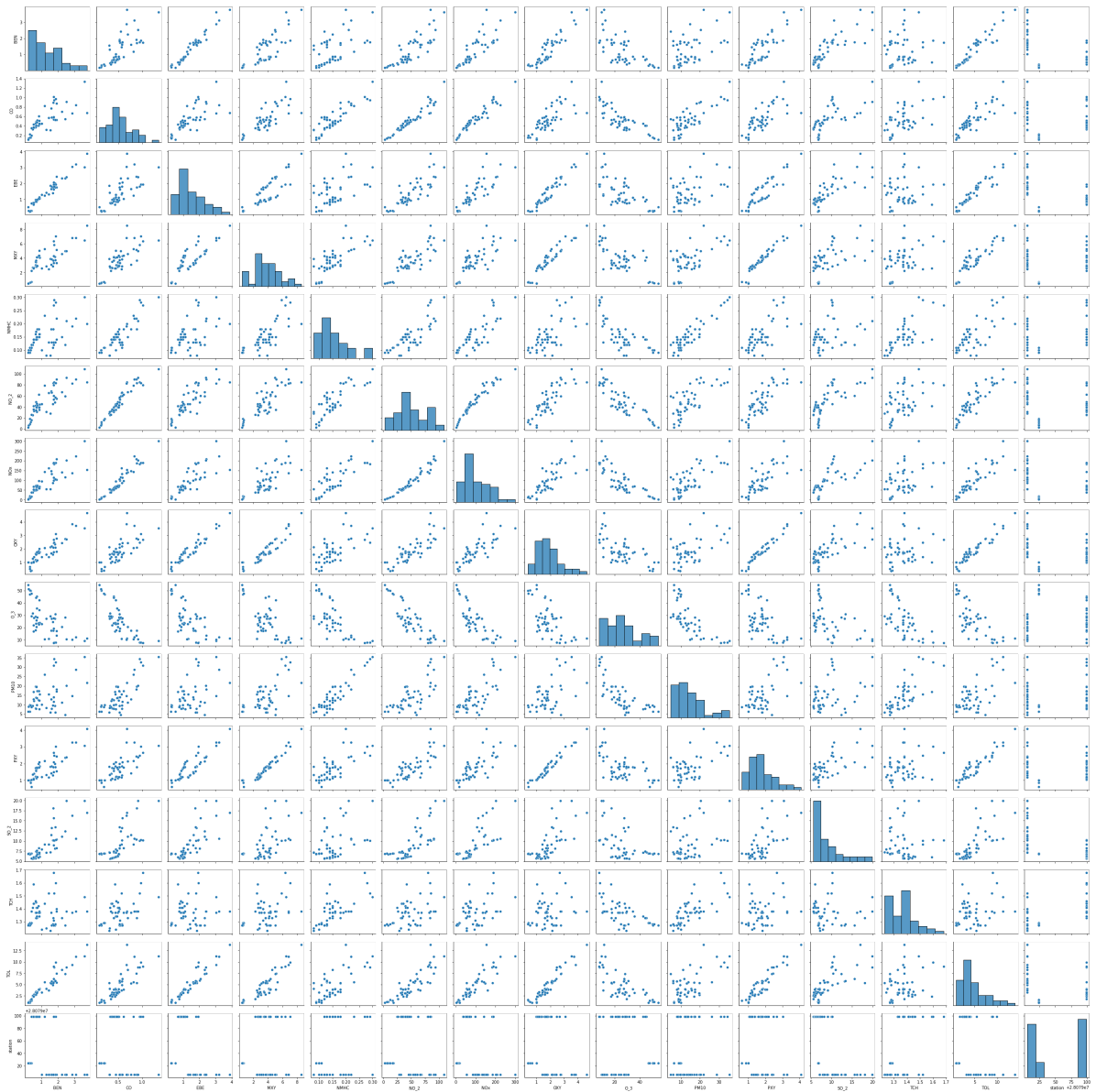
```
df1=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',  
        'PM10', 'PM25', 'SO_2', 'TSP', 'TCF', 'VOC', 'WINDSPEED']]
```

