

Importing Libraries

```
In [69]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

Importing Datasets

```
In [70]: df=pd.read_csv("madrid_2014.csv")
df
```

Out[70]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL
0	2014-06-01 01:00:00	NaN	0.2	NaN	NaN	3.0	10.0	NaN	NaN	NaN	3.0	NaN	NaN
1	2014-06-01 01:00:00	0.2	0.2	0.1	0.11	3.0	17.0	68.0	10.0	5.0	5.0	1.36	1.3
2	2014-06-01 01:00:00	0.3	NaN	0.1	NaN	2.0	6.0	NaN	NaN	NaN	NaN	NaN	1.1
3	2014-06-01 01:00:00	NaN	0.2	NaN	NaN	1.0	6.0	79.0	NaN	NaN	NaN	NaN	28
4	2014-06-01 01:00:00	NaN	NaN	NaN	NaN	1.0	6.0	75.0	NaN	NaN	4.0	NaN	NaN
...
210019	2014-09-01 00:00:00	NaN	0.5	NaN	NaN	20.0	84.0	29.0	NaN	NaN	NaN	NaN	28
210020	2014-09-01 00:00:00	NaN	0.3	NaN	NaN	1.0	22.0	NaN	15.0	NaN	6.0	NaN	NaN
210021	2014-09-01 00:00:00	NaN	NaN	NaN	NaN	1.0	13.0	70.0	NaN	NaN	NaN	NaN	28
210022	2014-09-01 00:00:00	NaN	NaN	NaN	NaN	3.0	38.0	42.0	NaN	NaN	NaN	NaN	28
210023	2014-09-01 00:00:00	NaN	NaN	NaN	NaN	1.0	26.0	65.0	11.0	NaN	NaN	NaN	28

210024 rows × 14 columns

Data Cleaning and Data Preprocessing

```
In [71]: df=df.dropna()
```

```
In [72]: df.columns
```

```
Out[72]: Index(['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',
      'SO_2', 'TCH', 'TOL', 'station'],
      dtype='object')
```

```
In [73]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 13946 entries, 1 to 210006
Data columns (total 14 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   date      13946 non-null   object 
 1   BEN        13946 non-null   float64
 2   CO         13946 non-null   float64
 3   EBE        13946 non-null   float64
 4   NMHC       13946 non-null   float64
 5   NO         13946 non-null   float64
 6   NO_2       13946 non-null   float64
 7   O_3        13946 non-null   float64
 8   PM10       13946 non-null   float64
 9   PM25       13946 non-null   float64
 10  SO_2       13946 non-null   float64
 11  TCH        13946 non-null   float64
 12  TOL        13946 non-null   float64
 13  station    13946 non-null   int64  
dtypes: float64(12), int64(1), object(1)
memory usage: 1.6+ MB
```

```
In [74]: data=df[['CO' , 'station']]  
data
```

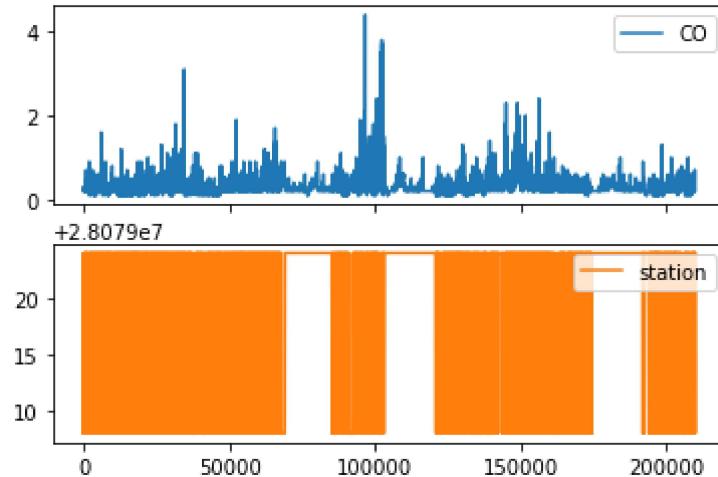
Out[74]:

	CO	station
1	0.2	28079008
6	0.2	28079024
25	0.2	28079008
30	0.2	28079024
49	0.2	28079008
...
209958	0.2	28079024
209977	0.7	28079008
209982	0.2	28079024
210001	0.4	28079008
210006	0.2	28079024

Line chart

