

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_2009.csv")
df
```

Out[2]:

	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3	
0	2009-10-01 01:00:00	NaN	0.27	NaN	NaN	NaN	39.889999	48.150002	NaN	50.680000	18.0
1	2009-10-01 01:00:00	NaN	0.22	NaN	NaN	NaN	21.230000	24.260000	NaN	55.880001	10.0
2	2009-10-01 01:00:00	NaN	0.18	NaN	NaN	NaN	31.230000	34.880001	NaN	49.060001	25.0
3	2009-10-01 01:00:00	0.95	0.33	1.43	2.68	0.25	55.180000	81.360001	1.57	36.669998	26.0
4	2009-10-01 01:00:00	NaN	0.41	NaN	NaN	0.12	61.349998	76.260002	NaN	38.090000	23.0
...
215683	2009-06-01 00:00:00	0.50	0.22	0.39	0.75	0.09	22.000000	24.510000	1.00	82.239998	10.0
215684	2009-06-01 00:00:00	NaN	0.31	NaN	NaN	NaN	76.110001	101.099998	NaN	41.220001	9.0
215685	2009-06-01 00:00:00	0.13	NaN	0.86	NaN	0.23	81.050003	99.849998	NaN	24.830000	12.0
215686	2009-06-01 00:00:00	0.21	NaN	2.96	NaN	0.10	72.419998	82.959999	NaN	NaN	13.0
215687	2009-06-01 00:00:00	0.37	0.32	0.99	1.36	0.14	54.290001	64.480003	1.06	56.919998	15.0

215688 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: 24717 entries, 3 to 215687
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        24717 non-null  object
1   BEN         24717 non-null  float64
2   CO          24717 non-null  float64
3   EBE         24717 non-null  float64
4   MXY         24717 non-null  float64
5   NMHC        24717 non-null  float64
6   NO_2        24717 non-null  float64
7   NOx         24717 non-null  float64
8   OXY         24717 non-null  float64
9   O_3         24717 non-null  float64
10  PM10        24717 non-null  float64
11  PM25        24717 non-null  float64
12  PXY         24717 non-null  float64
13  SO_2        24717 non-null  float64
14  TCH         24717 non-null  float64
15  TOL         24717 non-null  float64
16  station     24717 non-null  int64
dtypes: float64(15), int64(1), object(1)
memory usage: 3.4+ MB
```

In [6]:

```
data=df[['PXY', 'NOx', 'OXY']]
data
```

Out[6]:

	PXY	NOx	OXY
3	1.30	81.360001	1.57
20	0.84	19.240000	1.00
24	1.07	43.919998	1.28
28	1.04	48.869999	1.21
45	0.88	19.299999	1.00
...
215659	0.84	29.490000	0.86
215663	1.03	69.870003	1.26
215667	0.92	82.629997	1.13
215683	0.74	24.510000	1.00
215687	0.83	64.480003	1.06

24717 rows × 3 columns

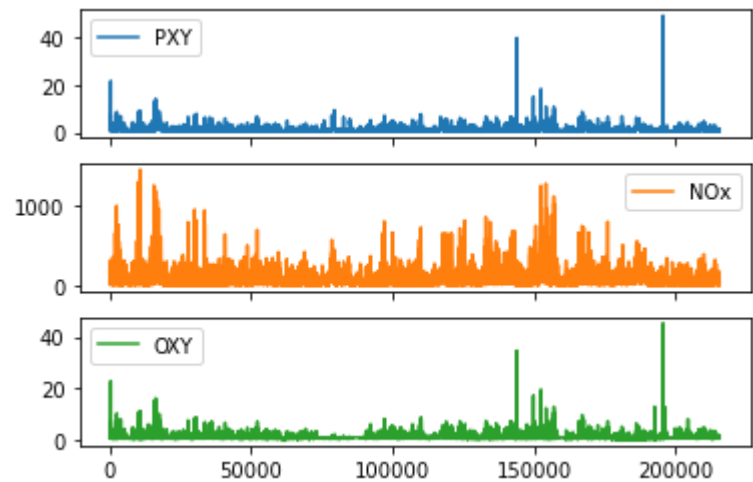
Line chart

In [7]:

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

Out[7]:

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



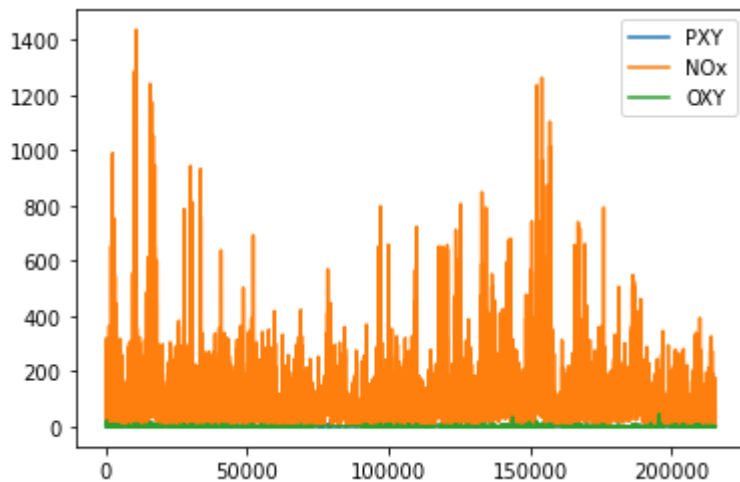
Line chart

In [8]:

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

Out[8]:

<AxesSubplot:>



Bar chart

In [9]:

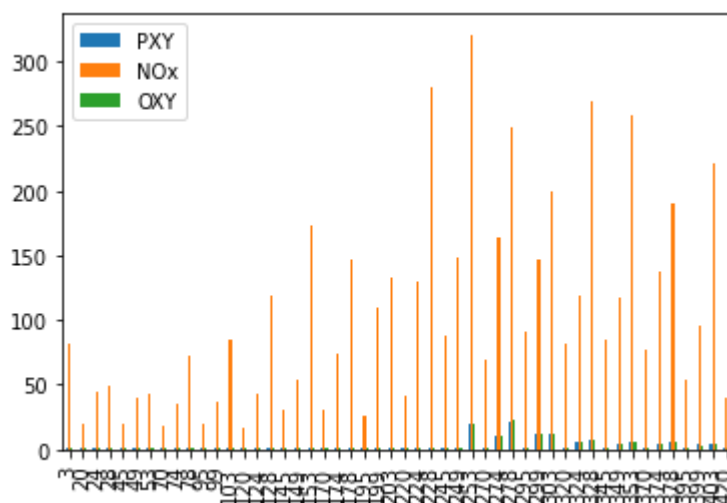
```
b=data[0:50]
```

In [10]:

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

Out[10]:

<AxesSubplot:>



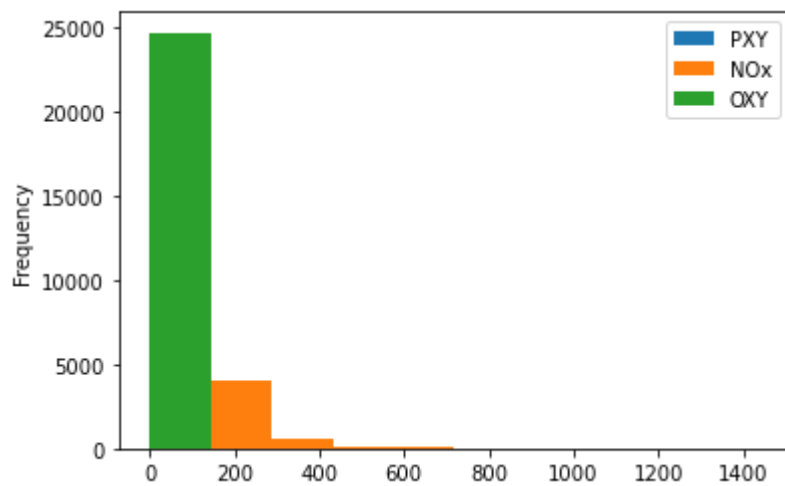
Histogram

In [11]:

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

Out[11]:

<AxesSubplot:ylabel='Frequency'>



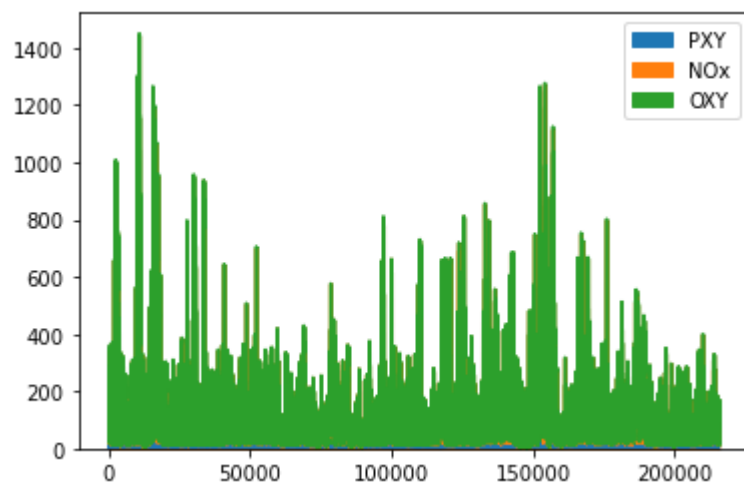
Area chart

In [12]:

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

Out[12]:

<AxesSubplot:>



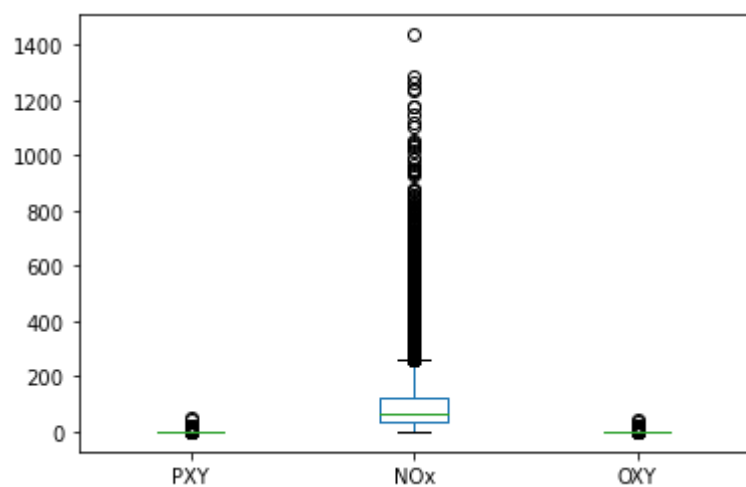
Box chart

In [13]:

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

Out[13]:

<AxesSubplot:>



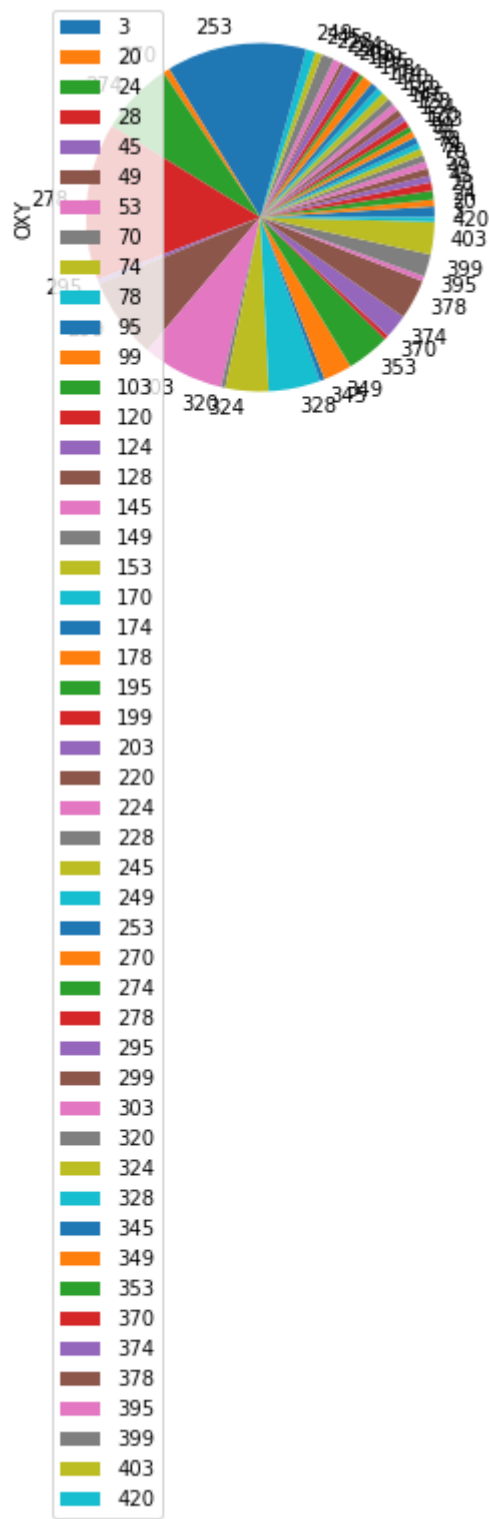
Pie chart

In [14]:

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

Out[14]:

<AxesSubplot:ylabel='OXY'>



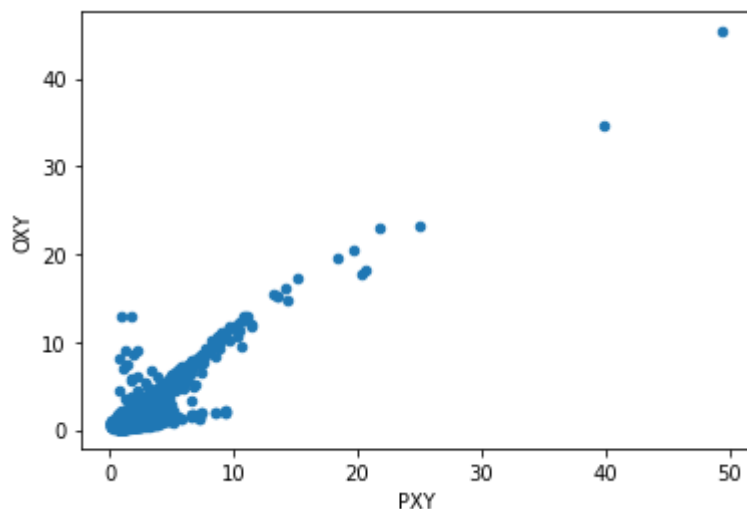
Scatter chart

In [15]:

```
data.plot.scatter(x='PXY', y='OXY')
```

Out[15]:

```
<AxesSubplot:xlabel='PXY', ylabel='OXY'>
```



In [16]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 24717 entries, 3 to 215687
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        24717 non-null  object
1   BEN         24717 non-null  float64
2   CO          24717 non-null  float64
3   EBE         24717 non-null  float64
4   MXY         24717 non-null  float64
5   NMHC        24717 non-null  float64
6   NO_2        24717 non-null  float64
7   NOx         24717 non-null  float64
8   OXY         24717 non-null  float64
9   O_3         24717 non-null  float64
10  PM10        24717 non-null  float64
11  PM25        24717 non-null  float64
12  PXY         24717 non-null  float64
13  SO_2        24717 non-null  float64
14  TCH         24717 non-null  float64
15  TOL         24717 non-null  float64
16  station     24717 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
count	24717.000000	24717.000000	24717.000000	24717.000000	24717.000000	24717.000000
mean	1.010583	0.448056	1.262430	2.244469	0.219582	55.563929
std	1.007345	0.291706	1.074768	2.242214	0.141661	38.911677
min	0.170000	0.060000	0.250000	0.240000	0.000000	0.600000
25%	0.460000	0.270000	0.720000	0.990000	0.140000	26.510000
50%	0.670000	0.370000	1.000000	1.490000	0.190000	47.930000
75%	1.180000	0.570000	1.430000	2.820000	0.260000	76.269997
max	22.379999	5.570000	47.669998	56.500000	2.580000	477.399994

In [18]:

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

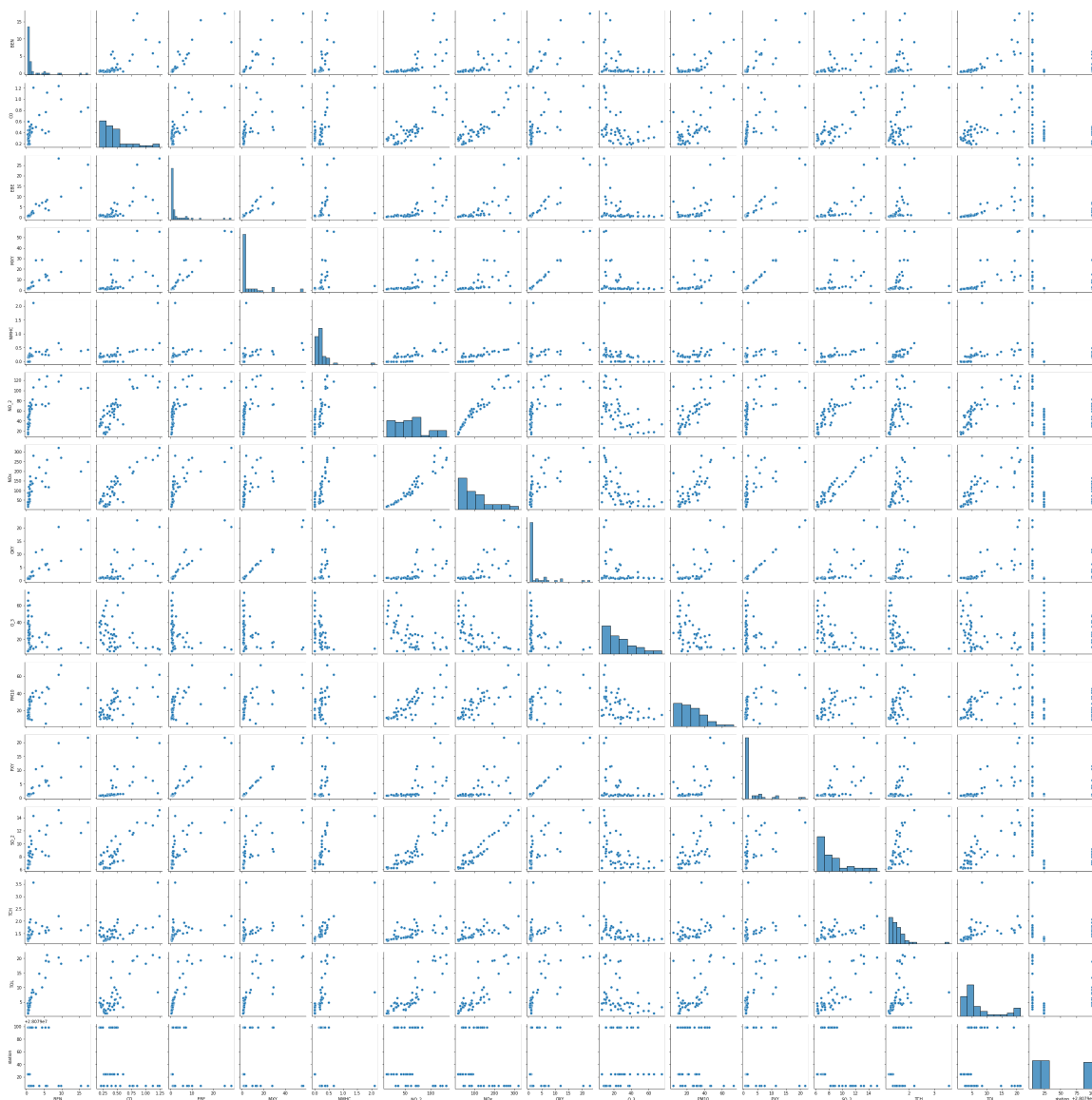
EDA AND VISUALIZATION

In [19]:

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

Out[19]:

<seaborn.axisgrid.PairGrid at 0x2183dee99d0>



In [20]:

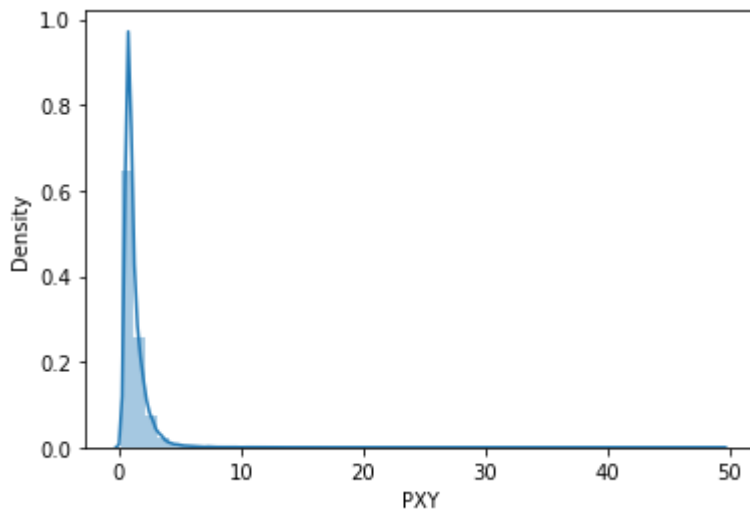
```
sns.distplot(df1['PXY'])
```

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 fun
ction for histograms).

```
warnings.warn(msg, FutureWarning)
```

Out[20]:

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

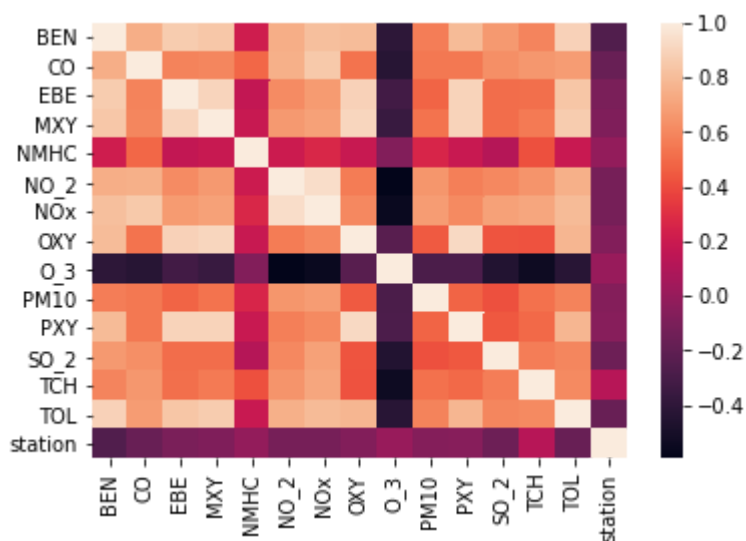


In [21]:

```
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]:

28078901.21705944

In [26]:

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

Out[26]:

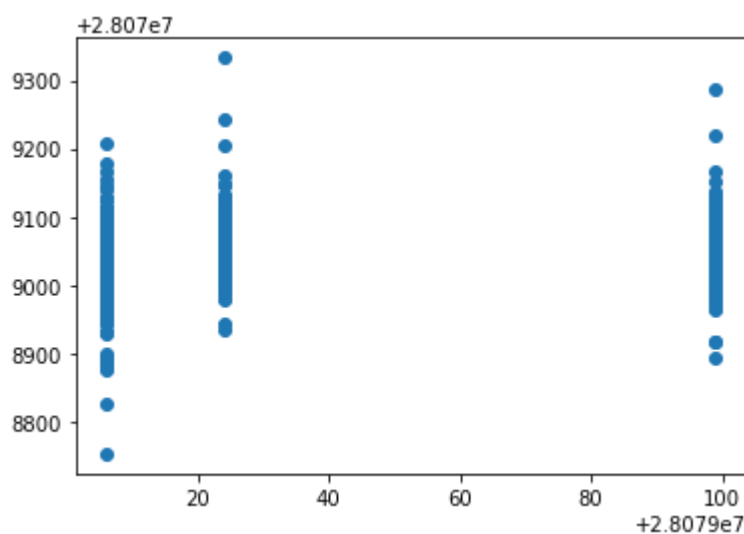
	Co-efficient
BEN	-37.569492
CO	-27.640183
EBE	7.444821
MXY	0.104949
NMHC	-18.836149
NO_2	-0.189433
NOx	0.210363
OXY	12.034927
O_3	0.022955
PM10	-0.056456
PXY	1.639321
SO_2	-0.369752
TCH	119.045865
TOL	-1.016695

In [27]:

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

Out[27]:

<matplotlib.collections.PathCollection at 0x2184fffb610>



ACCURACY

In [28]:

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

Out[28]:

0.2791163528241405

In [29]:

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

Out[29]:

0.2904829269822444

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]:

0.2789736363242975

In [33]:

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

Out[33]:

0.2901722163574638

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]:

```
0.03411461118355863
```

Accuracy(Lasso)

In [36]:

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

Out[36]:

```
0.03832089827644303
```

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([-7.09949154, -0.63183419,  0.29499383,  2.21487181,  0.          ,
        -0.23307397,  0.12859889,  1.00478444, -0.14991085,  0.07736015,
         1.84017177, -0.83222241,  1.50100091, -1.89102312])
```

In [39]:

```
en.intercept_
```

Out[39]:

```
28079064.81809502
```

In [40]:

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


In [41]:

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

Out[41]:

0.10463204276063354

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)))
```

35.84045092335339

1465.2505318850974

38.27859103840027

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]:

(24717, 14)

In [46]:

```
target_vector.shape
```

Out[46]:

(24717,)

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]:

```
0.8951733624630821
```

In [54]:

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

Out[54]:

```
5.447205522232353e-13
```

In [55]:

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

Out[55]:

```
array([[5.44720552e-13, 8.28692830e-44, 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]:

```
0.8947460115206034
```

In [61]:

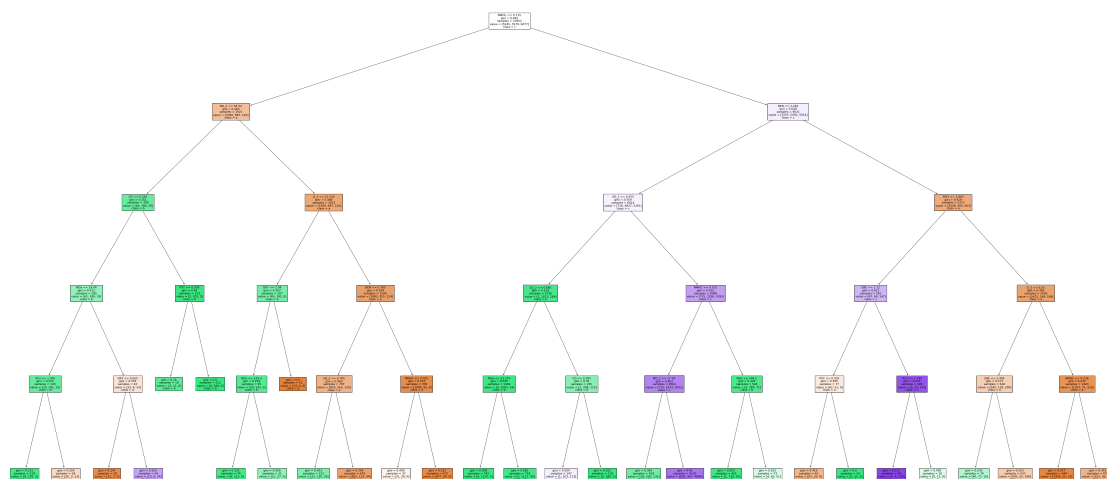
```
rfc_best=grid_search.best_estimator_
```

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
[1404, 25, 33]\n\nclass = a'),
Text(4381.333333333334, 181.19999999999982, 'gini = 0.453\nsamples = 47
6\nvalue = [523, 46, 183]\n\nclass = a')]
```



Conclusion

Accuracy

Linear Regression:0.2904829269822444

Ridge Regression:0.2901722163574638

Lasso Regression:0.03832089827644303

ElasticNet Regression:0.10463204276063354

Logistic Regression:0.8951733624630821

Random Forest:0.8947460115206034

Logistic Regression is suitable for this dataset