EDA AND VISUALIZATION

In [19]:

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

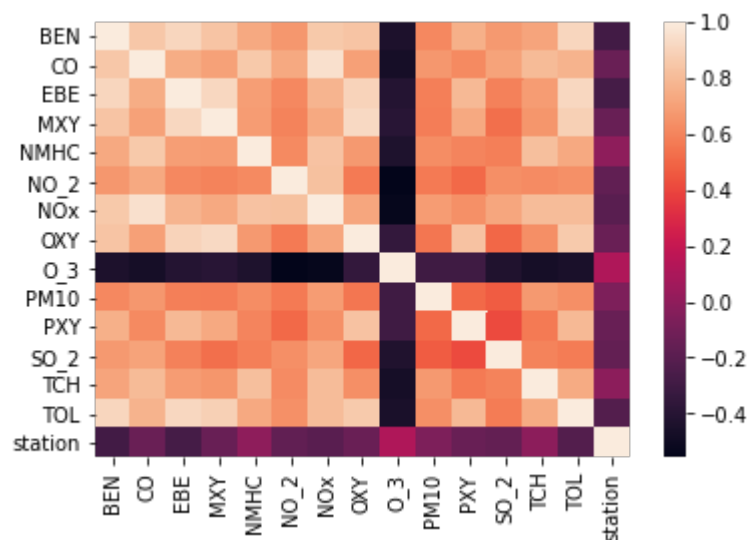
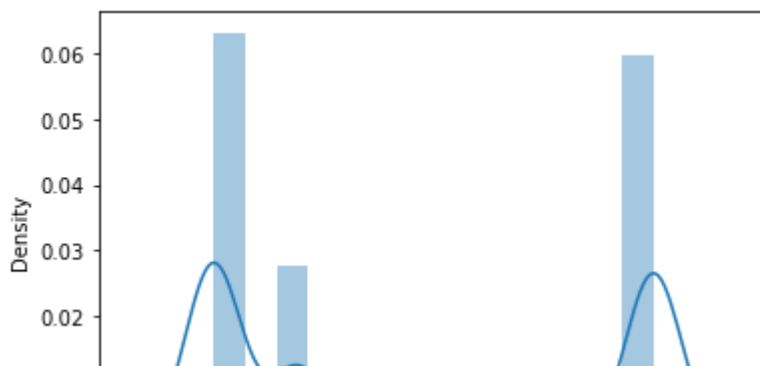
Out[19]: <seaborn.axisgrid.PairGrid at 0x1fb4e0b4d00>



10.1371/journal.pone.0164511.g001

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

```
<AxesSubplot:xlabel= station , ylabel= Density >
```



TO TRAIN THE MODEL AND MODEL BUILDING

```
x=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
      'PM10', 'PXY', 'SO_2', 'TCH', 'TOL']]
```

```
from sklearn.model_selection import train_test_split
```

Linear Regression

```
In [24]: from sklearn.linear_model import LinearRegression  
lr=LinearRegression()
```

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

```
In [25]:
```

```
Out[25]: 28078960.681743436
```

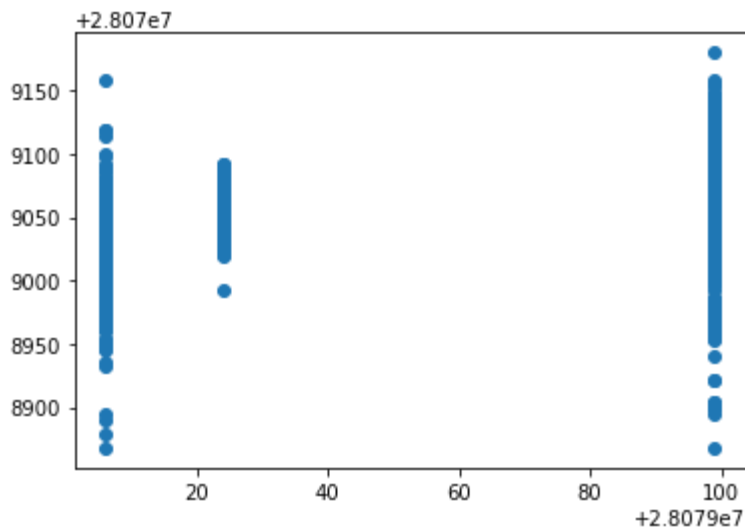
```
In [26]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
```

```
Out[26]:
```

