

# 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\C10_air\csvs_per_year\csvs(Dataset)\madrid_2016.
df
```

Out[2]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOI
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	14.0
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]:

```
df=df.dropna()
```

In [4]:

```
df.columns
```

Out[4]:

```
Index(['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'P  
M25',  
      'SO_2', 'TCH', 'TOL', 'station'],  
      dtype='object')
```

In [5]:

```
df.info()
```

```
<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[['BEN', 'TOL', 'TCH']]
data
```

Out[6]:

	BEN	TOL	TCH
1	3.1	14.4	2.44
6	0.7	5.0	1.35
25	2.7	15.0	2.30
30	0.7	5.0	1.35
49	1.7	10.7	1.95
...	...	...	...
209430	0.1	0.2	1.15
209449	0.6	1.9	1.48
209454	0.1	0.3	1.15
209473	0.6	1.9	1.50
209478	0.1	0.2	1.15

16932 rows × 3 columns

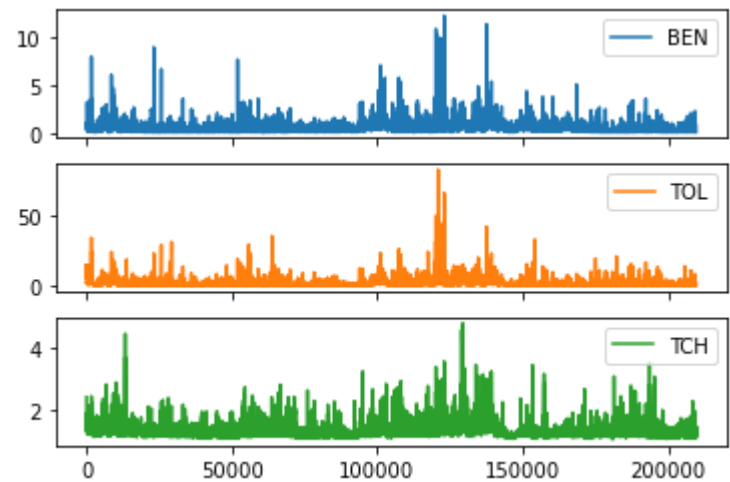
## Line chart

In [7]:

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

Out[7]:

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



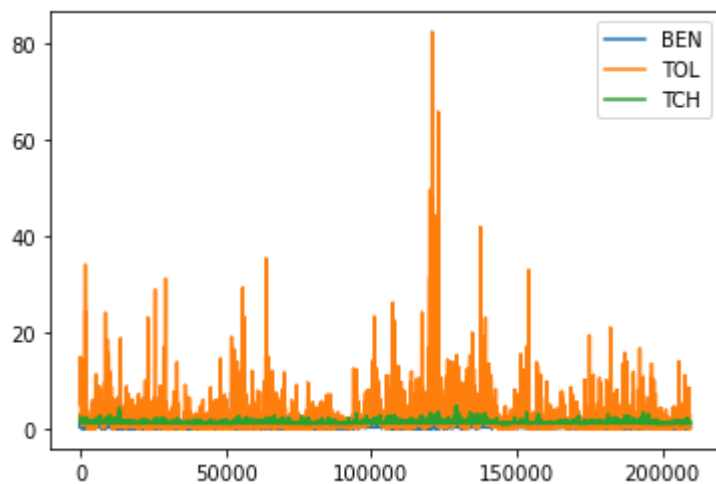
## Line chart

In [8]:

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

Out[8]:

&lt;AxesSubplot:&gt;



## Bar chart

In [9]:

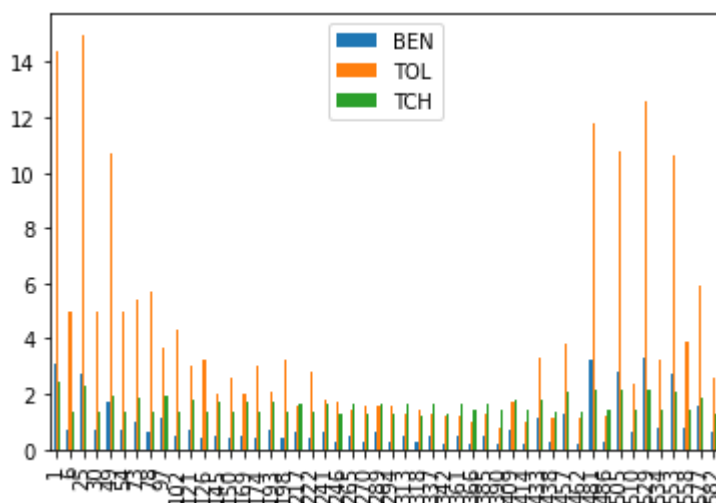
```
b=data[0:50]
```

In [10]:

```
b.plot.bar()
```

Out[10]:

&lt;AxesSubplot:&gt;



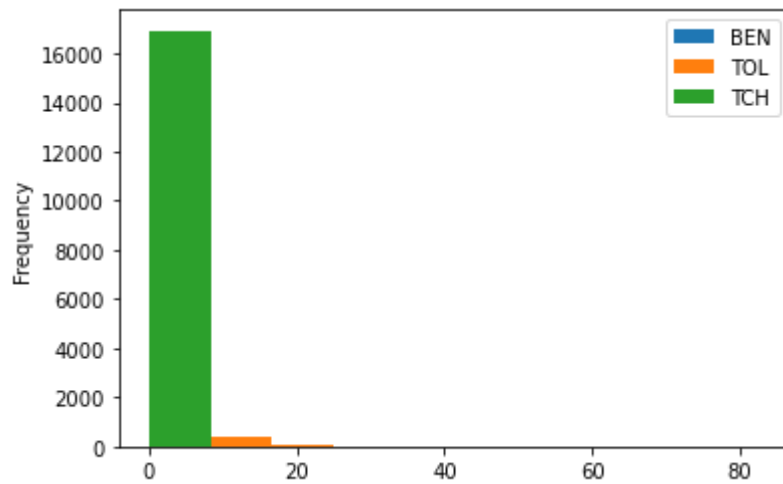
## Histogram

In [11]:

```
data.plot.hist()
```

Out[11]:

<AxesSubplot:ylabel='Frequency'>



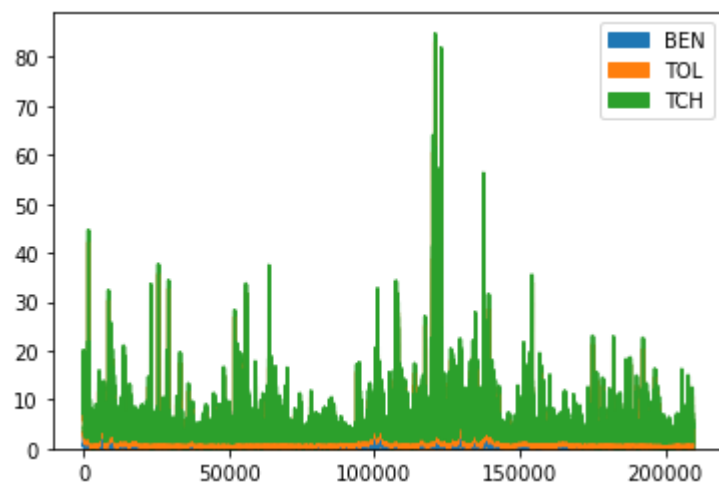
## Area chart

In [12]:

```
data.plot.area()
```

Out[12]:

<AxesSubplot:>



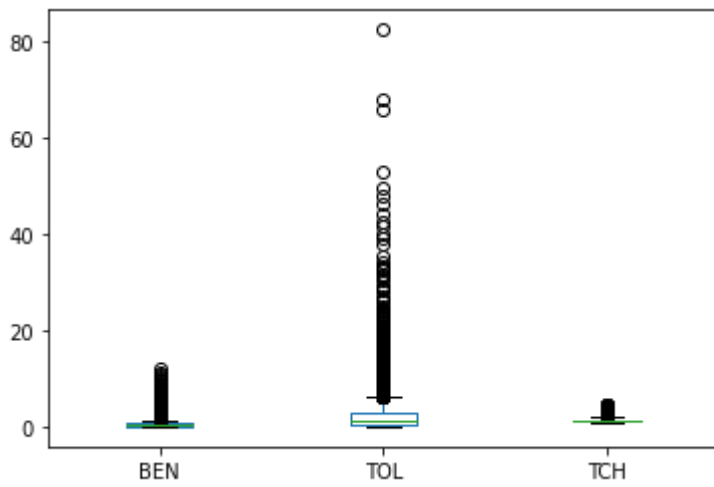
## Box chart

In [13]:

```
data.plot.box()
```

Out[13]:

<AxesSubplot:>



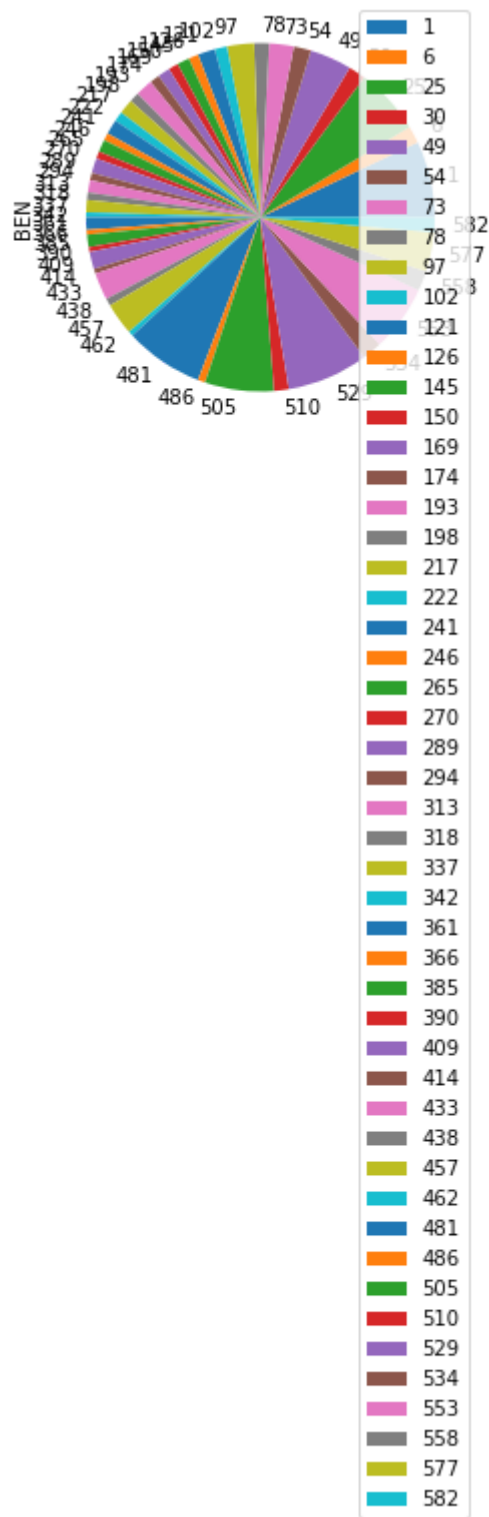
## Pie chart

In [14]:

```
b.plot.pie(y='BEN' )
```

Out[14]:

<AxesSubplot:ylabel='BEN'>



# Scatter chart

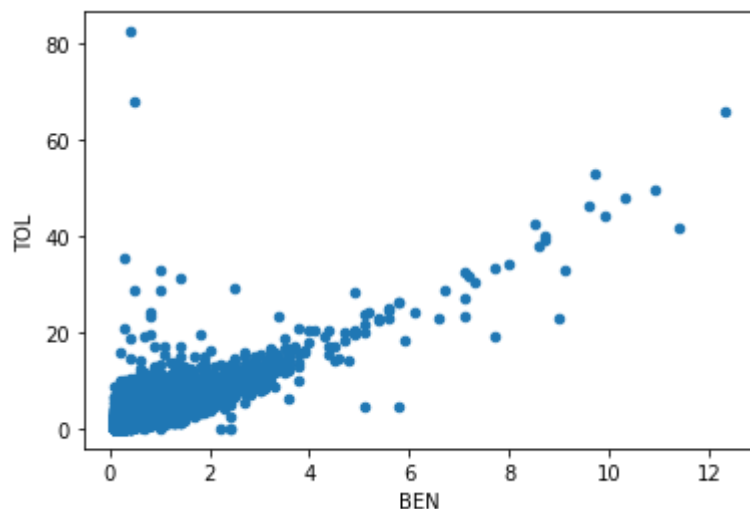


In [15]:

```
data.plot.scatter(x='BEN' ,y='TOL')
```

Out[15]:

<AxesSubplot:xlabel='BEN', ylabel='TOL'>



In [16]:

```
df.info()
```

```
<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 [17]:

```
df.describe()
```

Out[17]:

	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

In [18]:

```
df1=df[['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',  
        'SO_2', 'TCH', 'TOL', 'station']]
```

