

Importing Libraries

```
In [1]: 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\Desktop\csvs_per_year\csvs_per_year\madrid_2002.csv")
df
```

Out[2]:

	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3	PM10
0	2002-04-01 01:00:00	NaN	1.39	NaN	NaN	NaN	145.100006	352.100006	NaN	6.54	41.990002
1	2002-04-01 01:00:00	1.93	0.71	2.33	6.20	0.15	98.150002	153.399994	2.67	6.85	20.980000
2	2002-04-01 01:00:00	NaN	0.80	NaN	NaN	NaN	103.699997	134.000000	NaN	13.01	28.440001
3	2002-04-01 01:00:00	NaN	1.61	NaN	NaN	NaN	97.599998	268.000000	NaN	5.12	42.180000
4	2002-04-01 01:00:00	NaN	1.90	NaN	NaN	NaN	92.089996	237.199997	NaN	7.28	76.330002
...
217291	2002-11-01 00:00:00	4.16	1.14	NaN	NaN	NaN	81.080002	265.700012	NaN	7.21	36.750000
217292	2002-11-01 00:00:00	3.67	1.73	2.89	NaN	0.38	113.900002	373.100006	NaN	5.66	63.389999
217293	2002-11-01 00:00:00	1.37	0.58	1.17	2.37	0.15	65.389999	107.699997	1.30	9.11	9.640000
217294	2002-11-01 00:00:00	4.51	0.91	4.83	10.99	NaN	149.800003	202.199997	1.00	5.75	NaN
217295	2002-11-01 00:00:00	3.11	1.17	3.00	7.77	0.26	80.110001	180.300003	2.25	7.38	29.240000

217296 rows × 16 columns

Data Cleaning and Data Preprocessing

```
In [3]: df=df.dropna()
```

```
In [4]: df.columns
```

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

```
In [5]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 32381 entries, 1 to 217295  
Data columns (total 16 columns):  
#   Column      Non-Null Count  Dtype  
---  ---  
0   date        32381 non-null  object  
1   BEN         32381 non-null  float64  
2   CO          32381 non-null  float64  
3   EBE         32381 non-null  float64  
4   MXY         32381 non-null  float64  
5   NMHC        32381 non-null  float64  
6   NO_2        32381 non-null  float64  
7   NOx         32381 non-null  float64  
8   OXY         32381 non-null  float64  
9   O_3         32381 non-null  float64  
10  PM10        32381 non-null  float64  
11  PXY         32381 non-null  float64  
12  SO_2        32381 non-null  float64  
13  TCH         32381 non-null  float64  
14  TOL         32381 non-null  float64  
15  station     32381 non-null  int64  
dtypes: float64(14), int64(1), object(1)  
memory usage: 4.2+ MB
```

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

Out[6]:

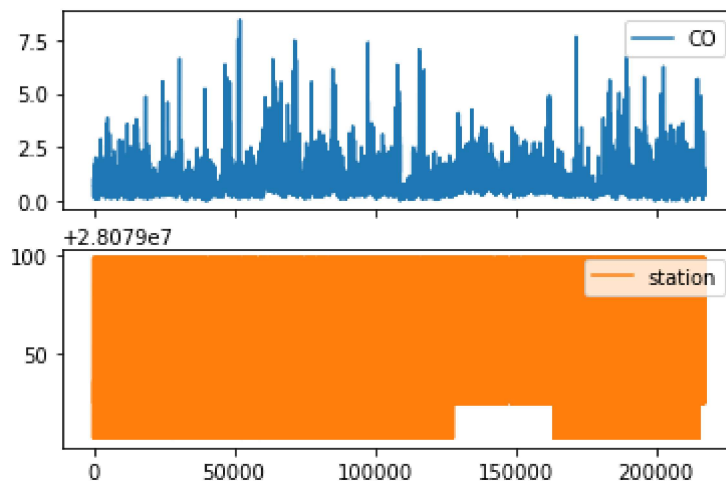
	CO	station
1	0.71	28079035
5	0.72	28079006
22	0.80	28079024
24	1.04	28079099
26	0.53	28079035
...
217269	0.28	28079024
217271	1.30	28079099
217273	0.97	28079035
217293	0.58	28079024
217295	1.17	28079099