	Co-efficient
BEN	-9.394363
CO	37.322343
EBE	-13.130422
MXY	3.576339
NMHC	78.993132
NO_2	0.131444
NOx	-0.260618
OXY	3.435149
O_3	0.018455
PM10	0.036763
PXY	2.590431
SO_2	0.192576
TCH	60.840751
TOL	-0.596073

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

Out[27]: `<matplotlib.collections.PathCollection at 0x1fb5c9c7580>`



ACCURACY

In [28]: `accuracy = accuracy_score(y_test, prediction)`

Out[28]: `0.3248407277105878`

In [29]: `accuracy = accuracy_score(y_test, prediction)`

Out[29]: `0.2950540583701028`

Ridge and Lasso

In [30]: `from sklearn.linear_model import Ridge, Lasso`

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

Out[31]: `Ridge(alpha=10)`

Accuracy(Ridge)

In [32]: `accuracy = accuracy_score(y_test, rr.predict(x_test))`

Out[32]: `0.32397714607462713`

In [33]: `accuracy = accuracy_score(y_test, rr.predict(x_test))`

Out[33]: `0.2948207426682077`

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

```
Out[34]: Lasso(alpha=10)
```

```
In [35]:
```

```
Out[35]: 0.06357408465154679
```

Accuracy(Lasso)

```
In [36]:
```

```
Out[36]: 0.06436861694506213
```

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

```
Out[37]: ElasticNet()
```

```
In [38]:
```

```
Out[38]: array([-5.39142485,  1.39997368, -7.3101742 ,  2.55109973,  0.85173608,  
                -0.04312806, -0.00851285,  1.97508995, -0.01565268,  0.21613232,  
                1.31201436,  0.11359589,  1.49553805, -0.771418  ])
```

```
In [39]:
```

```
Out[39]: 28079049.603366435
```

```
In [40]:
```

```
In [41]:
```

```
Out[41]: 0.1828750526236339
```

Evaluation Metrics

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

```
37.08273544256618
```

```
1550.5842939117824
```

```
39.37745921097224
```

Logistic Regression

In [43]:

```
In [44]: feature_matrix=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O  
          'PM10', 'PXY', 'SO_2', 'TCH', 'TOL']]
```

In [45]:

Out[45]: (20070, 14)

In [46]:

Out[46]: (20070,)

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)  
[28079006]
```

In [52]:

Out[52]: array([28079006, 28079024, 28079099], dtype=int64)

In [53]:

Out[53]: 0.879023418036871

In [54]:

Out[54]: 0.9998967601812779

In [55]:

Out[55]: array([[9.99896760e-01, 3.21124597e-30, 1.03239819e-04]])

Random Forest

In [56]:

```
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="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]: 0.8629793732115209
```