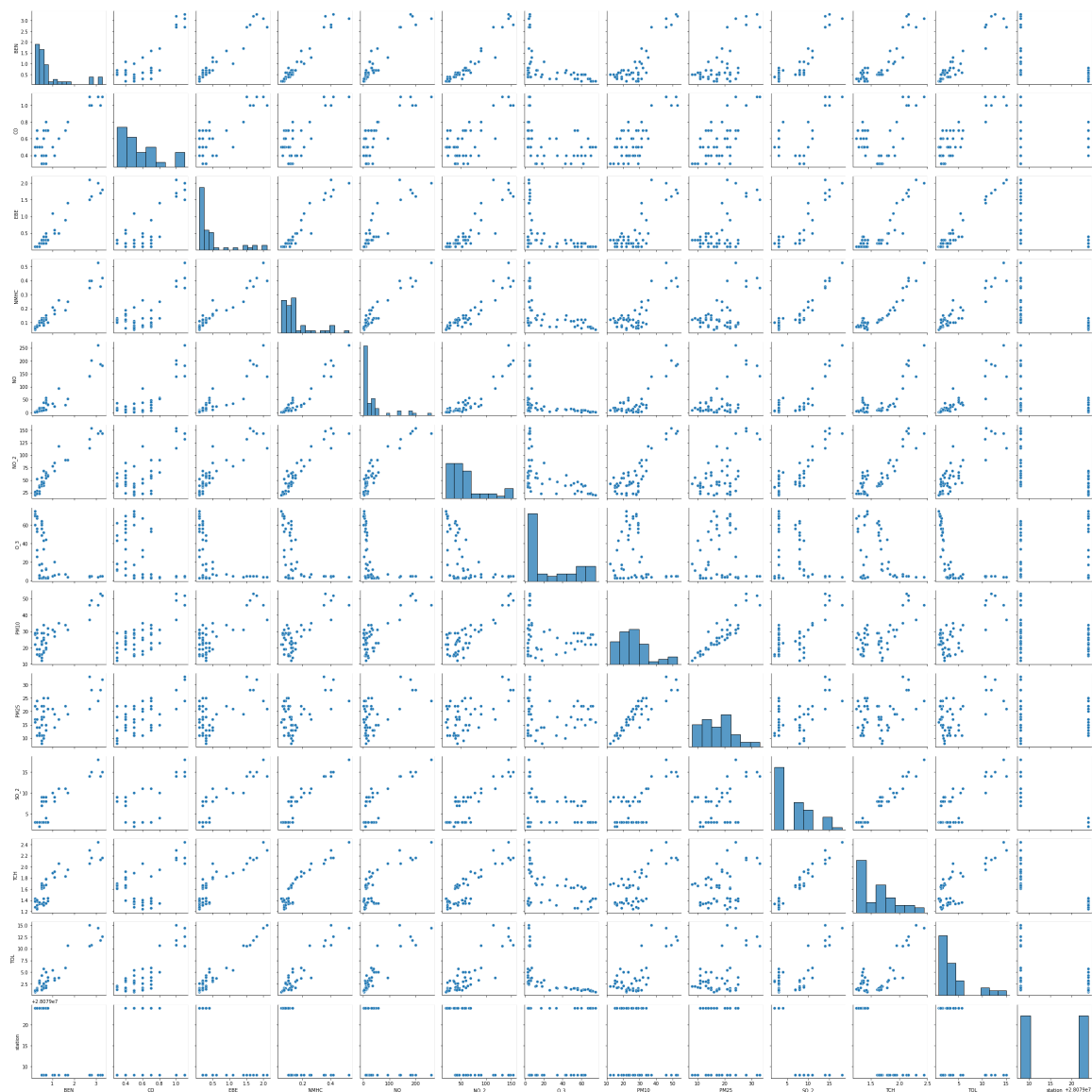
## EDA AND VISUALIZATION

In [19]:

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

Out[19]:

&lt;seaborn.axisgrid.PairGrid at 0x1b232bba220&gt;



In [20]:

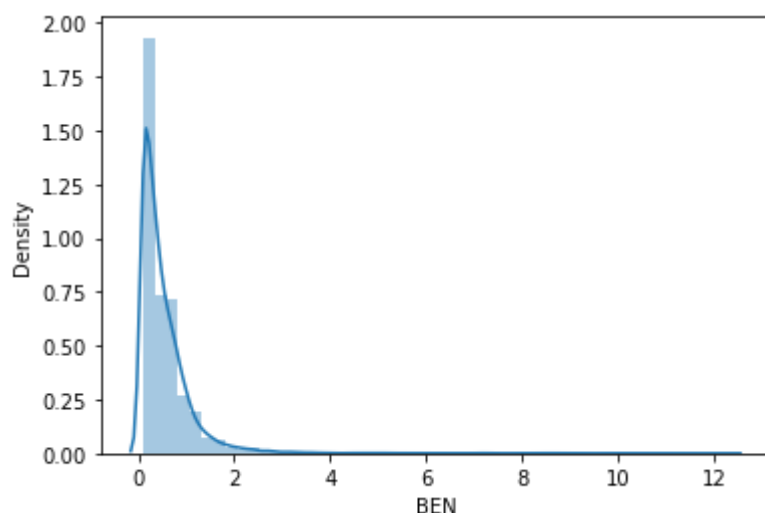
```
sns.distplot(df1['BEN'])
```

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='BEN', ylabel='Density'>
```

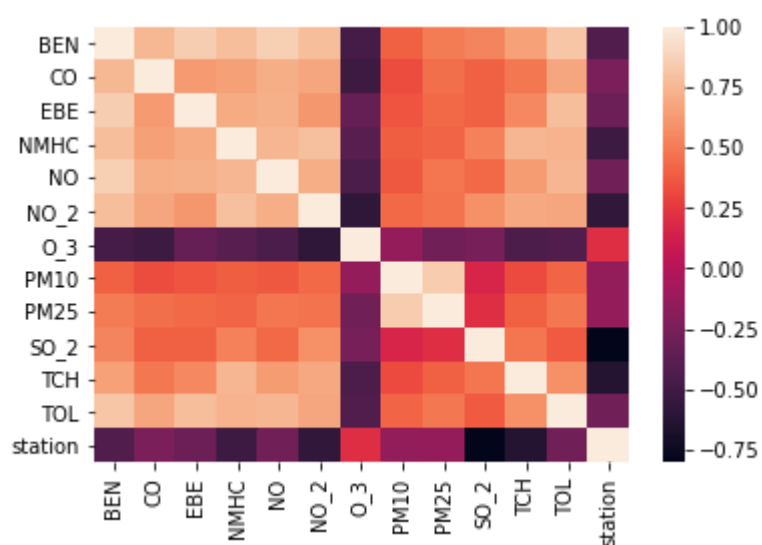


In [21]:

```
sns.heatmap(df1.corr())
```

Out[21]:

```
<AxesSubplot:>
```



## TO TRAIN THE MODEL AND MODEL BULDING

In [22]:

```
x=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',  
      '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]:

28079041.999899916

In [26]:

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

Out[26]:

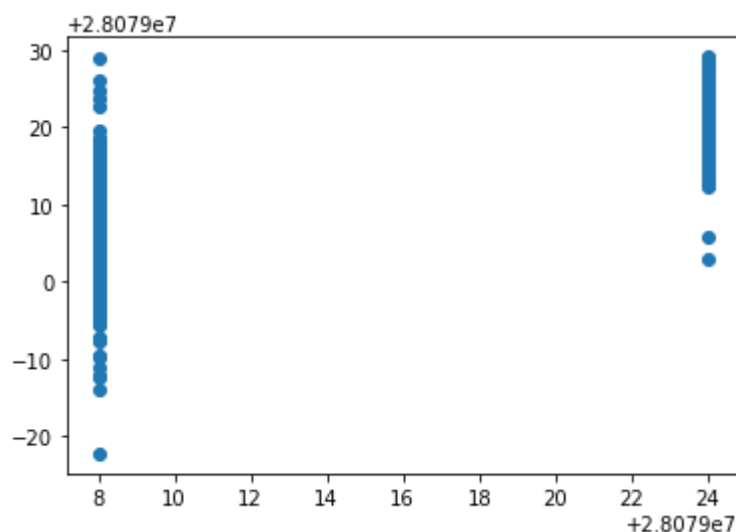
Co-efficient	
<b>BEN</b>	-1.807121
<b>CO</b>	4.467319
<b>EBE</b>	0.581989
<b>NMHC</b>	1.262170
<b>NO</b>	0.067188
<b>NO_2</b>	-0.067192
<b>O_3</b>	-0.024761
<b>PM10</b>	-0.012139
<b>PM25</b>	0.103008
<b>SO_2</b>	-0.803779
<b>TCH</b>	-13.998236
<b>TOL</b>	0.205204

In [27]:

```
prediction = lr.predict(x_test)
plt.scatter(y_test, prediction)
```

Out[27]:

<matplotlib.collections.PathCollection at 0x1b23ec55f70>



## ACCURACY

In [28]:

```
lr.score(x_test, y_test)
```

Out[28]:

0.826130159659848

In [29]:

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

Out[29]:

0.8286608817866095

## 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]:

```
0.8256727930423138
```

In [33]:

```
rr.score(x_train,y_train)
```

Out[33]:

```
0.828577461543129
```

In [34]:

```
la=Lasso(alpha=10)  
la.fit(x_train,y_train)
```

Out[34]:

```
Lasso(alpha=10)
```

In [35]:

```
la.score(x_test,y_test)
```

Out[35]:

```
0.6431475578851213
```

## Accuracy(Lasso)

In [36]:

```
la.score(x_train,y_train)
```

Out[36]:

```
0.6467262929917432
```

## Accuracy(Elastic Net)

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([-0.          ,  0.          , -0.          , -0.          ,  0.048158   ,
        -0.10696444, -0.02103236,  0.00315771,  0.04814422, -0.86268439,
        -0.03113773,  0.          ])
```

In [39]:

```
en.intercept_
```

Out[39]:

```
28079026.225688875
```

In [40]:

```
prediction=en.predict(x_test)
```

In [41]:

```
en.score(x_test,y_test)
```

Out[41]:

```
0.7049603330703134
```

## 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)))
```

```
3.352243386644606
```

```
18.8716480631203
```

```
4.344151017531538
```

## Logistic Regression

In [43]:

```
from sklearn.linear_model import LogisticRegression
```

In [44]:

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



In [45]:

```
feature_matrix.shape
```

Out[45]:

```
(16932, 10)
```

In [46]:

```
target_vector.shape
```

Out[46]:

```
(16932,)
```

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]]
```

In [51]:

```
prediction=logr.predict(observation)  
print(prediction)
```

```
[28079008]
```

In [52]:

```
logr.classes_
```

Out[52]:

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

In [53]:

```
logr.score(fs,target_vector)
```

Out[53]:

```
0.9923812898653437
```

In [54]:

```
logr.predict_proba(observation)[0][0]
```

Out[54]:

1.0

In [55]:

```
logr.predict_proba(observation)
```

Out[55]:

```
array([[1.0000000e+00, 1.6336121e-46]])
```

## 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]:

0.9946000674991562

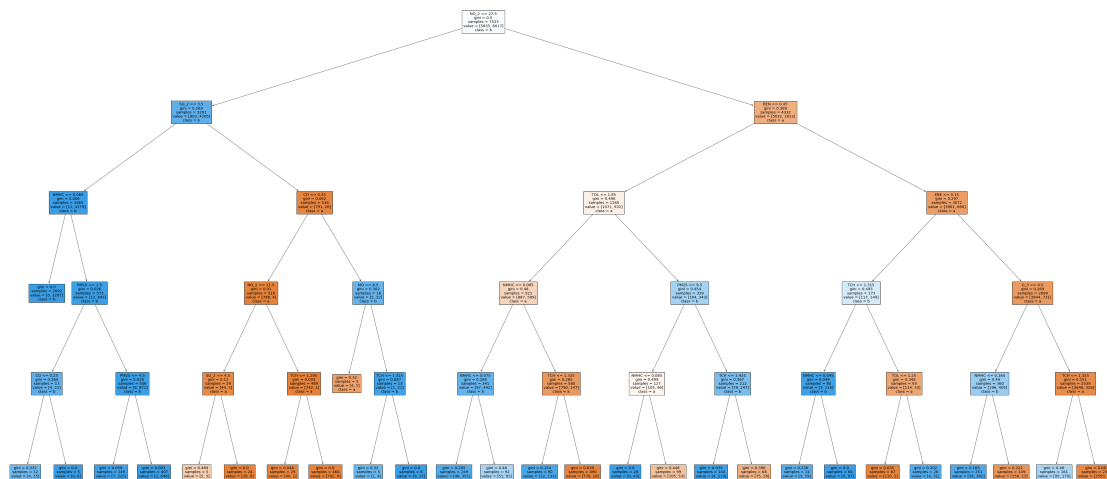
In [61]:

```
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'],f
\nvalue = [95, 170]\nnclass = b'),
Text(4378.153846153846, 181.19999999999982, 'gini = 0.081\nsamples = 23
78\nvalue = [3553, 156]\nnclass = a'))
```



## Conclusion

### Accuracy

**Linear Regression:0.8286608817866095**

**Ridge Regression:0.828577461543129**

**Lasso Regression:0.6467262929917432**

**ElasticNet Regression:0.7049603330703134**

**Logistic Regression:0.9923812898653437**

**Random Forest:0.9946000674991562**

**Random Forest is suitable for this dataset**

