

# 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_2002.
df
```

Out[2]:

	date	BEN	CO	EBE	MXY	NMHC	NO_2	NOx	OXY	O_3	P
0	2002-04-01 01:00:00	NaN	1.39	NaN	NaN	NaN	145.100006	352.100006	NaN	6.54	41.990
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	20.980
2	2002-04-01 01:00:00	NaN	0.80	NaN	NaN	NaN	103.699997	134.000000	NaN	13.01	28.440
3	2002-04-01 01:00:00	NaN	1.61	NaN	NaN	NaN	97.599998	268.000000	NaN	5.12	42.180
4	2002-04-01 01:00:00	NaN	1.90	NaN	NaN	NaN	92.089996	237.199997	NaN	7.28	76.330
...	...	...	...	...	...	...	...	...	...	...	...
217291	2002-11-01 00:00:00	4.16	1.14	NaN	NaN	NaN	81.080002	265.700012	NaN	7.21	36.750
217292	2002-11-01 00:00:00	3.67	1.73	2.89	NaN	0.38	113.900002	373.100006	NaN	5.66	63.380
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	9.640
217294	2002-11-01 00:00:00	4.51	0.91	4.83	10.99	NaN	149.800003	202.199997	1.00	5.75	1
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	29.240

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 [11]:

```
data=df[['BEN', 'TOL', 'PXY']]
data
```

Out[11]:

	BEN	TOL	PXY
1	1.93	10.98	2.53
5	3.19	15.60	2.98
22	2.02	7.32	1.48
24	3.02	11.42	2.18
26	2.02	10.60	2.45
...	...	...	...
217269	1.24	4.45	0.94
217271	3.13	15.10	3.40
217273	2.50	16.65	3.60
217293	1.37	4.33	0.94
217295	3.11	15.51	3.35

32381 rows × 3 columns

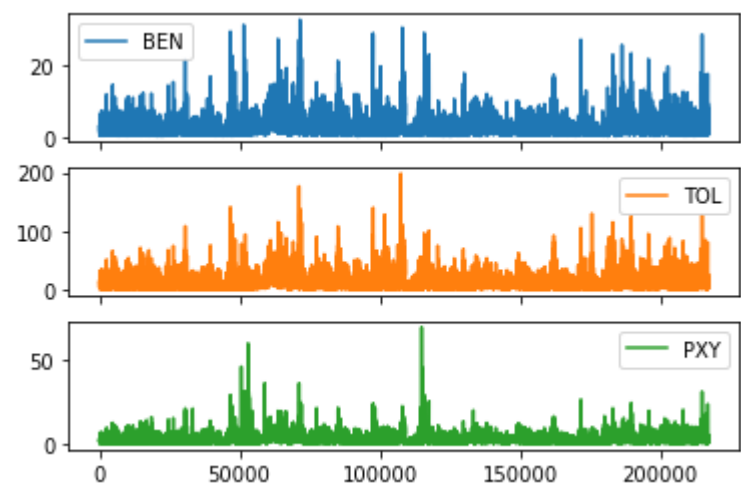
## Line chart

In [12]:

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

Out[12]:

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



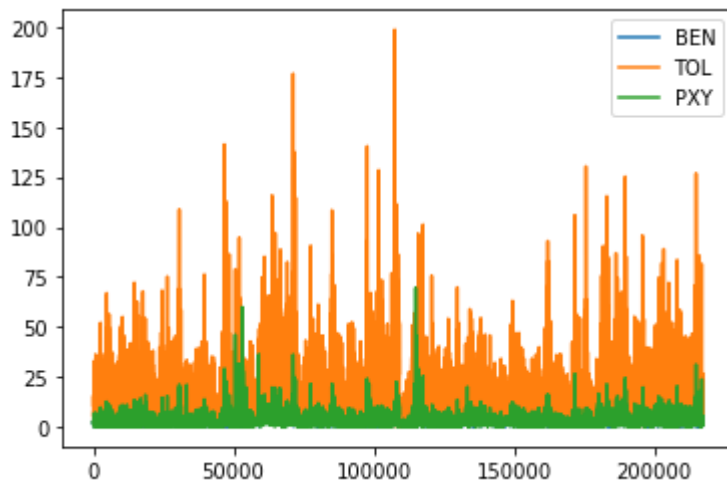
## Line chart

In [13]:

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

Out[13]:

<AxesSubplot:>



## Bar chart

In [14]:

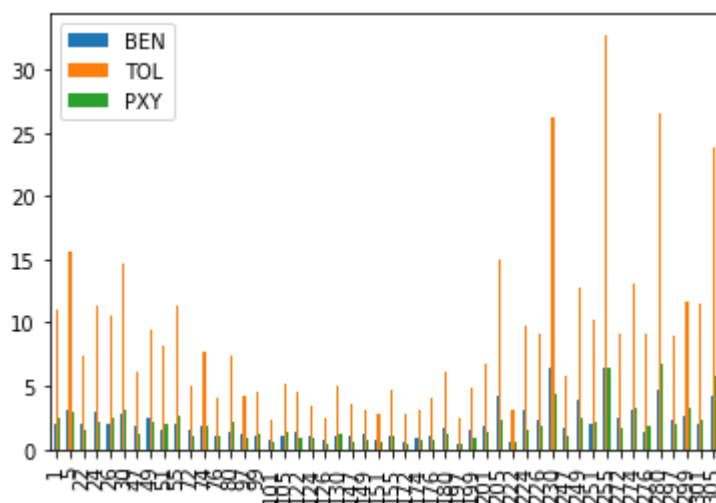
```
b=data[0:50]
```

In [15]:

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

Out[15]:

<AxesSubplot:>



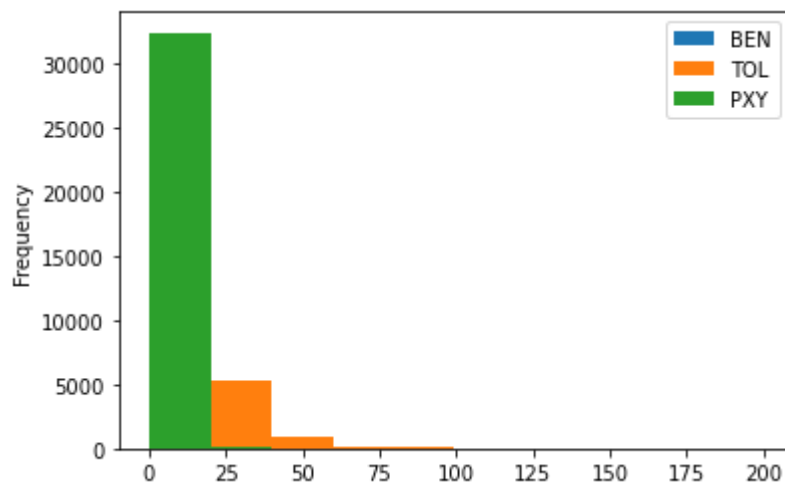
## Histogram

In [16]:

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

Out[16]:

<AxesSubplot:ylabel='Frequency'>



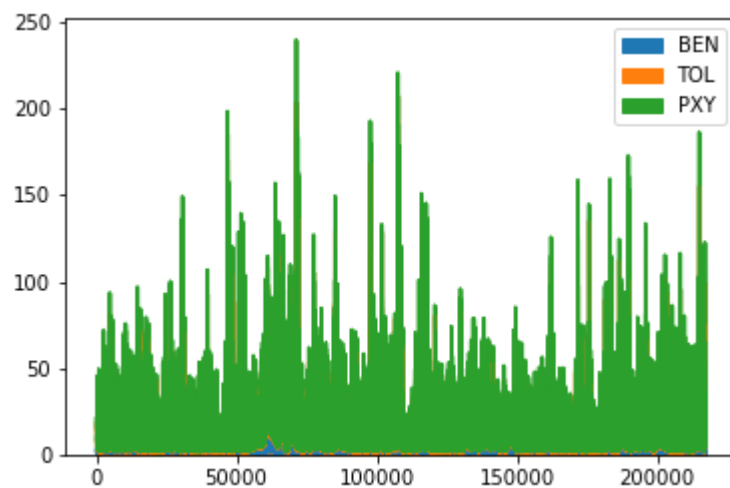
## Area chart

In [17]:

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

Out[17]:

<AxesSubplot:>



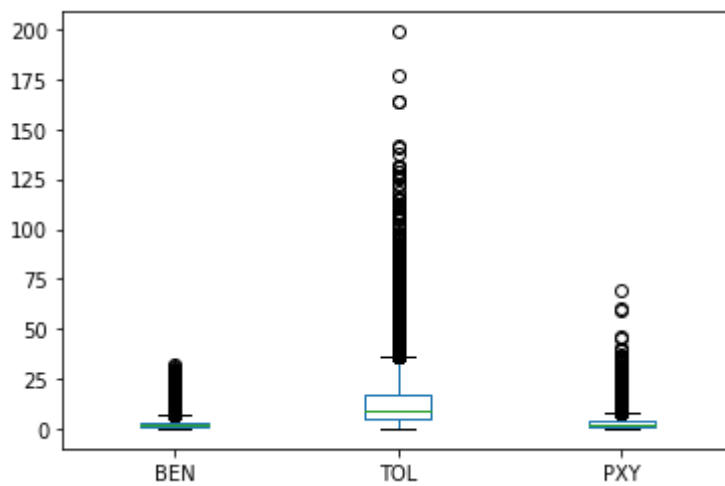
## Box chart

In [18]:

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

Out[18]:

<AxesSubplot:>



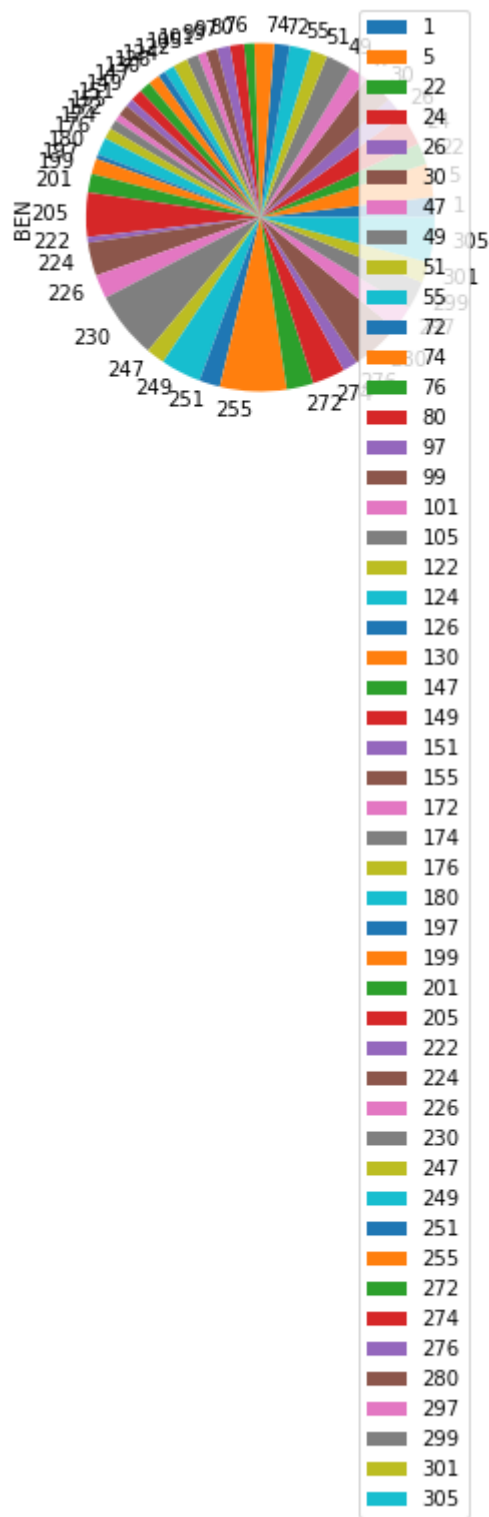
## Pie chart

In [20]:

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

Out[20]:

<AxesSubplot:ylabel='BEN'>



# Scatter chart

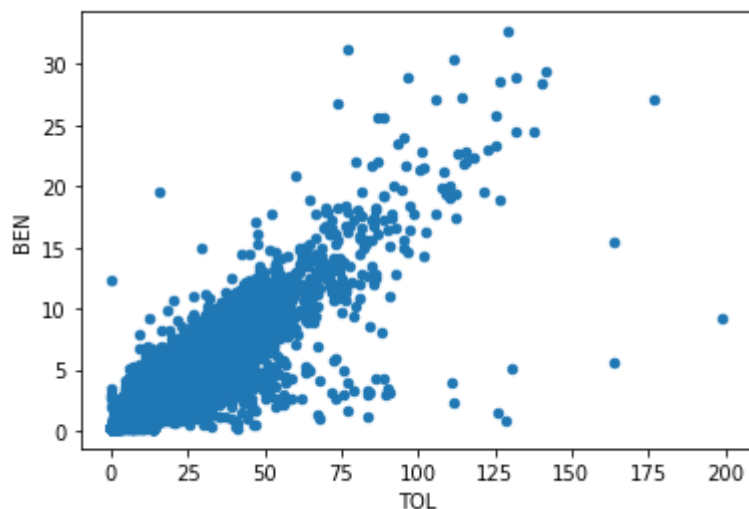


In [22]:

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

Out[22]:

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



In [23]:

```
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 [24]:

```
df.describe()
```

Out[24]:

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

In [25]:

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

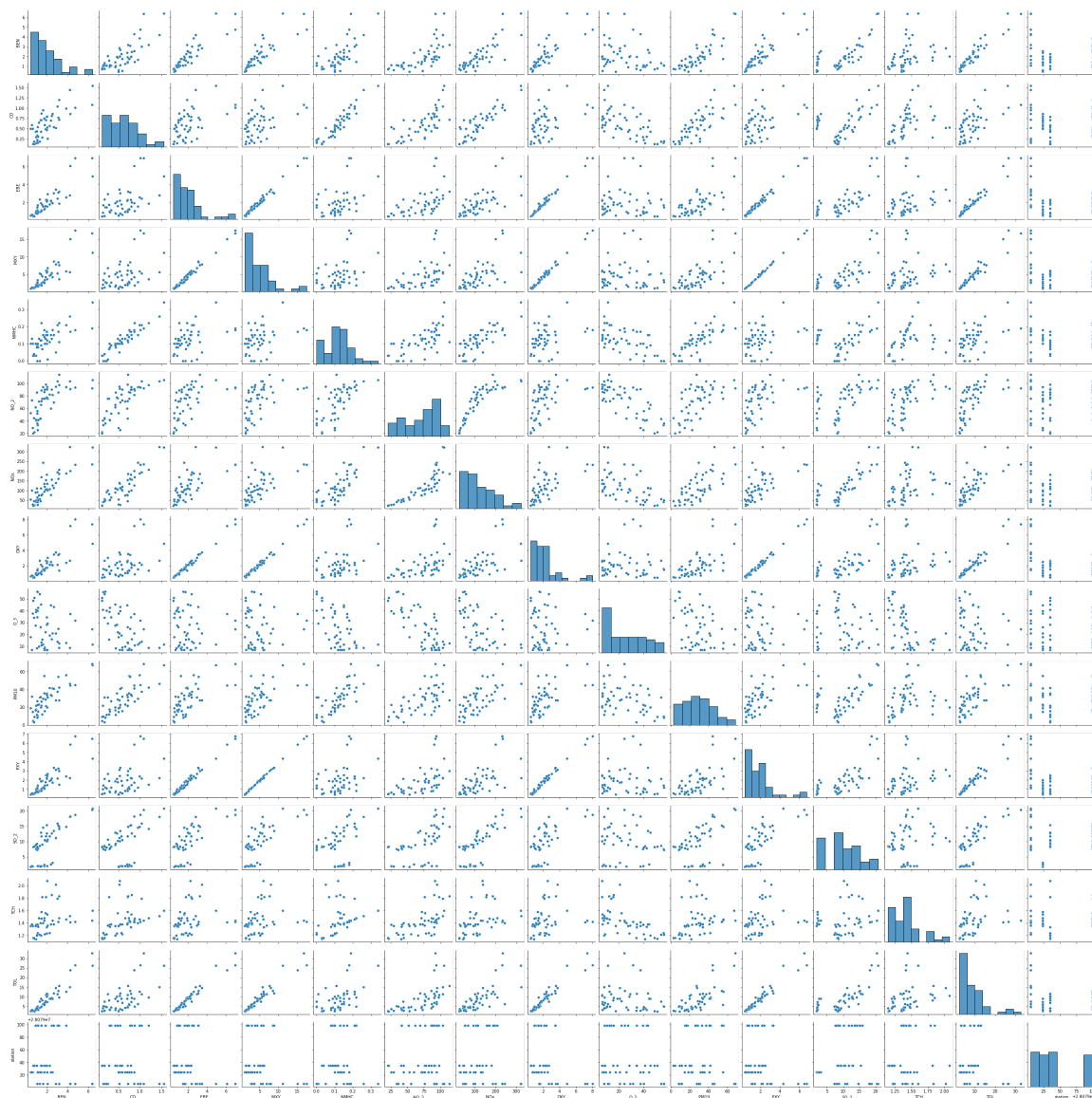
# EDA AND VISUALIZATION

In [26]:

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