```
In [61]:
```


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

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

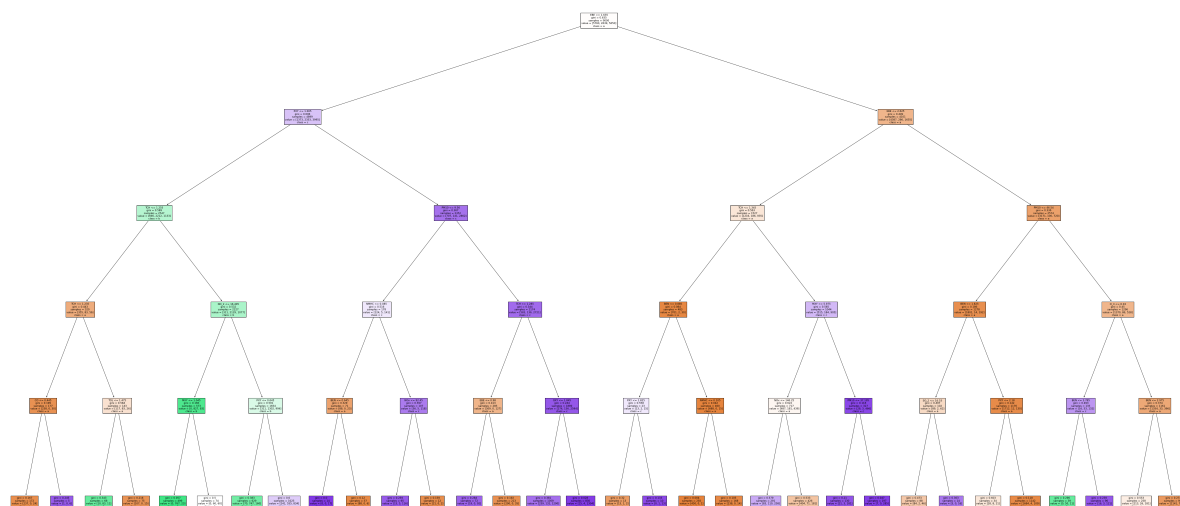
Out[62]: [Text(2232.0, 1993.2, 'EBE <= 1.635\ngini = 0.635\nsamples = 9000\nvalue = [5760, 2639, 5650]\nnclass = a'),
Text(1116.0, 1630.8000000000002, 'PXY <= 1.005\ngini = 0.608\nsamples = 4899\nvalue = [1373, 2353, 3995]\nnclass = c'),
Text(558.0, 1268.4, 'TCH <= 1.255\ngini = 0.589\nsamples = 2547\nvalue = [666, 2212, 1133]\nnclass = b'),
Text(279.0, 906.0, 'TCH <= 1.235\ngini = 0.443\nsamples = 320\nvalue = [355, 83, 56]\nnclass = a'),
Text(139.5, 543.5999999999999, 'CO <= 0.445\ngini = 0.199\nsamples = 177\nvalue = [238, 0, 30]\nnclass = a'),
Text(69.75, 181.19999999999982, 'gini = 0.167\nsamples = 172\nvalue = [237, 0, 24]\nnclass = a'),
Text(209.25, 181.19999999999982, 'gini = 0.245\nsamples = 5\nvalue = [1, 0, 6]\nnclass = c'),
Text(418.5, 543.5999999999999, 'TOL <= 1.475\ngini = 0.584\nsamples = 143\nvalue = [117, 83, 26]\nnclass = a'),
Text(348.75, 181.19999999999982, 'gini = 0.343\nsamples = 68\nvalue = [10, 83, 11]\nnclass = b'),
Text(488.25, 181.19999999999982, 'gini = 0.216\nsamples = 75\nvalue = [107, 0, 15]\nnclass = a'),
Text(837.0, 906.0, 'NO_2 <= 16.405\ngini = 0.532\nsamples = 2227\nvalue = [311, 2129, 1077]\nnclass = b'),
Text(697.5, 543.5999999999999, 'MXY <= 1.045\ngini = 0.166\nsamples = 573\nvalue = [0, 827, 83]\nnclass = b'),
Text(627.75, 181.19999999999982, 'gini = 0.057\nsamples = 499\nvalue = [0, 767, 23]\nnclass = b'),
Text(767.25, 181.19999999999982, 'gini = 0.5\nsamples = 74\nvalue = [0, 60, 60]\nnclass = b'),
Text(976.5, 543.5999999999999, 'PXY <= 0.645\ngini = 0.591\nsamples = 1654\nvalue = [311, 1302, 994]\nnclass = b'),
Text(906.75, 181.19999999999982, 'gini = 0.383\nsamples = 629\nvalue = [70, 747, 160]\nnclass = b'),
Text(1046.25, 181.19999999999982, 'gini = 0.6\nsamples = 1025\nvalue = [241, 555, 834]\nnclass = c'),
Text(1674.0, 1268.4, 'PM10 <= 9.56\ngini = 0.367\nsamples = 2352\nvalue = [707, 141, 2862]\nnclass = c'),
Text(1395.0, 906.0, 'NMHC <= 0.085\ngini = 0.516\nsamples = 178\nvalue = [124, 5, 141]\nnclass = c'),
Text(1255.5, 543.5999999999999, 'BEN <= 0.645\ngini = 0.329\nsamples = 71\nvalue = [88, 0, 23]\nnclass = a'),
Text(1185.75, 181.19999999999982, 'gini = 0.0\nsamples = 14\nvalue = [0, 0, 17]\nnclass = c'),
Text(1325.25, 181.19999999999982, 'gini = 0.12\nsamples = 57\nvalue = [88, 0, 6]\nnclass = a'),
Text(1534.5, 543.5999999999999, 'NOx <= 92.45\ngini = 0.397\nsamples = 107\nvalue = [36, 5, 118]\nnclass = c'),
Text(1464.75, 181.19999999999982, 'gini = 0.259\nsamples = 95\nvalue = [15, 5, 116]\nnclass = c'),
Text(1604.25, 181.19999999999982, 'gini = 0.159\nsamples = 12\nvalue = [21, 0, 2]\nnclass = a'),
Text(1953.0, 906.0, 'TCH <= 1.285\ngini = 0.344\nsamples = 2174\nvalue = [583, 136, 2721]\nnclass = c'),