32381 rows × 2 columns

Line chart

```
In [7]: data.plot.line(subplots=True)
```

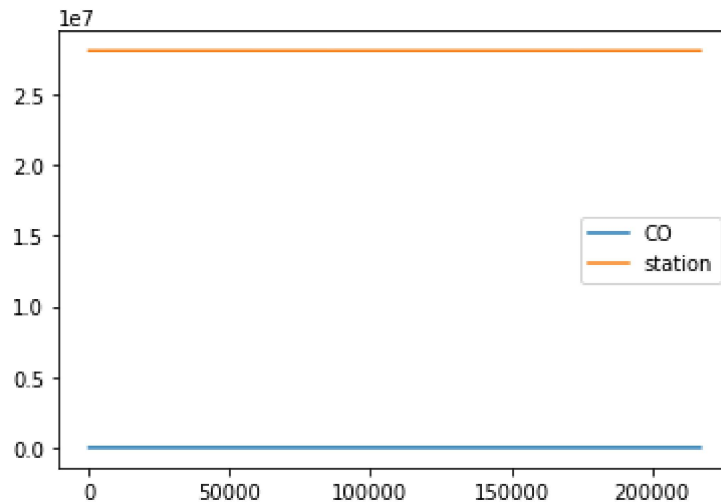
Out[7]: array([<AxesSubplot:>, <AxesSubplot:>], dtype=object)