Out[26]:

&lt;seaborn.axisgrid.PairGrid at 0x11456bcd30&gt;



In [27]:

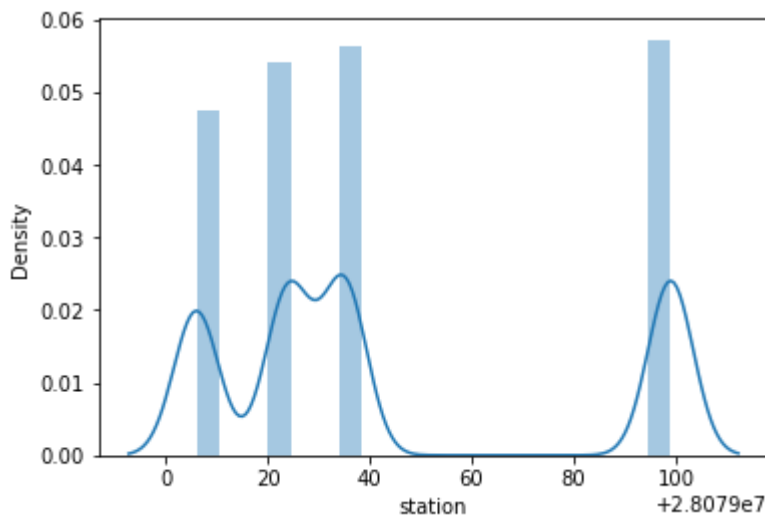
```
sns.distplot(df1['station'])
```

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[27]:

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

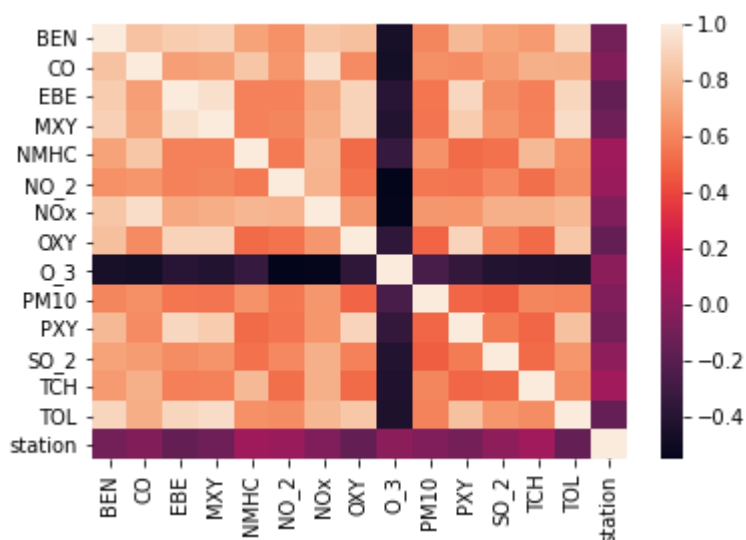


In [28]:

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

Out[28]:

```
<AxesSubplot:>
```



## TO TRAIN THE MODEL AND MODEL BUILDING

In [29]:

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

In [30]:

```
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 [31]:

```
from sklearn.linear_model import LinearRegression  
lr=LinearRegression()  
lr.fit(x_train,y_train)
```

Out[31]:

LinearRegression()

In [32]:

```
lr.intercept_
```

Out[32]:

28078991.39595005

In [33]:

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

Out[33]:

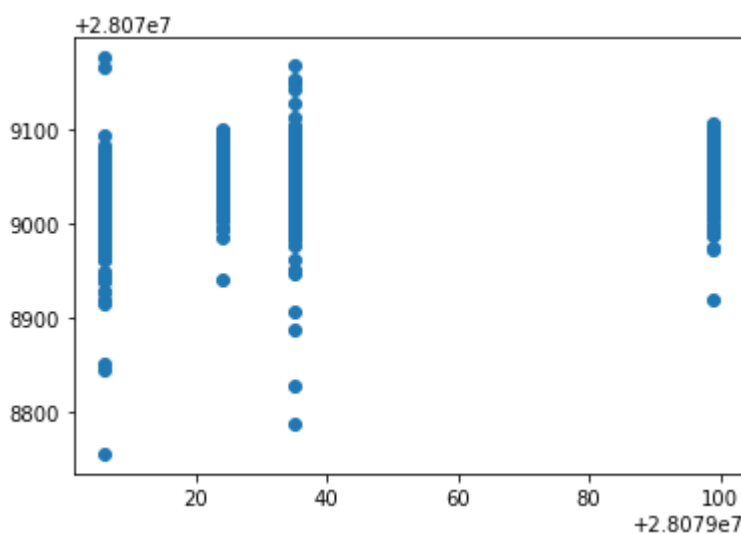
	Co-efficient
<b>BEN</b>	1.974642
<b>CO</b>	-13.144500
<b>EBE</b>	-12.590009
<b>MXY</b>	4.183719
<b>NMHC</b>	80.220883
<b>NO_2</b>	0.254409
<b>NOx</b>	-0.094810
<b>OXY</b>	-4.793770
<b>O_3</b>	-0.039461
<b>PM10</b>	-0.114979
<b>PXY</b>	8.165716
<b>SO_2</b>	0.551567
<b>TCH</b>	41.115474
<b>TOL</b>	-1.451087

In [34]:

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

Out[34]:

&lt;matplotlib.collections.PathCollection at 0x11467f38940&gt;



## ACCURACY

In [35]:

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

Out[35]:

0.19727401561681268

In [36]:

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

Out[36]:

0.19864042035800789

## Ridge and Lasso

In [37]:

```
from sklearn.linear_model import Ridge,Lasso
```

In [38]:

```
rr=Ridge(alpha=10)  
rr.fit(x_train,y_train)
```

Out[38]:

Ridge(alpha=10)

## Accuracy(Ridge)

In [39]:

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

Out[39]:

0.1961181197383136

In [40]:

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

Out[40]:

0.1984498616765833

In [41]:

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

Out[41]:

Lasso(alpha=10)

In [42]:

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

Out[42]:

```
0.059286006322965545
```

## Accuracy(Lasso)

In [43]:

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

Out[43]:

```
0.05399322398318607
```

## Accuracy(Elastic Net)

In [44]:

```
from sklearn.linear_model import ElasticNet
en=ElasticNet()
en.fit(x_train,y_train)
```

Out[44]:

```
ElasticNet()
```

In [45]:

```
en.coef_
```

Out[45]:

```
array([ 0.90807009,  0.          , -3.03778728,  1.50214494,  0.18520451,
         0.23397266, -0.0264939 , -2.25820265, -0.03191645,  0.00659687,
         2.32046868,  0.36749735,  1.03444336, -1.15987874])
```

In [46]:

```
en.intercept_
```

Out[46]:

```
28079038.682628848
```

In [47]:

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



In [48]:

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

Out[48]:

0.09587022840432957

## Evaluation Metrics

In [49]:

```
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)))
```

28.716055972336175

1131.179666418981

33.633014530650996

## Logistic Regression

In [50]:

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

In [51]:

```
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 [52]:

```
feature_matrix.shape
```

Out[52]:

(32381, 14)

In [53]:

```
target_vector.shape
```

Out[53]:

(32381,)

In [54]:

```
from sklearn.preprocessing import StandardScaler
```

In [55]:

```
fs=StandardScaler().fit_transform(feature_matrix)
```

In [56]:

```
logr=LogisticRegression(max_iter=10000)  
logr.fit(fs,target_vector)
```

Out[56]:

```
LogisticRegression(max_iter=10000)
```

In [57]:

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

In [58]:

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

```
[28079035]
```

In [59]:

```
logr.classes_
```

Out[59]:

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

In [60]:

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

Out[60]:

```
0.8480899292795158
```

In [61]:

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

Out[61]:

```
2.5638972732451705e-10
```

In [62]:

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

Out[62]:

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

## Random Forest

In [63]:

```
from sklearn.ensemble import RandomForestClassifier
```

In [64]:

```
rfc=RandomForestClassifier()  
rfc.fit(x_train,y_train)
```

Out[64]:

```
RandomForestClassifier()
```

In [65]:

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

In [66]:

```
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[66]:

```
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 [67]:

```
grid_search.best_score_
```

Out[67]:

```
0.7753022147710227
```

In [68]:

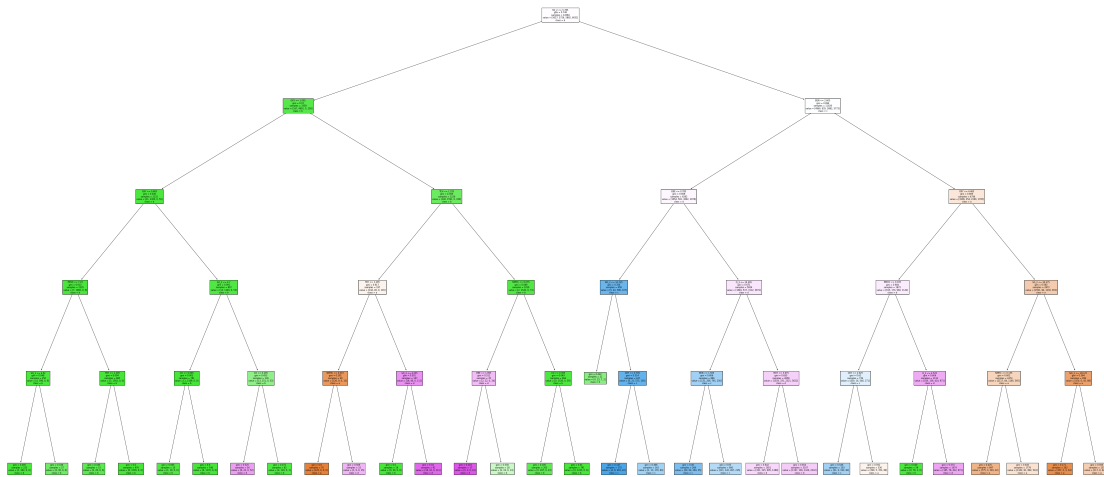
```
rfc_best=grid_search.best_estimator_
```

In [69]:

```

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
[357, 0, 3, 34]\nclclass = a'),
Text(4389.6, 181.199999999999982, 'gini = 0.595\nsamples = 137\nvalue =
[117, 0, 42, 54]\nclclass = a')]
```



## Conclusion

### Accuracy

**Linear Regression:0.19864042035800789**

**Ridge Regression:0.059286006322965545**

**Lasso Regression:0.05399322398318607**

**ElasticNet Regression:0.09587022840432957**

**Logistic Regression:0.8480899292795158**

**Random Forest:0.7753022147710227**

**Logistic Regression is suitable for this dataset**