```
Text(1813.5, 543.5999999999999, 'EBE <= 0.98\ngini = 0.413\nsamples = 285\nvalue = [309, 0, 127]\nclass = a'),
Text(1743.75, 181.19999999999982, 'gini = 0.284\nsamples = 72\nvalue = [19,
0, 92]\nclass = c'),
Text(1883.25, 181.19999999999982, 'gini = 0.192\nsamples = 213\nvalue = [29
0, 0, 35]\nclass = a'),
Text(2092.5, 543.5999999999999, 'OXY <= 1.845\ngini = 0.244\nsamples = 1889\nvalue = [274, 136, 2594]\nclass = c'),
Text(2022.75, 181.19999999999982, 'gini = 0.381\nsamples = 1059\nvalue = [25
9, 132, 1290]\nclass = c'),
Text(2162.25, 181.19999999999982, 'gini = 0.028\nsamples = 830\nvalue = [15,
4, 1304]\nclass = c'),
Text(3348.0, 1630.8000000000002, 'EBE <= 2.625\ngini = 0.449\nsamples = 410
1\nvalue = [4387, 286, 1655]\nclass = a'),
Text(2790.0, 1268.4, 'TCH <= 1.345\ngini = 0.563\nsamples = 1527\nvalue = [1
216, 186, 935]\nclass = a'),
Text(2511.0, 906.0, 'BEN <= 0.885\ngini = 0.084\nsamples = 483\nvalue = [70
1, 2, 30]\nclass = a'),
Text(2371.5, 543.5999999999999, 'PXY <= 1.825\ngini = 0.558\nsamples = 23\nvalue = [13, 2, 15]\nclass = c'),
Text(2301.75, 181.19999999999982, 'gini = 0.32\nsamples = 13\nvalue = [13,
1, 2]\nclass = a'),
Text(2441.25, 181.19999999999982, 'gini = 0.133\nsamples = 10\nvalue = [0,
1, 13]\nclass = c'),
Text(2650.5, 543.5999999999999, 'NMHC <= 0.105\ngini = 0.042\nsamples = 460\nvalue = [688, 0, 15]\nclass = a'),
Text(2580.75, 181.19999999999982, 'gini = 0.004\nsamples = 292\nvalue = [45
0, 0, 1]\nclass = a'),
Text(2720.25, 181.19999999999982, 'gini = 0.105\nsamples = 168\nvalue = [23
8, 0, 14]\nclass = a'),
Text(3069.0, 906.0, 'MXV <= 5.975\ngini = 0.565\nsamples = 1044\nvalue = [51
5, 184, 905]\nclass = c'),
Text(2929.5, 543.5999999999999, 'NOx <= 146.25\ngini = 0.622\nsamples = 717\nvalue = [487, 181, 439]\nclass = a'),
Text(2859.75, 181.19999999999982, 'gini = 0.579\nsamples = 291\nvalue = [83,
110, 258]\nclass = c'),
Text(2999.25, 181.19999999999982, 'gini = 0.533\nsamples = 426\nvalue = [40
4, 71, 181]\nclass = a'),
Text(3208.5, 543.5999999999999, 'PM10 <= 37.195\ngini = 0.118\nsamples = 32
7\nvalue = [28, 3, 466]\nclass = c'),
Text(3138.75, 181.19999999999982, 'gini = 0.21\nsamples = 130\nvalue = [21,
3, 181]\nclass = c'),
Text(3278.25, 181.19999999999982, 'gini = 0.047\nsamples = 197\nvalue = [7,
0, 285]\nclass = c'),
Text(3906.0, 1268.4, 'PM10 <= 46.54\ngini = 0.336\nsamples = 2574\nvalue =
[3171, 100, 720]\nclass = a'),
Text(3627.0, 906.0, 'BEN <= 1.825\ngini = 0.186\nsamples = 1278\nvalue = [18
01, 14, 192]\nclass = a'),
Text(3487.5, 543.5999999999999, 'SO_2 <= 16.53\ngini = 0.497\nsamples = 102\nvalue = [89, 2, 62]\nclass = a'),
Text(3417.75, 181.19999999999982, 'gini = 0.473\nsamples = 88\nvalue = [84,
2, 46]\nclass = a'),
Text(3557.25, 181.19999999999982, 'gini = 0.363\nsamples = 14\nvalue = [5,
0, 16]\nclass = c'),
Text(3766.5, 543.5999999999999, 'PXY <= 2.18\ngini = 0.142\nsamples = 1176\nvalue = [1712, 12, 130]\nclass = a'),
```