Line chart

```
In [8]: data.plot.line()
```

```
Out[8]: <AxesSubplot:>
```

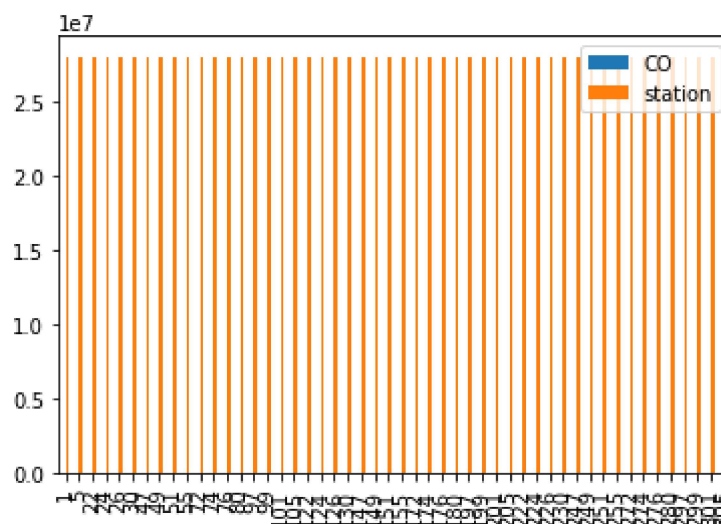


Bar chart

```
In [9]: b=data[0:50]
```

```
In [10]: b.plot.bar()
```

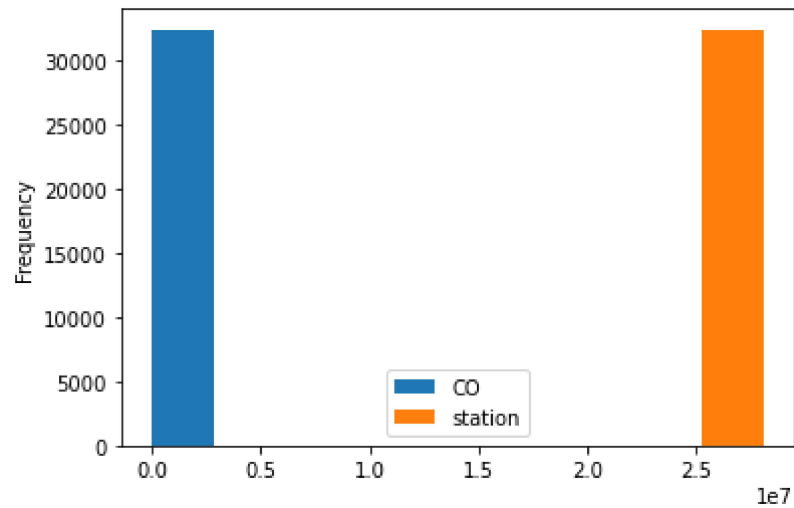
```
Out[10]: <AxesSubplot:>
```



Histogram

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

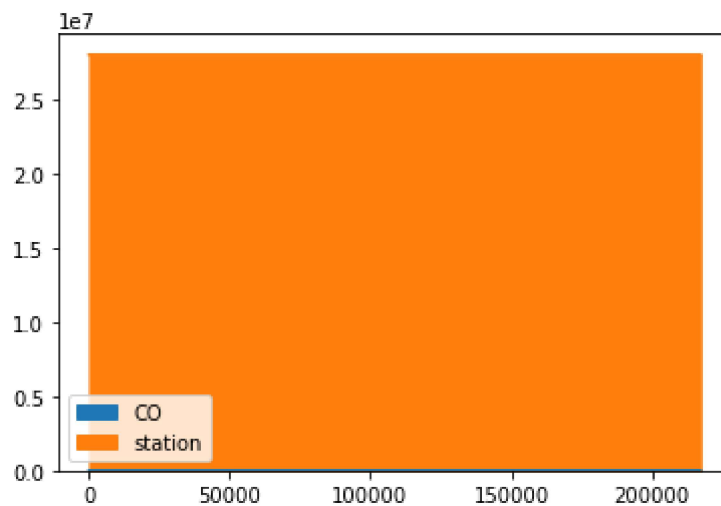
```
Out[11]: <AxesSubplot:ylabel='Frequency'>
```