```
In [75]: data.plot.line(subplots=True)
```

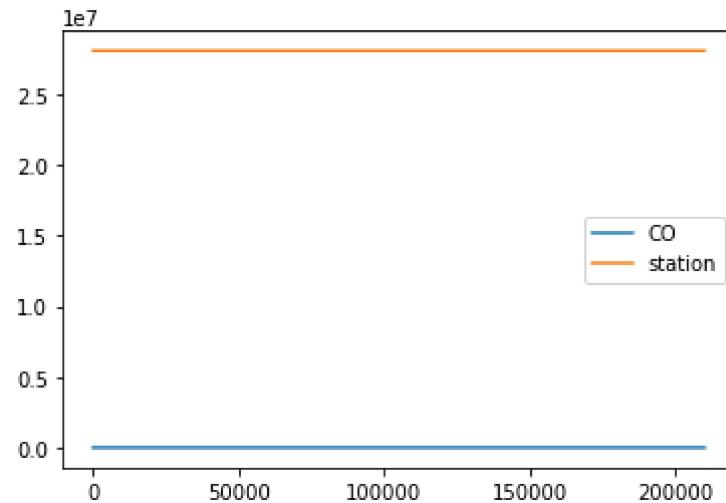
Out[75]: array([<AxesSubplot:>, <AxesSubplot:>], dtype=object)



Line chart

```
In [76]: data.plot.line()
```

```
Out[76]: <AxesSubplot:>
```

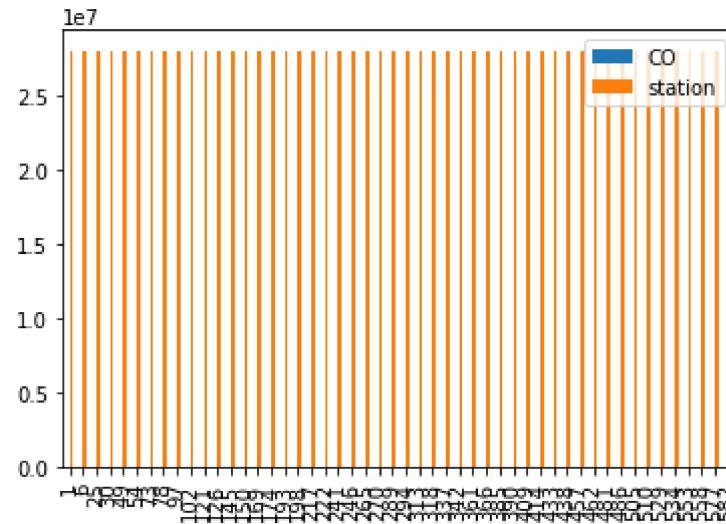


Bar chart

```
In [77]: b=data[0:50]
```

```
In [78]: b.plot.bar()
```

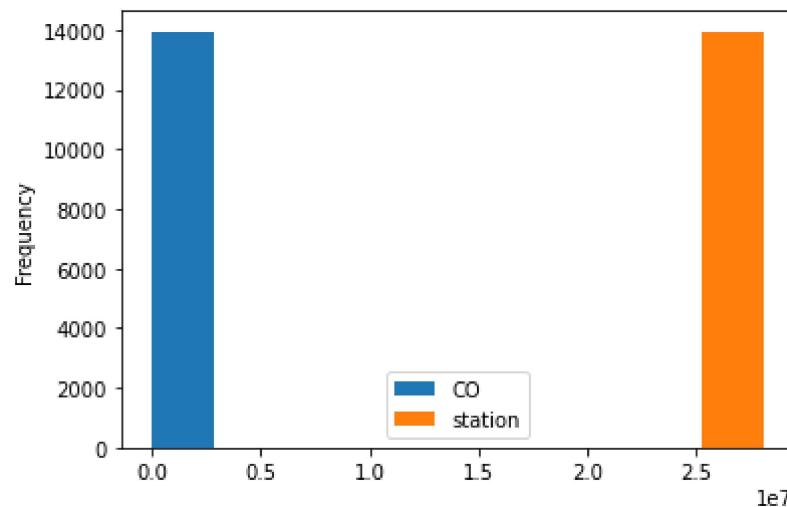
```
Out[78]: <AxesSubplot:>
```



Histogram

```
In [79]: data.plot.hist()
```

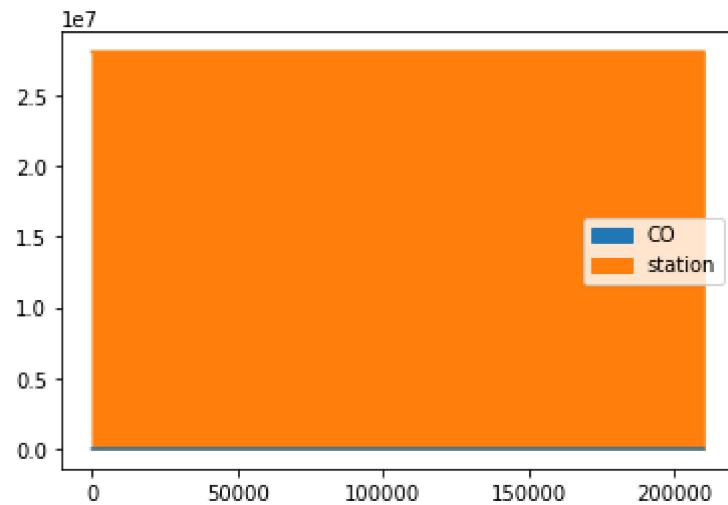
```
Out[79]: <AxesSubplot:ylabel='Frequency'>
```



Area chart

```
In [80]: data.plot.area()
```

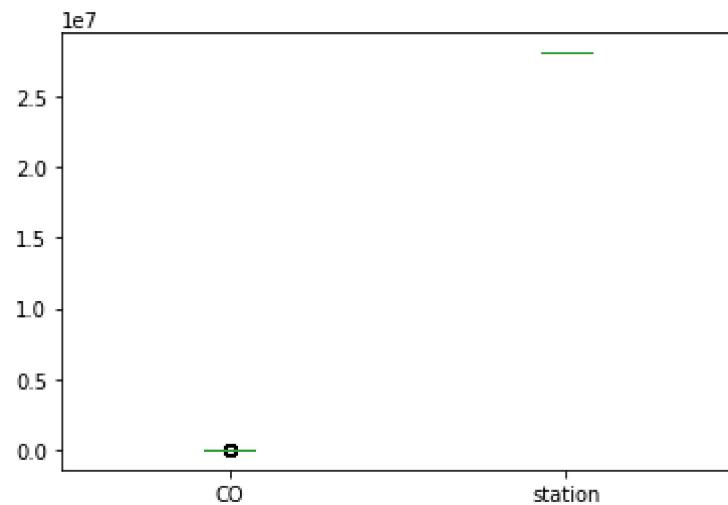
```
Out[80]: <AxesSubplot:>
```



Box chart

```
In [81]: data.plot.box()
```

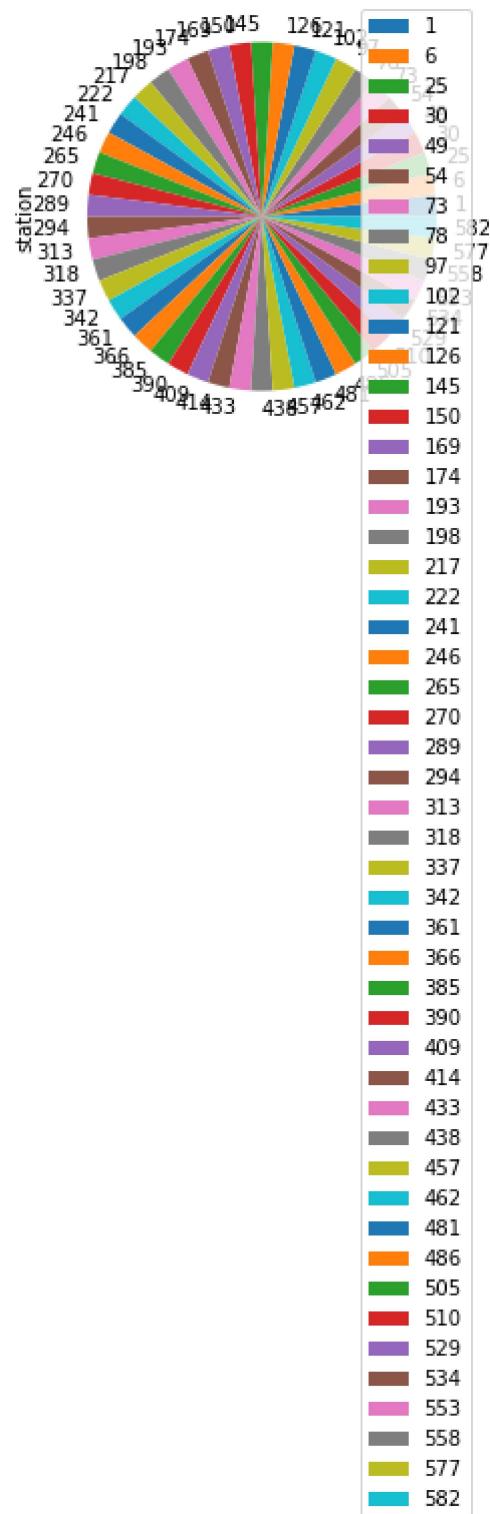
```
Out[81]: <AxesSubplot:>
```



Pie chart

```
In [82]: b.plot.pie(y='station' )
```

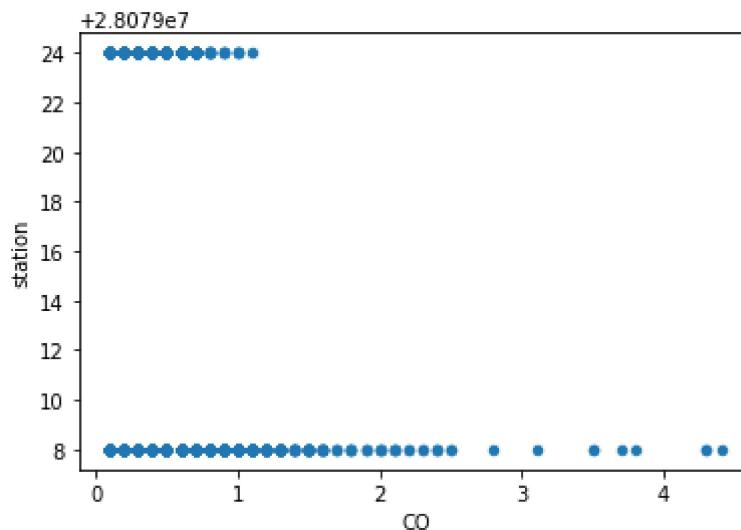
Out[82]: <AxesSubplot:ylabel='station'>



Scatter chart