```

Text(3696.75, 181.19999999999982, 'gini = 0.603\nsamples = 34\nvalue = [28,
8, 21]\nnclass = a'),
Text(3836.25, 181.19999999999982, 'gini = 0.118\nsamples = 1142\nvalue = [16
84, 4, 109]\nnclass = a'),
Text(4185.0, 906.0, 'O_3 <= 6.03\ngini = 0.45\nsamples = 1296\nvalue = [137
0, 86, 528]\nnclass = a'),
Text(4045.5, 543.5999999999999, 'BEN <= 3.785\ngini = 0.493\nsamples = 135\n
value = [16, 53, 132]\nnclass = c'),
Text(3975.75, 181.19999999999982, 'gini = 0.296\nsamples = 39\nvalue = [0, 5
0, 11]\nnclass = b'),
Text(4115.25, 181.19999999999982, 'gini = 0.239\nsamples = 96\nvalue = [16,
3, 121]\nnclass = c'),
Text(4324.5, 543.5999999999999, 'BEN <= 2.975\ngini = 0.374\nsamples = 1161\
nvalue = [1354, 33, 396]\nnclass = a'),
Text(4254.75, 181.19999999999982, 'gini = 0.553\nsamples = 258\nvalue = [21
3, 26, 181]\nnclass = a'),
Text(4394.25, 181.19999999999982, 'gini = 0.274\nsamples = 903\nvalue = [114
1, 7, 215]\nnclass = a')]

```



Conclusion

Accuracy

Linear Regression : 0.2950540583701028

Ridge Regression : 0.06357408465154679

Lasso Regression : 0.06436861694506213

ElasticNet Regression : 0.1828750526236339

Logistic Regression : 0.9998967601812779

Logistic Regression is suitable for this dataset