

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\Downloads\madrid_2006.csv")
df
```

		date	BEN	CO	EBE	MXY	NMHC	NO_2	NOx	OXY	O_3	PM10	
		01-											
0	02-	NaN	1.84	NaN	NaN	NaN	155.100006	490.100006	NaN	4.880000	97.570000	40.2	
	2006	01:00											
		01-											
1	02-	1.68	1.01	2.38	6.36	0.32	94.339996	229.699997	3.04	7.100000	25.820000		
	2006	01:00											
		01-											
2	02-	NaN	1.25	NaN	NaN	NaN	66.800003	192.000000	NaN	4.430000	34.419998		
	2006	01:00											
		01-											
3	02-	NaN	1.68	NaN	NaN	NaN	103.000000	407.799988	NaN	4.830000	28.260000		
	2006	01:00											
		01-											
4	02-	NaN	1.31	NaN	NaN	NaN	105.400002	269.200012	NaN	6.990000	54.180000		
	2006	01:00											
...
		01-											
230563	05-	5.88	0.83	6.23	NaN	0.20	112.500000	218.000000	NaN	24.389999	93.120003		
	2006	00:00											
		01-											
230564	05-	0.76	0.32	0.48	1.09	0.08	51.900002	54.820000	0.61	48.410000	29.469999	15.6	
	2006	00:00											
		01-											
230565	05-	0.96	NaN	0.69	NaN	0.19	135.100006	179.199997	NaN	11.460000	64.680000	35.0	
	2006	00:00											

	date	BEN	CO	EBE	MXY	NMHC	NO_2	NOx	OXY	O_3	PM10	PM25	PXY	SO_2	TCH	TOL	station
230566	01-05-2006 00:00	0.50	NaN	0.67	NaN	0.10	82.599998	105.599998	NaN	NaN	94.360001						
230567	01-05-2006 00:00	1.95	0.74	1.99	4.00	0.24	107.300003	160.199997	2.01	17.730000	52.490002	27.9					

230568 rows × 17 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', 'PM25', 'PXY', 'SO_2', 'TCH', 'TOL', 'station'], dtype='object')`In [5]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 24758 entries, 5 to 230567
Data columns (total 17 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   date     24758 non-null   object 
 1   BEN      24758 non-null   float64
 2   CO       24758 non-null   float64
 3   EBE      24758 non-null   float64
 4   MXY      24758 non-null   float64
 5   NMHC     24758 non-null   float64
 6   NO_2     24758 non-null   float64
 7   NOx      24758 non-null   float64
 8   OXY      24758 non-null   float64
 9   O_3      24758 non-null   float64
 10  PM10     24758 non-null   float64
 11  PM25     24758 non-null   float64
 12  PXY      24758 non-null   float64
 13  SO_2     24758 non-null   float64
 14  TCH      24758 non-null   float64
 15  TOL      24758 non-null   float64
 16  station   24758 non-null   int64  
dtypes: float64(15), int64(1), object(1)
memory usage: 3.4+ MB
```

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

Out[6]:

	CO	station
5	1.69	28079006
22	0.79	28079024
25	1.47	28079099
31	0.85	28079006
48	0.79	28079024
...
230538	0.40	28079024
230541	0.94	28079099
230547	1.06	28079006
230564	0.32	28079024
230567	0.74	28079099

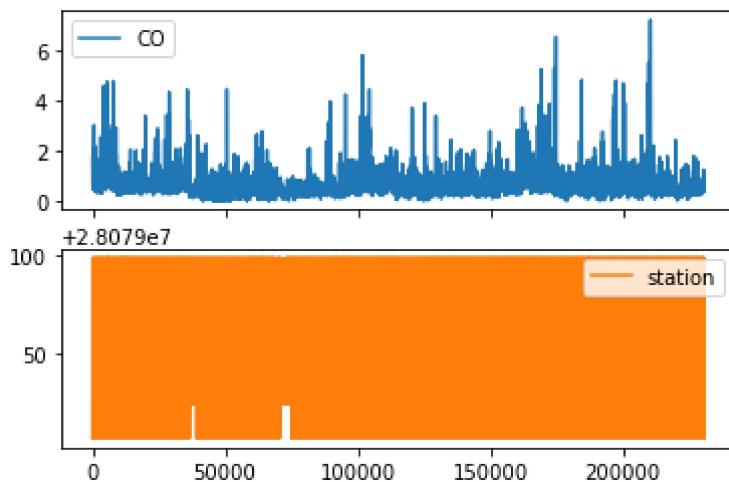
24758 rows × 2 columns

Line chart

In [7]:

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

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

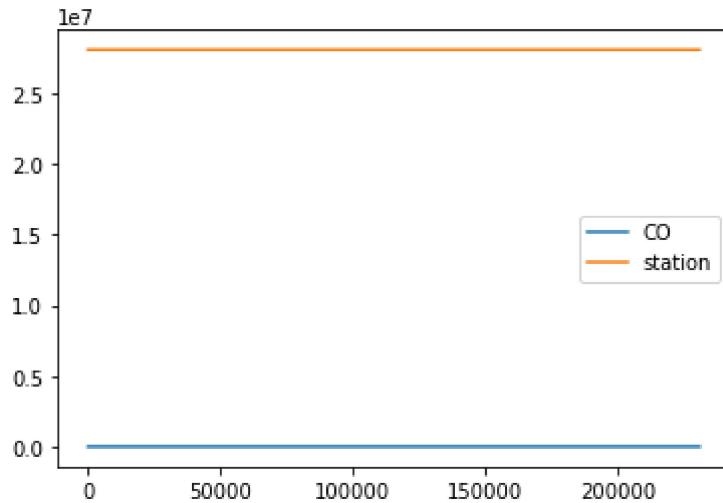


Line chart

In [8]:

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

Out[8]: <AxesSubplot:>

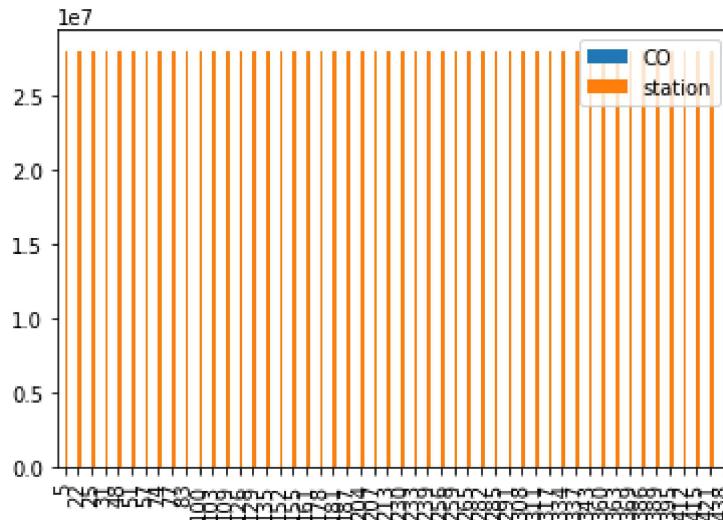


Bar chart

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

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

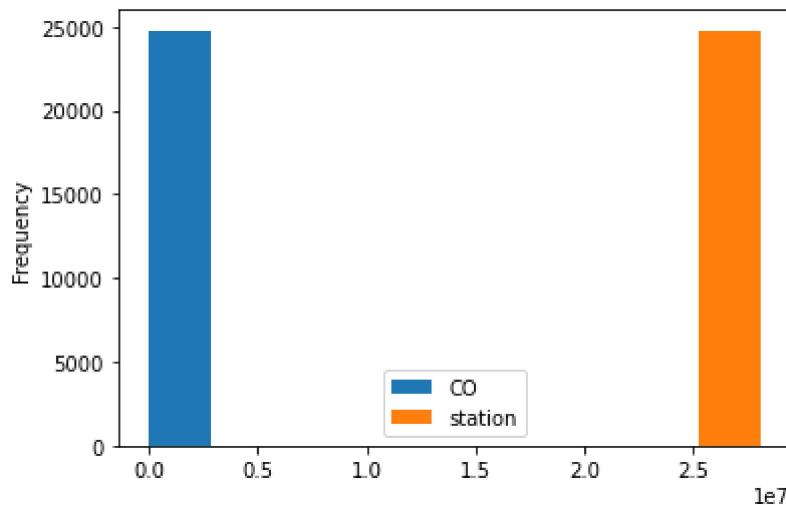
```
Out[10]: <AxesSubplot:>
```



Histogram

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

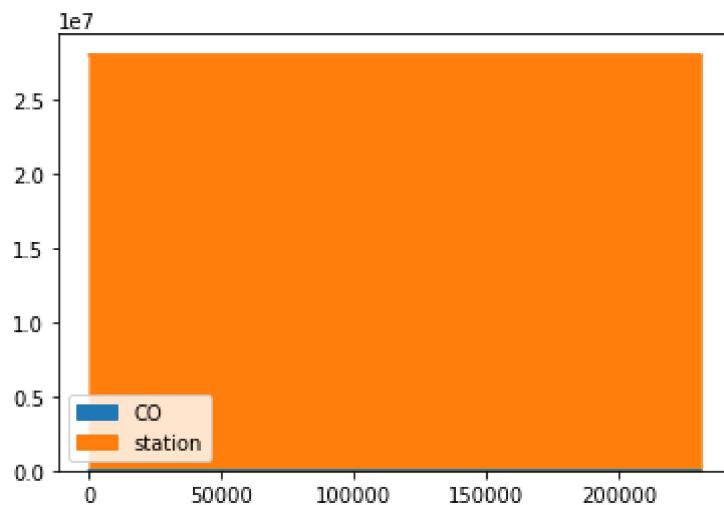
```
Out[11]: <AxesSubplot:ylabel='Frequency'>
```



Area chart

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

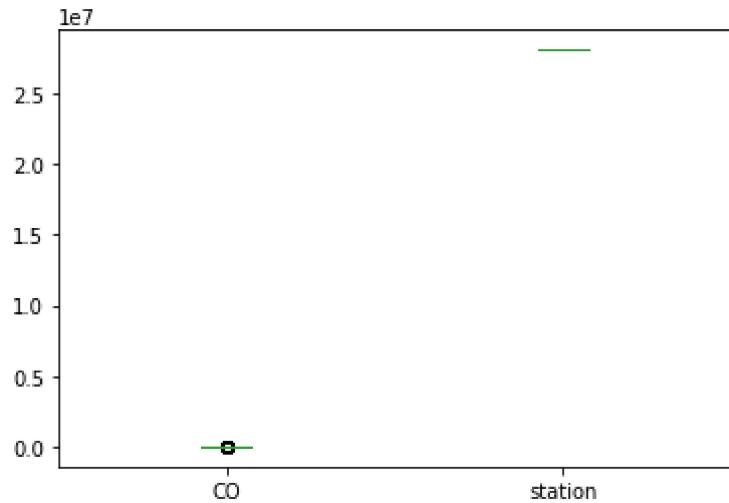
```
Out[12]: <AxesSubplot:>
```



Box chart

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

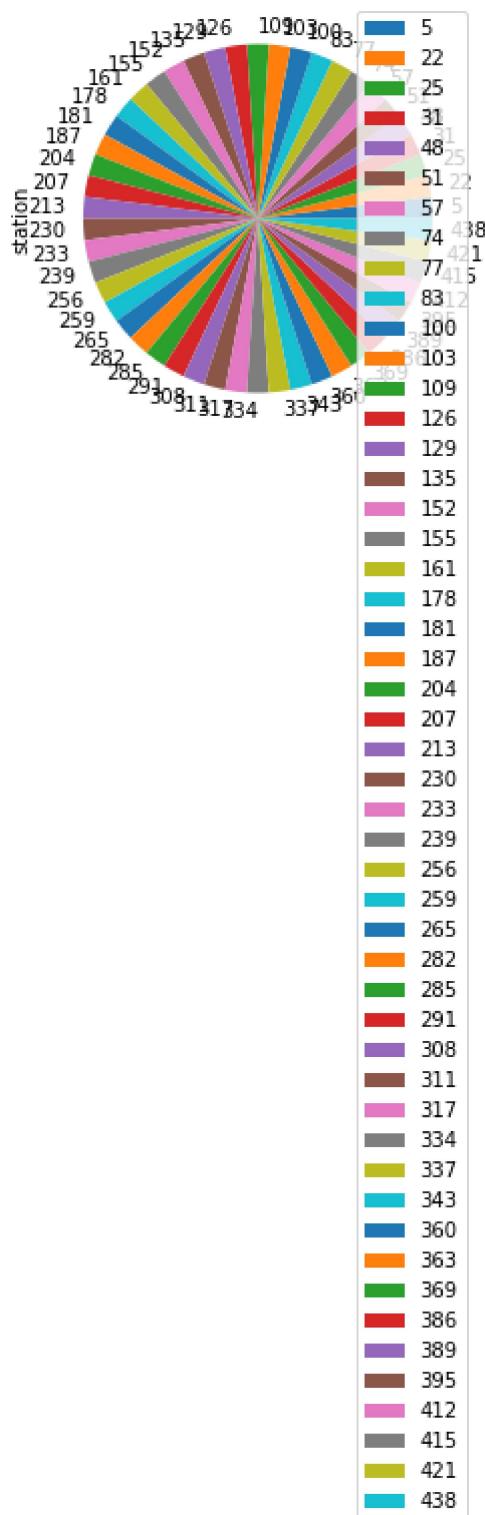
```
Out[13]: <AxesSubplot:>
```



Pie chart

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

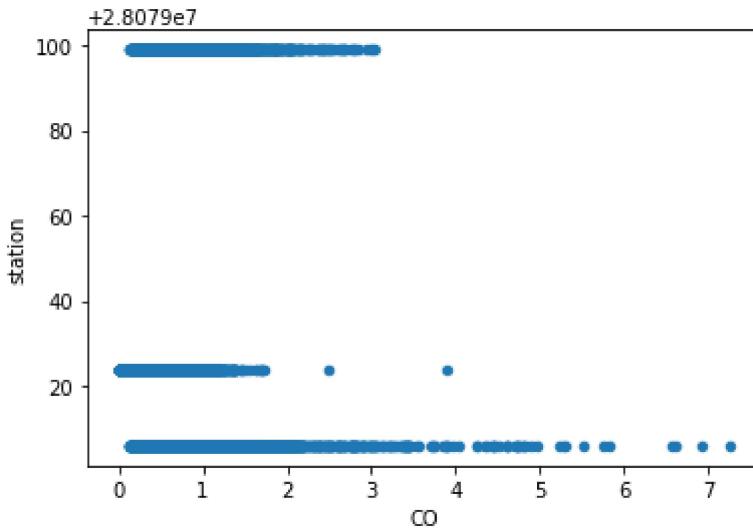
```
Out[14]: <AxesSubplot:ylabel='station'>
```



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: 24758 entries, 5 to 230567
Data columns (total 17 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   date      24758 non-null   object 
 1   BEN        24758 non-null   float64
 2   CO         24758 non-null   float64
 3   EBE        24758 non-null   float64
 4   MXY        24758 non-null   float64
 5   NMHC       24758 non-null   float64
 6   NO_2       24758 non-null   float64
 7   NOX        24758 non-null   float64
 8   OXY        24758 non-null   float64
 9   O_3         24758 non-null   float64
 10  PM10       24758 non-null   float64
 11  PM25       24758 non-null   float64
 12  PXY        24758 non-null   float64
 13  SO_2       24758 non-null   float64
 14  TCH        24758 non-null   float64
 15  TOL        24758 non-null   float64
 16  station    24758 non-null   int64  
dtypes: float64(15), int64(1), object(1)
memory usage: 3.4+ MB
```

In [17]:

`df.describe()`

Out[17]:

	BEN	CO	EBE	MXY	NMHC	NO_2	NO
count	24758.000000	24758.000000	24758.000000	24758.000000	24758.000000	24758.000000	24758.000000
mean	1.350624	0.600713	1.824534	3.835034	0.176546	58.333481	116.419091
std	1.541636	0.419048	1.868939	4.069036	0.126683	40.529382	117.557064
min	0.110000	0.000000	0.170000	0.150000	0.000000	1.680000	2.020000
25%	0.450000	0.360000	0.810000	1.060000	0.100000	28.450001	36.88250
50%	0.850000	0.500000	1.130000	2.500000	0.150000	52.959999	85.180001

	BEN	CO	EBE	MXY	NMHC	NO_2	NOx
75%	1.680000	0.720000	2.160000	5.090000	0.220000	79.347498	158.300000
max	45.430000	7.250000	57.799999	66.900002	2.020000	461.299988	1680.000000

In [18]:

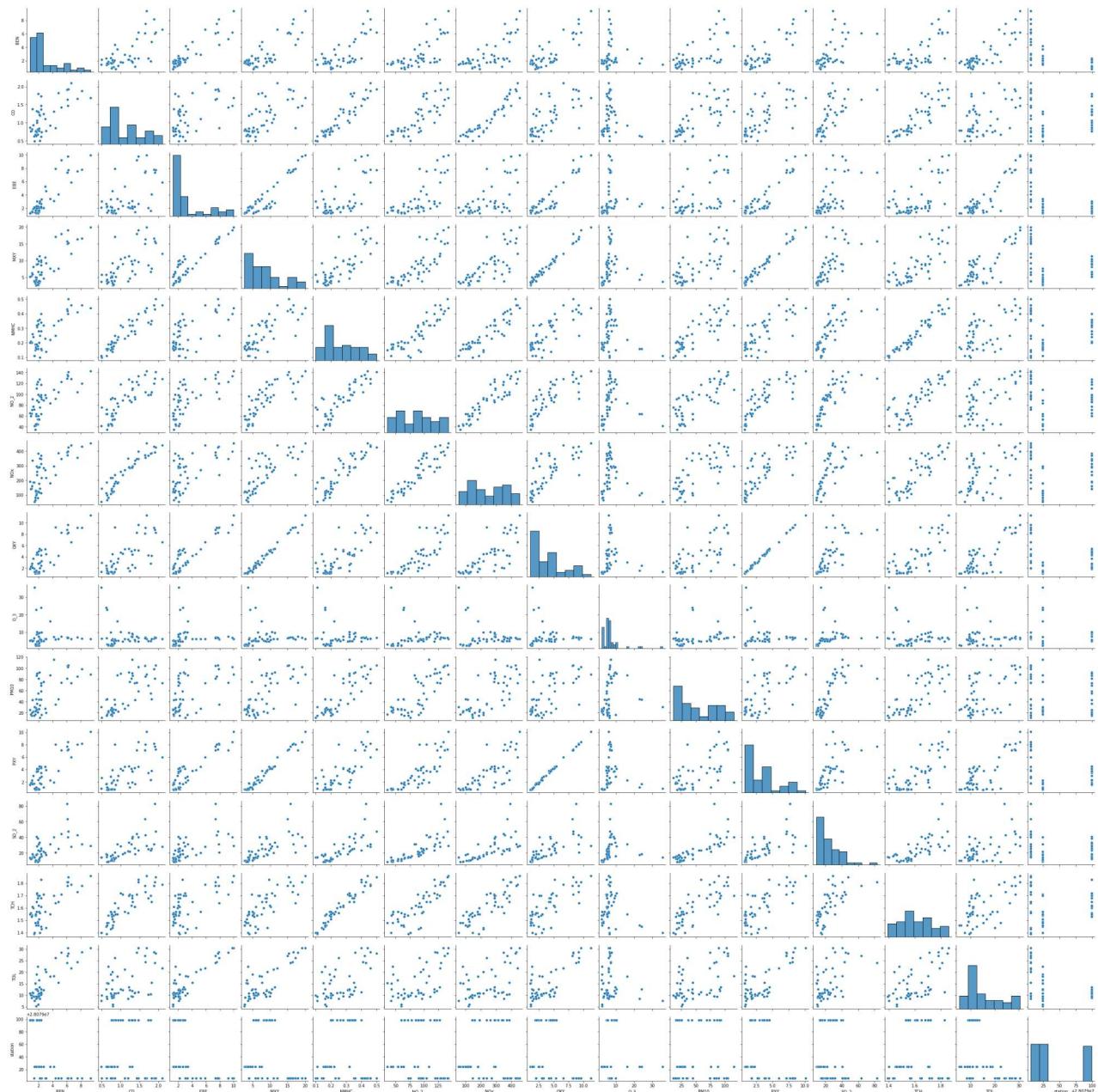
```
df1=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
       'PM10', 'PXY', 'SO_2', 'TCH', 'TOL', 'station']]
```

EDA AND VISUALIZATION

In [19]:

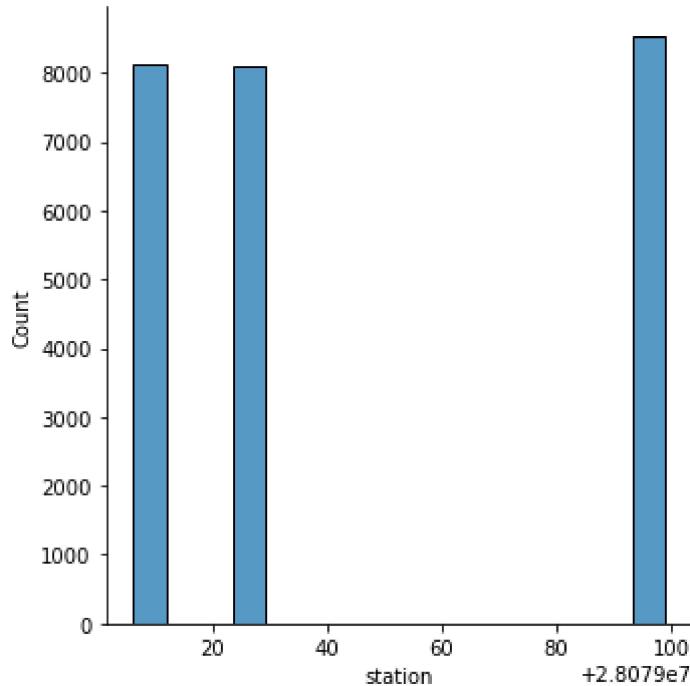
```
sns.pairplot(df1[0:50])
```

Out[19]: <seaborn.axisgrid.PairGrid at 0x17ec0337d90>



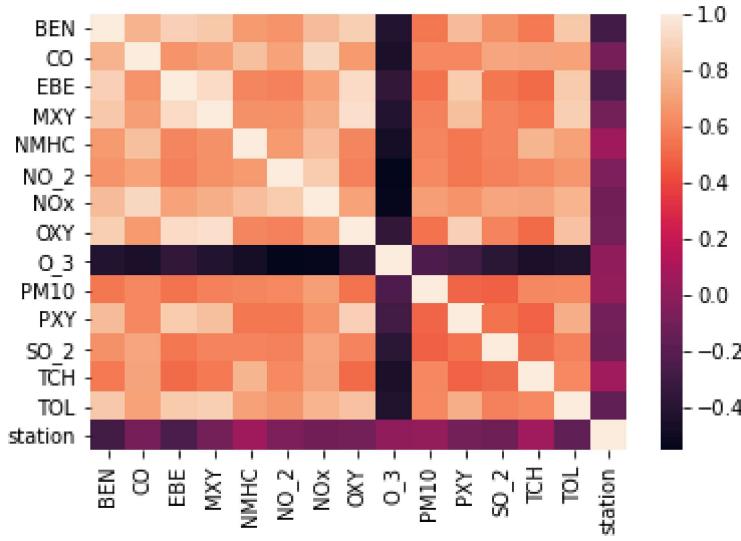
In [20]: `sns.displot(df1['station'])`

Out[20]: <seaborn.axisgrid.FacetGrid at 0x17ecc064b80>



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_
```

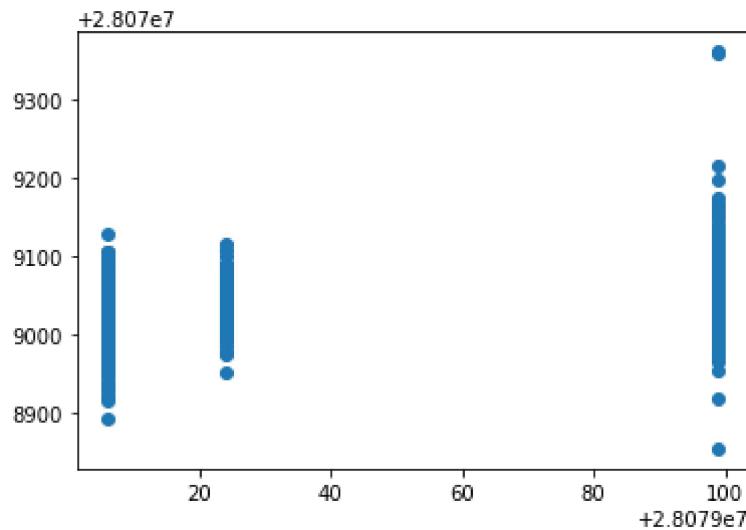
Out[25]: 28079015.645236906

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

	Co-efficient
BEN	-18.705413
CO	-12.229780
EBE	-21.328150
MXY	2.550989
NMHC	123.225183
NO_2	-0.016937
NOx	-0.001274
OXY	18.751170
O_3	-0.076270
PM10	0.152285
PXY	5.956187
SO_2	-0.723212
TCH	22.618359
TOL	-0.683651

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

Out[27]: <matplotlib.collections.PathCollection at 0x17ecedf76a0>



ACCURACY

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

```
Out[28]: 0.3942253498114515
```

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

```
Out[29]: 0.39194541363193647
```

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]: 0.3923734972503121
```

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

```
Out[33]: 0.39136386690973246
```

```
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]: 0.06316793635391815

Accuracy(Lasso)

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

Out[36]: 0.05941509611611351

Elastic Net regression

```
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([-8.35533203e+00, 0.00000000e+00, -8.59840911e+00, 3.09185743e+00,
 4.01177598e-01, -8.28635026e-03, 4.04294315e-03, 3.99986179e+00,
 -1.35922657e-01, 3.11629520e-01, 2.88533244e+00, -5.01608944e-01,
 6.01762095e-01, -1.14802686e+00])

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

Out[39]: 28079052.59413356

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

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

Out[41]: 0.23414894775926287

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)))
```

32.23789819181627
1247.706497002361
35.32288913724868

Logistic Regression

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

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

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

Out[45]: (24758, 14)

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

Out[46]: (24758,)

```
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)
```

[28079099]

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

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

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

```
Out[53]: 0.8741416915744405
```

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

```
Out[54]: 3.5470091844587883e-15
```

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

```
Out[55]: array([[3.54700918e-15, 7.82458875e-29, 1.00000000e+00]])
```

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]: 0.8773225620311598
```

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

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

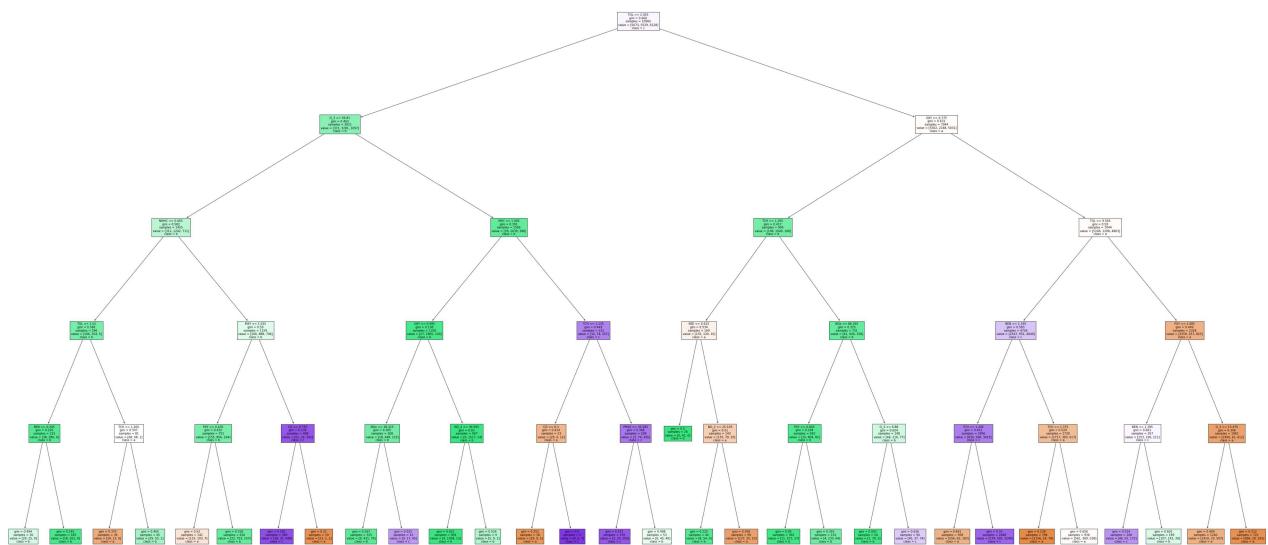
```
plt.figure(figsize=(80,40))
plot_tree/rfc_best.estimators_[5],feature_names=x.columns,class_names=['a','b','c','d'])
```

```
Out[62]: [Text(2241.3, 1993.2, 'TOL <= 2.055\ngini = 0.666\nsamples = 10965\nvalue = [5673, 5529, 6128]\nnclass = c'),
Text(1190.4, 1630.800000000002, 'O_3 <= 56.81\ngini = 0.463\nsamples = 3021\nvalue = [371, 3281, 1097]\nnclass = b'),
Text(595.2, 1268.4, 'NMHC <= 0.065\ngini = 0.582\nsamples = 1455\nvalue = [312, 1242, 711]\nnclass = b'),
Text(297.6, 906.0, 'TOL <= 1.52\ngini = 0.368\nsamples = 296\nvalue = [106, 354, 5]\nnclass = b'),
Text(148.8, 543.599999999999, 'BEN <= 0.205\ngini = 0.226\nsamples = 215\nvalue = [38, 286, 4]\nnclass = b'),
Text(74.4, 181.1999999999982, 'gini = 0.494\nsamples = 30\nvalue = [20, 25, 0]\nnclass = b'),
Text(223.200000000002, 181.1999999999982, 'gini = 0.145\nsamples = 185\nvalue = [18, 261, 4]\nnclass = b'),
Text(446.400000000003, 543.599999999999, 'TCH <= 1.265\ngini = 0.507\nsamples = 81\nvalue = [68, 68, 1]\nnclass = a'),
Text(372.0, 181.1999999999982, 'gini = 0.375\nsamples = 36\nvalue = [39, 13, 0]\nnclass = a'),
Text(520.800000000001, 181.1999999999982, 'gini = 0.465\nsamples = 45\nvalue = [29, 55, 1]\nnclass = b'),
Text(892.800000000001, 906.0, 'MXY <= 1.035\ngini = 0.59\nsamples = 1159\nvalue = [206, 888, 706]\nnclass = b'),
Text(744.0, 543.599999999999, 'PXY <= 0.435\ngini = 0.432\nsamples = 751\nvalue = [155, 856, 164]\nnclass = b'),
Text(669.6, 181.1999999999982, 'gini = 0.52\nsamples = 141\nvalue = [133, 103, 7]\nnclass = a'),
Text(818.400000000001, 181.1999999999982, 'gini = 0.318\nsamples = 610\nvalue = [22, 753, 157]\nnclass = b'),
Text(1041.600000000001, 543.599999999999, 'CO <= 0.755\ngini = 0.239\nsamples = 408\nvalue = [51, 32, 542]\nnclass = c'),
Text(967.2, 181.1999999999982, 'gini = 0.182\nsamples = 389\nvalue = [28, 31, 540]\nnclass = c'),
Text(1116.0, 181.1999999999982, 'gini = 0.21\nsamples = 19\nvalue = [23, 1, 2]\nnclass = a'),
Text(1785.600000000001, 1268.4, 'MXY <= 1.005\ngini = 0.301\nsamples = 1566\nvalue = [59, 2039, 386]\nnclass = b'),
Text(1488.0, 906.0, 'OXY <= 0.995\ngini = 0.138\nsamples = 1335\nvalue = [27, 1965, 129]\nnclass = b'),
Text(1339.2, 543.599999999999, 'NOx <= 38.315\ngini = 0.365\nsamples = 368\nvalue = [18, 448, 115]\nnclass = b'),
Text(1264.800000000002, 181.1999999999982, 'gini = 0.267\nsamples = 325\nvalue = [9, 431, 70]\nnclass = b'),
Text(1413.600000000001, 181.1999999999982, 'gini = 0.525\nsamples = 43\nvalue = [9, 17, 45]\nnclass = c'),
Text(1636.800000000002, 543.599999999999, 'NO_2 <= 36.945\ngini = 0.03\nsamples = 967\nvalue = [9, 1517, 14]\nnclass = b'),
Text(1562.4, 181.1999999999982, 'gini = 0.022\nsamples = 958\nvalue = [4, 1508, 13]\nnclass = b'),
Text(1711.2, 181.1999999999982, 'gini = 0.524\nsamples = 9\nvalue = [5, 9, 1]\nnclass = b'),
Text(2083.200000000003, 906.0, 'TCH <= 1.235\ngini = 0.449\nsamples = 231\nvalue = [32, 74, 257]\nnclass = c'),
Text(1934.4, 543.599999999999, 'CO <= 0.3\ngini = 0.414\nsamples = 23\nvalue = [29, 0, 12]\nnclass = a'),
Text(1860.000000000002, 181.1999999999982, 'gini = 0.251\nsamples = 18\nvalue = [29, 0, 5]\nnclass = a'),
Text(2008.800000000002, 181.1999999999982, 'gini = 0.0\nsamples = 5\nvalue = [0, 0, 7]\nnclass = c'),
Text(2232.0, 543.599999999999, 'PM10 <= 35.585\ngini = 0.368\nsamples = 208\nvalue = [3, 74, 245]\nnclass = c'),
Text(2157.600000000004, 181.1999999999982, 'gini = 0.237\nsamples = 155\nvalue = [3, 29, 205]\nnclass = c'),
```

```

Text(2306.4, 181.1999999999982, 'gini = 0.498\nsamples = 53\nvalue = [0, 45, 40]\nclas
s = b'),
Text(3292.200000000003, 1630.800000000002, 'OXY <= 0.775\ngini = 0.631\nsamples = 794
4\nvalue = [5302, 2248, 5031]\nklass = a'),
Text(2715.600000000004, 1268.4, 'TCH <= 1.295\ngini = 0.417\nsamples = 900\nvalue = [1
96, 1040, 168]\nklass = b'),
Text(2455.200000000003, 906.0, 'EBE <= 0.615\ngini = 0.534\nsamples = 169\nvalue = [13
5, 120, 10]\nklass = a'),
Text(2380.8, 543.5999999999999, 'gini = 0.0\nsamples = 26\nvalue = [0, 41, 0]\nklass =
b'),
Text(2529.600000000004, 543.5999999999999, 'NO_2 <= 25.105\ngini = 0.51\nsamples = 143
\nvalue = [135, 79, 10]\nklass = a'),
Text(2455.200000000003, 181.1999999999982, 'gini = 0.225\nsamples = 44\nvalue = [8, 5
4, 0]\nklass = b'),
Text(2604.0, 181.1999999999982, 'gini = 0.358\nsamples = 99\nvalue = [127, 25, 10]\ncl
ass = a'),
Text(2976.0, 906.0, 'NOx <= 86.265\ngini = 0.325\nsamples = 731\nvalue = [61, 920, 158]
\nklass = b'),
Text(2827.200000000003, 543.5999999999999, 'PXY <= 0.605\ngini = 0.194\nsamples = 583
\nvalue = [15, 804, 81]\nklass = b'),
Text(2752.8, 181.1999999999982, 'gini = 0.09\nsamples = 392\nvalue = [11, 571, 17]\ncl
ass = b'),
Text(2901.600000000004, 181.1999999999982, 'gini = 0.355\nsamples = 191\nvalue = [4,
233, 64]\nklass = b'),
Text(3124.8, 543.5999999999999, 'O_3 <= 6.88\ngini = 0.624\nsamples = 148\nvalue = [46,
116, 77]\nklass = b'),
Text(3050.4, 181.1999999999982, 'gini = 0.093\nsamples = 54\nvalue = [1, 79, 3]\nklass
= b'),
Text(3199.200000000003, 181.1999999999982, 'gini = 0.636\nsamples = 94\nvalue = [45,
37, 74]\nklass = c'),
Text(3868.8, 1268.4, 'TOL <= 9.565\ngini = 0.59\nsamples = 7044\nvalue = [5106, 1208, 4
863]\nklass = a'),
Text(3571.200000000003, 906.0, 'BEN <= 1.195\ngini = 0.583\nsamples = 4726\nvalue = [2
547, 951, 4040]\nklass = c'),
Text(3422.4, 543.5999999999999, 'TCH <= 1.295\ngini = 0.452\nsamples = 2996\nvalue = [8
30, 568, 3423]\nklass = c'),
Text(3348.000000000005, 181.1999999999982, 'gini = 0.461\nsamples = 508\nvalue = [55
6, 63, 183]\nklass = a'),
Text(3496.8, 181.1999999999982, 'gini = 0.33\nsamples = 2488\nvalue = [274, 505, 3240]
\nklass = c'),
Text(3720.000000000005, 543.5999999999999, 'TCH <= 1.375\ngini = 0.529\nsamples = 1730
\nvalue = [1717, 383, 617]\nklass = a'),
Text(3645.600000000004, 181.1999999999982, 'gini = 0.139\nsamples = 794\nvalue = [115
6, 14, 79]\nklass = a'),
Text(3794.4, 181.1999999999982, 'gini = 0.656\nsamples = 936\nvalue = [561, 369, 538]
\nklass = a'),
Text(4166.400000000001, 906.0, 'PXY <= 2.085\ngini = 0.449\nsamples = 2318\nvalue = [25
59, 257, 823]\nklass = a'),
Text(4017.600000000004, 543.5999999999999, 'BEN <= 1.395\ngini = 0.661\nsamples = 357
\nvalue = [153, 196, 211]\nklass = c'),
Text(3943.200000000003, 181.1999999999982, 'gini = 0.534\nsamples = 168\nvalue = [46,
55, 172]\nklass = c'),
Text(4092.000000000005, 181.1999999999982, 'gini = 0.601\nsamples = 189\nvalue = [10
7, 141, 39]\nklass = b'),
Text(4315.200000000001, 543.5999999999999, 'O_3 <= 13.475\ngini = 0.349\nsamples = 1961
\nvalue = [2406, 61, 612]\nklass = a'),
Text(4240.8, 181.1999999999982, 'gini = 0.409\nsamples = 1240\nvalue = [1410, 33, 507]
\nklass = a'),
Text(4389.6, 181.1999999999982, 'gini = 0.212\nsamples = 721\nvalue = [996, 28, 105]\n
class = a')]

```



Conclusion

Accuracy

Linear Regression:0.3942253498114515

Ridge Regression:0.3923734972503121

Lasso Regression:0.05941509611611351

ElasticNet Regression:0.23414894775926287

Logistic Regression:0.8741416915744405

Random Forest:0.8773225620311598

Accuracy for random forest is higher so it is the best fit model