Area chart

```
In [12]: data.plot.area()
```

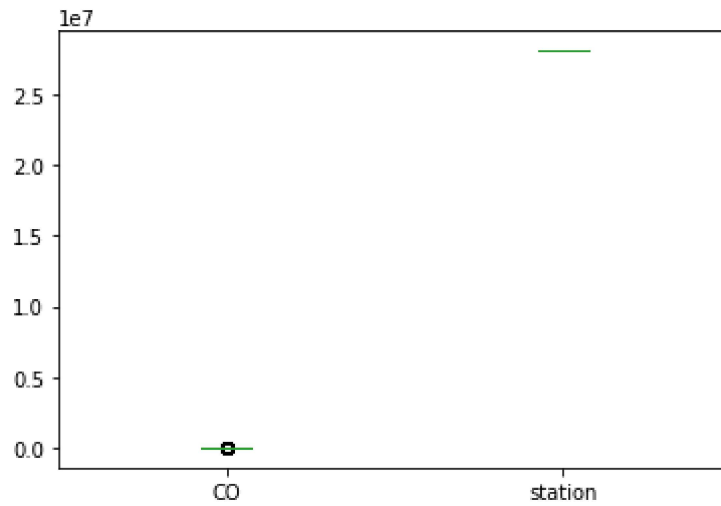
```
Out[12]: <AxesSubplot:>
```



Box chart

```
In [13]: data.plot.box()
```

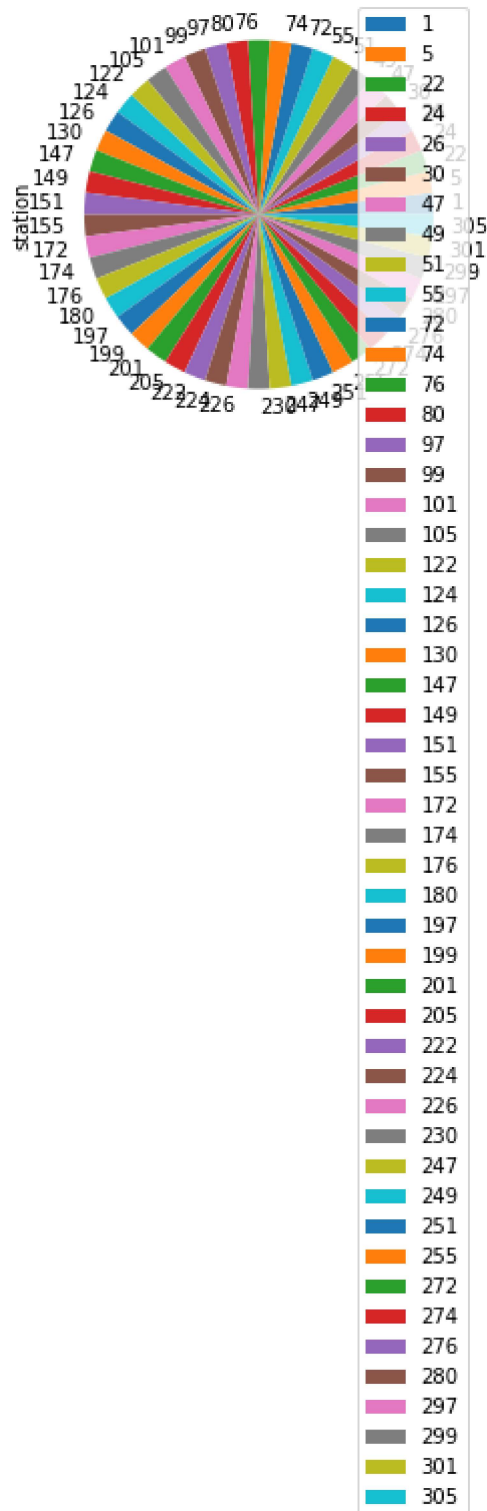
```
Out[13]: <AxesSubplot:>
```



Pie chart

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

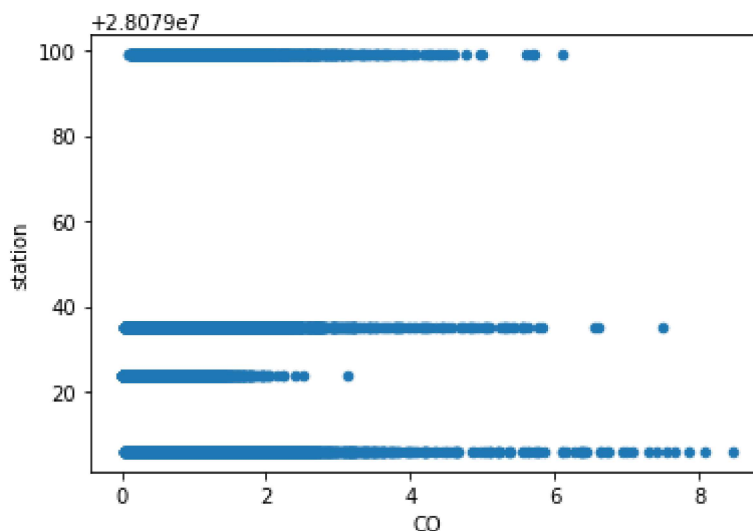
```
Out[14]: <AxesSubplot:ylabel='station'>
```