```
In [83]: data.plot.scatter(x='CO' ,y='station')
```

```
Out[83]: <AxesSubplot:xlabel='CO', ylabel='station'>
```



```
In [84]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 13946 entries, 1 to 210006
Data columns (total 14 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   date      13946 non-null   object 
 1   BEN       13946 non-null   float64
 2   CO        13946 non-null   float64
 3   EBE       13946 non-null   float64
 4   NMHC      13946 non-null   float64
 5   NO        13946 non-null   float64
 6   NO_2      13946 non-null   float64
 7   O_3       13946 non-null   float64
 8   PM10      13946 non-null   float64
 9   PM25      13946 non-null   float64
 10  SO_2      13946 non-null   float64
 11  TCH       13946 non-null   float64
 12  TOL       13946 non-null   float64
 13  station   13946 non-null   int64
```

In [85]: df.describe()

Out[85]:

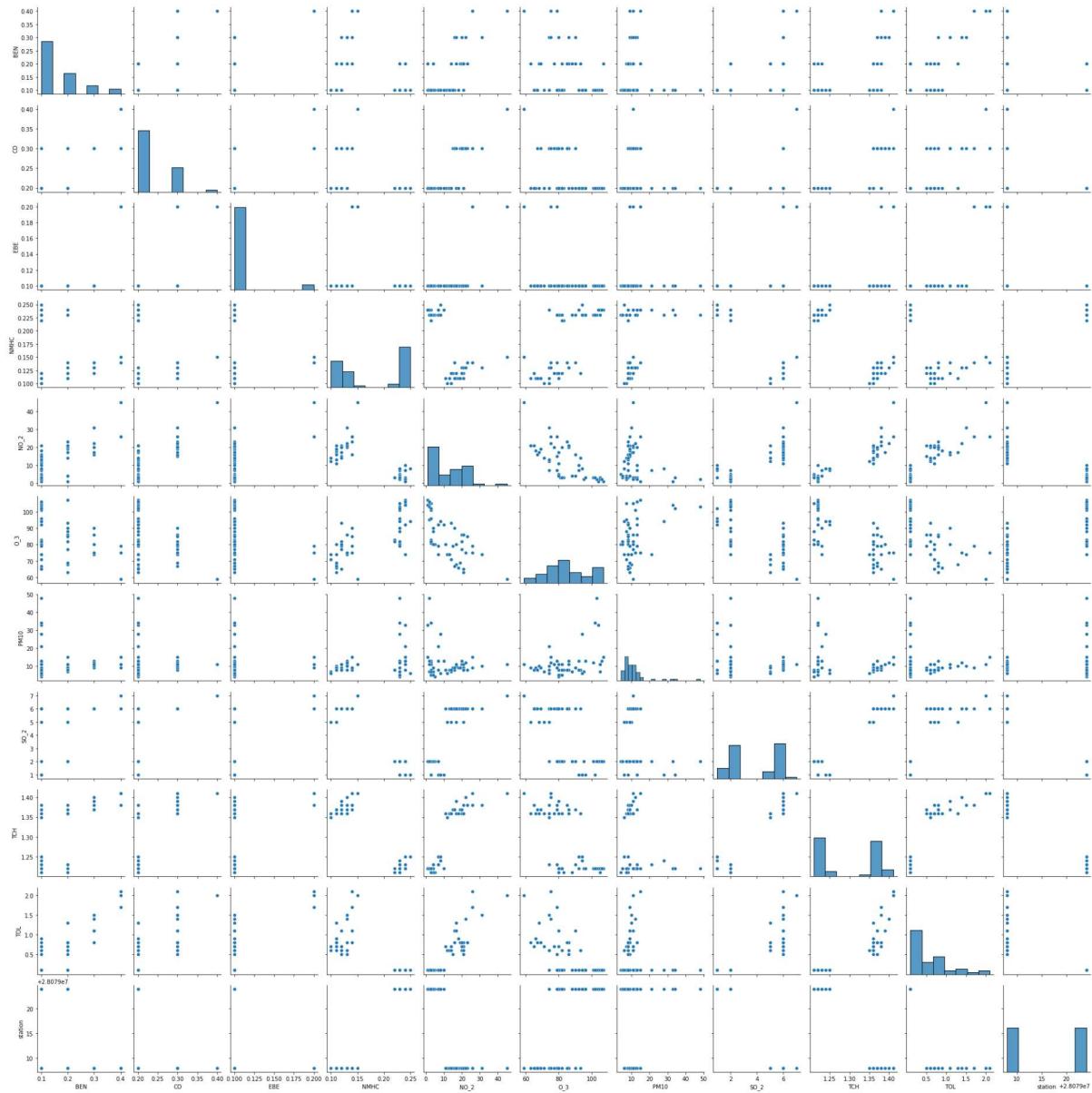
	BEN	CO	EBE	NMHC	NO	NO_2	PM10	SO_2	TCH	TOL	station
count	13946.000000	13946.000000	13946.000000	13946.000000	13946.000000	13946.000000	13946.000000	13946.000000	13946.000000	13946.000000	13946.000000
mean	0.375921	0.314793	0.306016	0.222302	17.589129	34.240929	1.000000	1.000000	1.000000	1.000000	1.000000
std	0.555093	0.207375	0.635475	0.082403	39.432216	30.654229	1.000000	1.000000	1.000000	1.000000	1.000000
min	0.100000	0.100000	0.100000	0.060000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
25%	0.100000	0.200000	0.100000	0.160000	1.000000	10.000000	1.000000	1.000000	1.000000	1.000000	1.000000
50%	0.200000	0.300000	0.100000	0.230000	4.000000	27.000000	1.000000	1.000000	1.000000	1.000000	1.000000
75%	0.400000	0.400000	0.300000	0.260000	18.000000	51.000000	1.000000	1.000000	1.000000	1.000000	1.000000
max	9.400000	4.400000	16.200001	1.290000	725.000000	346.000000	2.000000	2.000000	2.000000	2.000000	2.000000

In [86]: df1=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',
'PM10', 'SO_2', 'TCH', 'TOL', 'station']]

EDA AND VISUALIZATION

```
In [87]: sns.pairplot(df1[0:50])
```

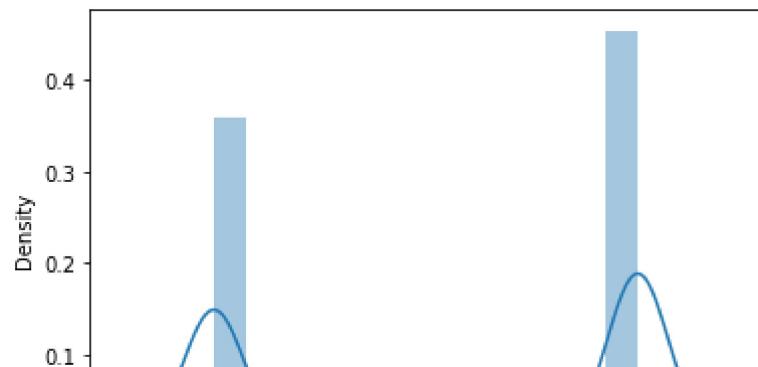
```
Out[87]: <seaborn.axisgrid.PairGrid at 0x22f78c016d0>
```



In [88]: `sns.distplot(df1['station'])`

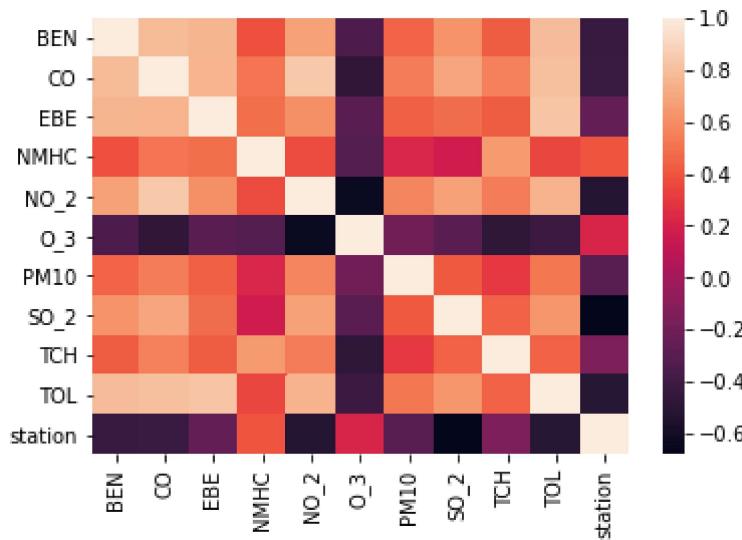
```
C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function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

Out[88]: <AxesSubplot:xlabel='station', ylabel='Density'>



In [89]: `sns.heatmap(df1.corr())`

Out[89]: <AxesSubplot:>



TO TRAIN THE MODEL AND MODEL BUILDING

In [90]: `x=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3', 'PM10', 'SO_2', 'TCH', 'TOL']]
y=df['station']`

```
In [91]: 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 [92]: from sklearn.linear_model import LinearRegression  
lr=LinearRegression()  
lr.fit(x_train,y_train)
```

```
Out[92]: LinearRegression()
```

```
In [93]: lr.intercept_
```

```
Out[93]: 28079022.421686016
```

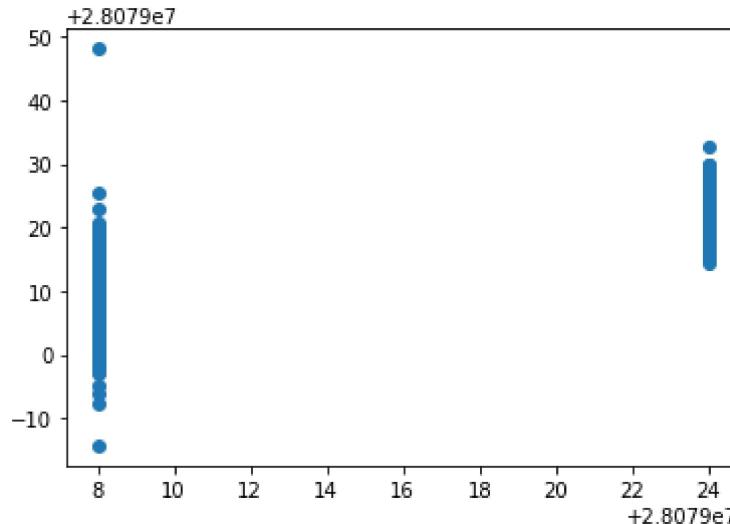
```
In [94]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])  
coeff
```

```
Out[94]:
```

	Co-efficient
BEN	-1.342339
CO	-5.992288
EBE	0.549392
NMHC	83.287967
NO_2	-0.032899
O_3	0.001669
PM10	0.018065
SO_2	-0.871010
TCH	-11.724090
TOL	-0.461004

```
In [95]: prediction = lr.predict(x_test)
plt.scatter(y_test,prediction)
```

```
Out[95]: <matplotlib.collections.PathCollection at 0x22f04167a00>
```



ACCURACY

```
In [96]: lr.score(x_test,y_test)
```

```
Out[96]: 0.8846176160658556
```

```
In [97]: lr.score(x_train,y_train)
```

```
Out[97]: 0.8840230718546449
```

Ridge and Lasso

```
In [98]: from sklearn.linear_model import Ridge,Lasso
```

```
In [99]: rr=Ridge(alpha=10)
rr.fit(x_train,y_train)
```

```
Out[99]: Ridge(alpha=10)
```

Accuracy(Ridge)

```
In [100]: rr.score(x_test,y_test)
```

```
Out[100]: 0.8612062412347833
```

```
In [101]: rr.score(x_train,y_train)
```

```
Out[101]: 0.8609527068908107
```

Accuracy(Lasso)

```
In [102]: la=Lasso(alpha=10)  
la.fit(x_train,y_train)
```

```
Out[102]: Lasso(alpha=10)
```

```
In [103]: la.score(x_train,y_train)
```

```
Out[103]: 0.27130096610743504
```

ElasticNet

```
In [104]: la.score(x_test,y_test)
```

```
Out[104]: 0.2781066212648108
```

```
In [105]: from sklearn.linear_model import ElasticNet  
en=ElasticNet()  
en.fit(x_train,y_train)
```

```
Out[105]: ElasticNet()
```

```
In [106]: en.coef_
```

```
Out[106]: array([ 0.          ,  0.          ,  0.21871313,  0.          , -0.04338765,  
       -0.01130181,  0.02008062, -1.25661205,  0.          , -0.18083073])
```

```
In [107]: en.intercept_
```

```
Out[107]: 28079024.76468402
```

```
In [108]: prediction=en.predict(x_test)
```

```
In [109]: en.score(x_test,y_test)
```

```
Out[109]: 0.47733513340041633
```

Evaluation Metrics

```
In [110]: 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)))
```

4.974607989503637
33.0248511905904
5.746725257970003

Logistic Regression

```
In [111]: from sklearn.linear_model import LogisticRegression
```

```
In [112]: feature_matrix=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',  
    'PM10', 'SO_2', 'TCH', 'TOL']]  
target_vector=df['station']
```

```
In [113]: feature_matrix.shape
```

```
Out[113]: (13946, 10)
```

```
In [114]: target_vector.shape
```

```
Out[114]: (13946,)
```

```
In [115]: from sklearn.preprocessing import StandardScaler
```

```
In [116]: fs=StandardScaler().fit_transform(feature_matrix)
```

```
In [117]: logr=LogisticRegression(max_iter=10000)  
logr.fit(fs,target_vector)
```

```
Out[117]: LogisticRegression(max_iter=10000)
```

```
In [118]: observation=[[1,2,3,4,5,6,7,8,9,10]]
```

```
In [119]: prediction=logr.predict(observation)
```

```
print(prediction)
```

```
[28079008]
```

```
In [120]: logr.classes_
```

```
Out[120]: array([28079008, 28079024], dtype=int64)
```

```
In [121]: logr.score(fs,target_vector)
```

```
Out[121]: 0.9926143697117453
```

```
In [122]: logr.predict_proba(observation)[0][0]
```

```
Out[122]: 1.0
```

```
In [123]: logr.predict_proba(observation)
```

```
Out[123]: array([[1.0, 5.27113072e-18]])
```

Random Forest

```
In [124]: from sklearn.ensemble import RandomForestClassifier
```

```
In [125]: rfc=RandomForestClassifier()  
rfc.fit(x_train,y_train)
```

```
Out[125]: RandomForestClassifier()
```

```
In [126]: parameters={'max_depth':[1,2,3,4,5],  
                     'min_samples_leaf':[5,10,15,20,25],  
                     'n_estimators':[10,20,30,40,50]  
}
```

```
In [127]: 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[127]: 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 [128]: grid_search.best_score_
```

```
Out[128]: 0.9962097930751895
```

```
In [129]: rfc_best=grid_search.best_estimator_
```

```
In [130]: from sklearn.tree import plot_tree
```

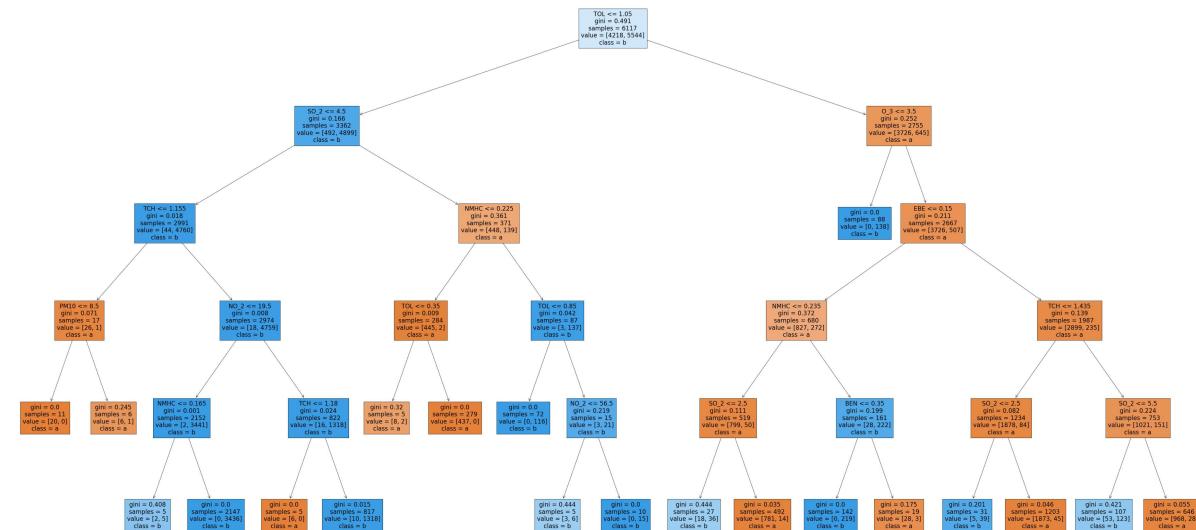
```
plt.figure(figsize=(80,40))
plot_tree(rfc_best.estimators_[5], feature_names=x.columns, class_names=['a', 'b',
```

```
Out[130]: [Text(2247.942857142857, 1993.2, 'TOL <= 1.05\ngini = 0.491\nsamples = 6117\nvalue = [4218, 5544]\nnclass = b'),  
 Text(1179.7714285714285, 1630.8000000000002, 'SO_2 <= 4.5\ngini = 0.166\nsamples = 3362\nvalue = [492, 4899]\nnclass = b'),  
 Text(573.9428571428572, 1268.4, 'TCH <= 1.155\ngini = 0.018\nsamples = 2991\nvalue = [44, 4760]\nnclass = b'),  
 Text(255.0857142857143, 906.0, 'PM10 <= 8.5\ngini = 0.071\nsamples = 17\nvalue = [26, 1]\nnclass = a'),  
 Text(127.54285714285714, 543.5999999999999, 'gini = 0.0\nsamples = 11\nvalue = [20, 0]\nnclass = a'),  
 Text(382.62857142857143, 543.5999999999999, 'gini = 0.245\nsamples = 6\nvalue = [6, 1]\nnclass = a'),  
 Text(892.8, 906.0, 'NO_2 <= 19.5\ngini = 0.008\nsamples = 2974\nvalue = [18, 4759]\nnclass = b'),  
 Text(637.7142857142858, 543.5999999999999, 'NMHC <= 0.165\ngini = 0.001\nsamples = 2152\nvalue = [2, 3441]\nnclass = b'),  
 Text(510.1714285714286, 181.1999999999982, 'gini = 0.408\nsamples = 5\nvalue = [2, 5]\nnclass = b'),  
 Text(765.2571428571429, 181.1999999999982, 'gini = 0.0\nsamples = 2147\nvalue = [0, 3436]\nnclass = b'),  
 Text(1147.8857142857144, 543.5999999999999, 'TCH <= 1.18\ngini = 0.024\nsamples = 822\nvalue = [16, 1318]\nnclass = b'),  
 Text(1020.3428571428572, 181.1999999999982, 'gini = 0.0\nsamples = 5\nvalue = [6, 0]\nnclass = a'),  
 Text(1275.4285714285716, 181.1999999999982, 'gini = 0.015\nsamples = 817\nvalue = [10, 1318]\nnclass = b'),  
 Text(1785.6, 1268.4, 'NMHC <= 0.225\ngini = 0.361\nsamples = 371\nvalue = [448, 139]\nnclass = a'),  
 Text(1530.5142857142857, 906.0, 'TOL <= 0.35\ngini = 0.009\nsamples = 284\nvalue = [445, 2]\nnclass = a'),  
 Text(1402.9714285714285, 543.5999999999999, 'gini = 0.32\nsamples = 5\nvalue = [8, 2]\nnclass = a'),  
 Text(1658.057142857143, 543.5999999999999, 'gini = 0.0\nsamples = 279\nvalue = [437, 0]\nnclass = a'),  
 Text(2040.6857142857143, 906.0, 'TOL <= 0.85\ngini = 0.042\nsamples = 87\nvalue = [3, 137]\nnclass = b'),  
 Text(1913.142857142857, 543.5999999999999, 'gini = 0.0\nsamples = 72\nvalue = [0, 116]\nnclass = b'),  
 Text(2168.2285714285713, 543.5999999999999, 'NO_2 <= 56.5\ngini = 0.219\nsamples = 15\nvalue = [3, 21]\nnclass = b'),  
 Text(2040.6857142857143, 181.1999999999982, 'gini = 0.444\nsamples = 5\nvalue = [3, 6]\nnclass = b'),  
 Text(2295.7714285714287, 181.1999999999982, 'gini = 0.0\nsamples = 10\nvalue = [0, 15]\nnclass = b'),  
 Text(3316.114285714286, 1630.8000000000002, 'O_3 <= 3.5\ngini = 0.252\nsamples = 2755\nvalue = [3726, 645]\nnclass = a'),  
 Text(3188.5714285714284, 1268.4, 'gini = 0.0\nsamples = 88\nvalue = [0, 138]\nnclass = b'),  
 Text(3443.657142857143, 1268.4, 'EBE <= 0.15\ngini = 0.211\nsamples = 2667\nvalue = [3726, 507]\nnclass = a'),  
 Text(2933.4857142857145, 906.0, 'NMHC <= 0.235\ngini = 0.372\nsamples = 680\nvalue = [827, 272]\nnclass = a'),  
 Text(2678.4, 543.5999999999999, 'SO_2 <= 2.5\ngini = 0.111\nsamples = 519\nvalue = [799, 50]\nnclass = a'),  
 Text(2550.857142857143, 181.1999999999982, 'gini = 0.444\nsamples = 27\nvalue = [18, 36]\nnclass = b'),  
 Text(2805.942857142857, 181.1999999999982, 'gini = 0.035\nsamples = 492\nvalue = [18, 36]\nnclass = b')]
```

```

lue = [781, 14]\nclass = a'),
Text(3188.5714285714284, 543.5999999999999, 'BEN <= 0.35\ngini = 0.199\nsamples = 161\nvalue = [28, 222]\nclass = b'),
Text(3061.0285714285715, 181.1999999999982, 'gini = 0.0\nnsamples = 142\nvalue = [0, 219]\nclass = b'),
Text(3316.114285714286, 181.1999999999982, 'gini = 0.175\nnsamples = 19\nvalue = [28, 3]\nclass = a'),
Text(3953.8285714285716, 906.0, 'TCH <= 1.435\ngini = 0.139\nnsamples = 1987\nvalue = [2899, 235]\nclass = a'),
Text(3698.7428571428572, 543.5999999999999, 'SO_2 <= 2.5\ngini = 0.082\nsamples = 1234\nvalue = [1878, 84]\nclass = a'),
Text(3571.2, 181.1999999999982, 'gini = 0.201\nnsamples = 31\nvalue = [5, 39]\nclass = b'),
Text(3826.285714285714, 181.1999999999982, 'gini = 0.046\nnsamples = 1203\nvalue = [1873, 45]\nclass = a'),
Text(4208.914285714286, 543.5999999999999, 'SO_2 <= 5.5\ngini = 0.224\nsamples = 753\nvalue = [1021, 151]\nclass = a'),
Text(4081.3714285714286, 181.1999999999982, 'gini = 0.421\nnsamples = 107\nvalue = [53, 123]\nclass = b'),
Text(4336.457142857143, 181.1999999999982, 'gini = 0.055\nnsamples = 646\nvalue = [968, 28]\nclass = a')

```



Conclusion

Accuracy

In [131]: `lr.score(x_train,y_train)`

Out[131]: 0.8840230718546449

In [132]: `rr.score(x_train,y_train)`

Out[132]: 0.8609527068908107

```
In [133]: la.score(x_train,y_train)
```

```
Out[133]: 0.27130096610743504
```

```
In [134]: en.score(x_test,y_test)
```

```
Out[134]: 0.47733513340041633
```

```
In [135]: logr.score(fs,target_vector)
```

```
Out[135]: 0.9926143697117453
```

```
In [136]: grid_search.best_score_
```

```
Out[136]: 0.9962097930751895
```

Random forest is the suitable for this dataset