Scatter chart

```
In [15]: data.plot.scatter(x='CO' ,y='station')
```

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



```
In [16]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 32381 entries, 1 to 217295
Data columns (total 16 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        32381 non-null  object
1   BEN         32381 non-null  float64
2   CO          32381 non-null  float64
3   EBE         32381 non-null  float64
4   MXY         32381 non-null  float64
5   NMHC        32381 non-null  float64
6   NO_2        32381 non-null  float64
7   NOx         32381 non-null  float64
8   OXY         32381 non-null  float64
9   O_3         32381 non-null  float64
10  PM10        32381 non-null  float64
11  PXY         32381 non-null  float64
12  SO_2        32381 non-null  float64
13  TCH         32381 non-null  float64
14  TOL         32381 non-null  float64
15  station     32381 non-null  int64
dtypes: float64(14), int64(1), object(1)
memory usage: 4.2+ MB
```



```
In [17]: df.describe()
```

```
Out[17]:
```

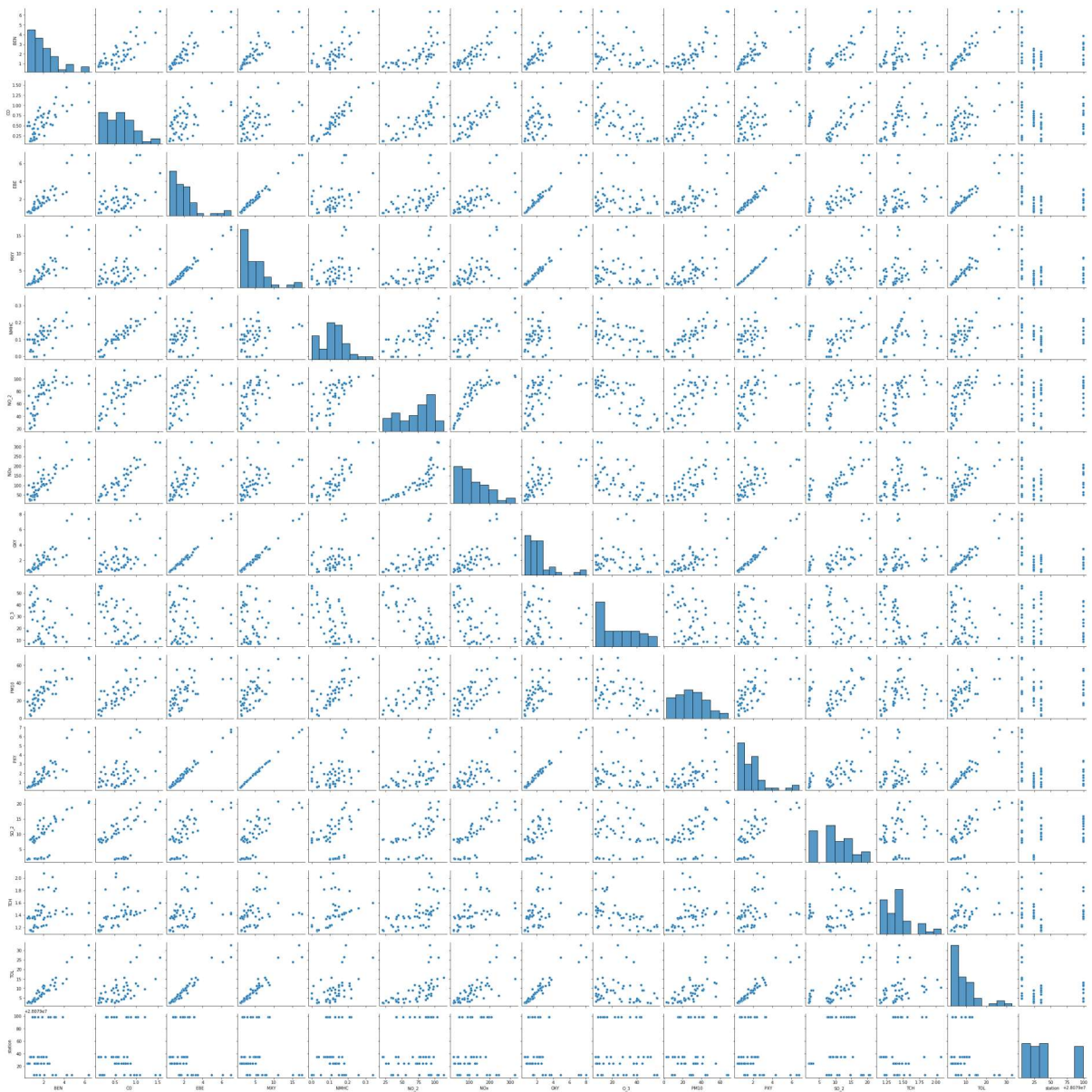
	BEN	CO	EBE	MXY	NMHC	NO_2	
count	32381.000000	32381.000000	32381.000000	32381.000000	32381.000000	32381.000000	32381.000000
mean	2.479155	0.787323	2.914004	7.013636	0.155827	58.936796	13.479997
std	2.280959	0.610810	2.667881	6.774365	0.135731	31.472733	13.479997
min	0.180000	0.000000	0.180000	0.190000	0.000000	0.890000	0.180000
25%	0.970000	0.420000	1.140000	2.420000	0.080000	35.660000	4.140000
50%	1.840000	0.620000	2.130000	5.140000	0.130000	57.160000	9.130000
75%	3.250000	0.980000	3.830000	9.420000	0.200000	78.769997	13.479997
max	32.660000	8.460000	41.740002	99.879997	2.700000	263.600006	134.799997

```
In [18]: df1=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',  
                'PM10', 'PXY', 'SO_2', 'TCH', 'TOL', 'station']]
```

EDA AND VISUALIZATION

```
In [19]: sns.pairplot(df1[0:50])
```

```
Out[19]: <seaborn.axisgrid.PairGrid at 0x1cd706f87c0>
```

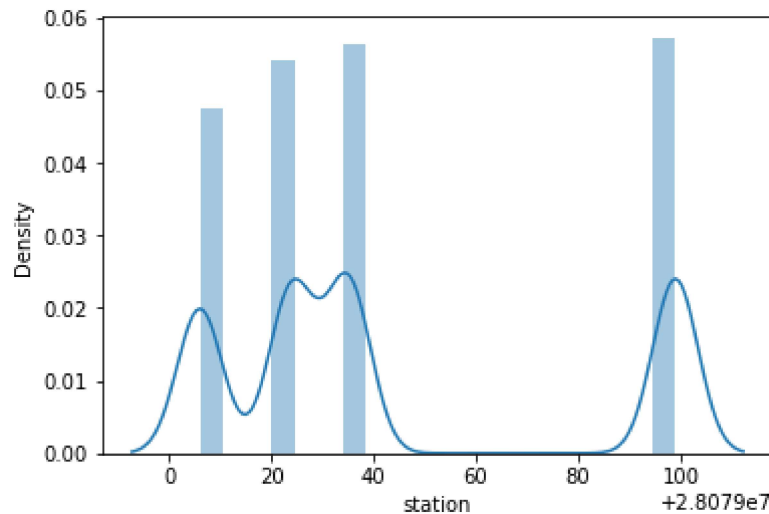


```
In [20]: sns.distplot(df1['station'])
```

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).

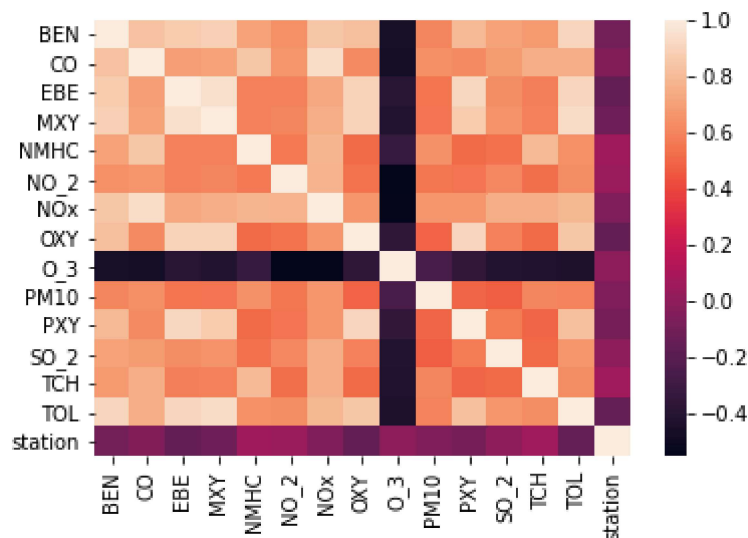
warnings.warn(msg, FutureWarning)

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



```
In [21]: sns.heatmap(df1.corr())
```

```
Out[21]: <AxesSubplot:>
```



TO TRAIN THE MODEL AND MODEL BUILDING

```
In [22]: x=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
              'PM10', 'PXY', 'SO_2', 'TCH', 'TOL']]
y=df['station']
```

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

Out[24]: LinearRegression()

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

Out[25]: 28078986.118767593

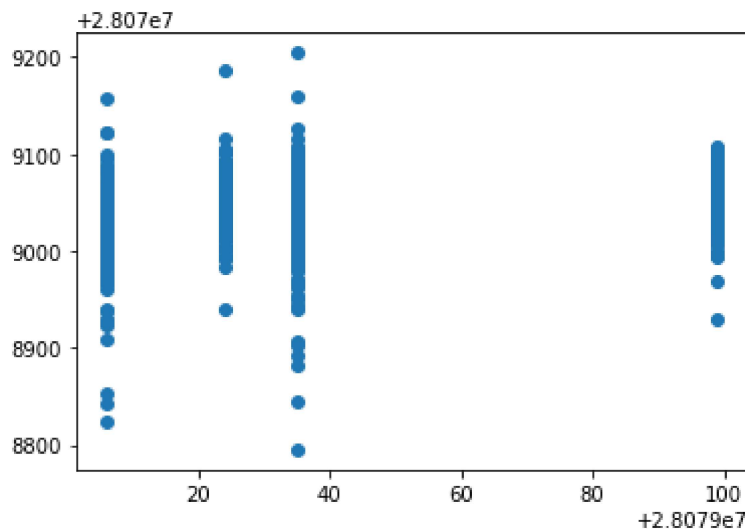
```
In [26]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

Out[26]:

	Co-efficient
BEN	2.274450
CO	-13.582455
EBE	-11.456936
MXY	4.132700
NMHC	78.599697
NO_2	0.248235
NOx	-0.094317
OXY	-5.283710
O_3	-0.025774
PM10	-0.108382
PXY	7.735897
SO_2	0.600407
TCH	45.166445
TOL	-1.564016

```
In [27]: prediction =lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[27]: <matplotlib.collections.PathCollection at 0x1cd00f6ab50>



ACCURACY

```
In [28]: lr.score(x_test,y_test)
```

Out[28]: 0.1937170619900085

```
In [29]: lr.score(x_train,y_train)
```

Out[29]: 0.20059181248043334

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]: 0.19266041562106995
```

```
In [33]: rr.score(x_train,y_train)
```

```
Out[33]: 0.2004088944667991
```

```
In [34]: la=Lasso(alpha=10)  
la.fit(x_train,y_train)
```

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

```
In [35]: la.score(x_train,y_train)
```

```
Out[35]: 0.057740736606634924
```

Accuracy(Lasso)

```
In [36]: la.score(x_test,y_test)
```

```
Out[36]: 0.05834962560339074
```

```
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([ 1.14363723,  0.          , -2.8545531 ,  1.62939227,  0.21486699,  
                0.22571418, -0.02632462, -2.47441929, -0.02159756,  0.01963492,  
                2.25749049,  0.38928751,  1.11558424, -1.26974741])
```

```
In [39]: en.intercept_
```

```
Out[39]: 28079038.361825775
```

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

```
In [41]: en.score(x_test,y_test)
```

```
Out[41]: 0.09734806326924794
```

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)))
```

```
28.44943271497292
1112.1725165149585
33.34925061399369
```

Logistic Regression

```
In [43]: from sklearn.linear_model import LogisticRegression
```

```
In [44]: feature_matrix=df[['BEN', 'CO', 'EBE', 'MXV', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
'PM10', 'PXY', 'SO_2', 'TCH', 'TOL']]
target_vector=df[ 'station']
```

```
In [45]: feature_matrix.shape
```

```
Out[45]: (32381, 14)
```

```
In [46]: target_vector.shape
```

```
Out[46]: (32381,)
```

```
In [47]: from sklearn.preprocessing import StandardScaler
```

```
In [48]: fs=StandardScaler().fit_transform(feature_matrix)
```

```
In [49]: logr=LogisticRegression(max_iter=10000)
logr.fit(fs,target_vector)
```

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

```
In [50]: observation=[[1,2,3,4,5,6,7,8,9,10,11,12,13,14]]
```

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

```
[28079035]
```

```
In [52]: logr.classes_
```

```
Out[52]: array([28079006, 28079024, 28079035, 28079099], dtype=int64)
```

```
In [53]: logr.score(fs,target_vector)
```

```
Out[53]: 0.8480899292795158
```

```
In [54]: logr.predict_proba(observation)[0][0]
```

```
Out[54]: 2.5638972732451705e-10
```

```
In [55]: logr.predict_proba(observation)
```

```
Out[55]: array([[2.56389727e-10, 3.44199742e-71, 1.00000000e+00, 1.43898646e-13]])
```

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]: 0.7723462454778082
```

```
In [61]: rfc_best = grid_search.best_estimator_
```

```
In [62]: 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'],
```

```
Out[62]: [Text(2232.0, 1993.2, 'NO_2 <= 26.995\ngini = 0.749\nsamples = 14325\nvalue
= [5014, 5610, 5961, 6081]\nclass = d'),
Text(1116.0, 1630.8000000000002, 'TCH <= 1.225\ngini = 0.503\nsamples = 23
46\nvalue = [369, 2492, 450, 352]\nclass = b'),
Text(558.0, 1268.4, 'SO_2 <= 6.335\ngini = 0.66\nsamples = 562\nvalue = [2
48, 103, 431, 108]\nclass = c'),
Text(279.0, 906.0, 'NMHC <= 0.015\ngini = 0.638\nsamples = 238\nvalue = [1
74, 102, 3, 85]\nclass = a'),
Text(139.5, 543.5999999999999, 'OXY <= 0.865\ngini = 0.011\nsamples = 113
\nvalue = [173, 0, 1, 0]\nclass = a'),
Text(69.75, 181.19999999999982, 'gini = 0.111\nsamples = 10\nvalue = [16,
0, 1, 0]\nclass = a'),
Text(209.25, 181.19999999999982, 'gini = 0.0\nsamples = 103\nvalue = [157,
0, 0, 0]\nclass = a'),
Text(418.5, 543.5999999999999, 'SO_2 <= 4.805\ngini = 0.512\nsamples = 125
\nvalue = [1, 102, 2, 85]\nclass = b'),
Text(348.75, 181.19999999999982, 'gini = 0.059\nsamples = 64\nvalue = [0,
96, 0, 3]\nclass = b'),
Text(488.25, 181.19999999999982, 'gini = 0.183\nsamples = 61\nvalue = [1,
```

Conclusion

Accuracy

Linear Regression::0.1994901005045402

Ridge Regression:0.19877953137155935

Lasso Regression::0.058987010961436104

ElasticNet Regression:0.09846653794618532

Logistic Regression:0.8087229094340894

Random Forest:0.7732727433159798

From the above data, we can conclude that logistic regression is preferable to other regression

In []: