# Importing Libraries

```
In [109]: import numpy as np
          import pandas as pd
          import seaborn as sns
          import matplotlib.pyplot as plt
```

# Importing Datasets

```
In [110]: df=pd.read_csv(r"C:\Users\user\Desktop\csvs_per_year\csvs_per_year\madrid_2004.
          df
```

Out[110]:

| | date | BEN | CO | EBE | MXY | NMHC | NO_2 | NOx | OXY | O_3 | PI |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2004-08-01 01:00:00 | NaN | 0.66 | NaN | NaN | NaN | 89.550003 | 118.900002 | NaN | 40.020000 | 39.990 |
| 1 | 2004-08-01 01:00:00 | 2.66 | 0.54 | 2.99 | 6.08 | 0.18 | 51.799999 | 53.860001 | 3.28 | 51.689999 | 22.950 |
| 2 | 2004-08-01 01:00:00 | NaN | 1.02 | NaN | NaN | NaN | 93.389999 | 138.600006 | NaN | 20.860001 | 49.480 |
| 3 | 2004-08-01 01:00:00 | NaN | 0.53 | NaN | NaN | NaN | 87.290001 | 105.000000 | NaN | 36.730000 | 31.070 |
| 4 | 2004-08-01 01:00:00 | NaN | 0.17 | NaN | NaN | NaN | 34.910000 | 35.349998 | NaN | 86.269997 | 54.080 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 245491 | 2004-06-01 00:00:00 | 0.75 | 0.21 | 0.85 | 1.55 | 0.07 | 59.580002 | 64.389999 | 0.66 | 33.029999 | 30.900 |
| 245492 | 2004-06-01 00:00:00 | 2.49 | 0.75 | 2.44 | 4.57 | NaN | 97.139999 | 146.899994 | 2.34 | 7.740000 | 37.689 |
| 245493 | 2004-06-01 00:00:00 | NaN | NaN | NaN | NaN | 0.13 | 102.699997 | 132.600006 | NaN | 17.809999 | 22.840 |
| 245494 | 2004-06-01 00:00:00 | NaN | NaN | NaN | NaN | 0.09 | 82.599998 | 102.599998 | NaN | NaN | 45.630 |
| 245495 | 2004-06-01 00:00:00 | 3.01 | 0.67 | 2.78 | 5.12 | 0.20 | 92.550003 | 141.000000 | 2.60 | 11.460000 | 24.389 |

245496 rows × 17 columns

# Data Cleaning and Data Preprocessing

In [111]:
```python
df=df.dropna()
```

In [112]:
```python
df.columns
```

Out[112]: Index(['date', 'BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_
3',
       'PM10', 'PM25', 'PXY', 'SO_2', 'TCH', 'TOL', 'station'],
      dtype='object')

In [113]:
```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 19397 entries, 5 to 245495
Data columns (total 17 columns):
 #    Column   Non-Null Count   Dtype
---   ------   --------------   -----
 0    date     19397 non-null   object
 1    BEN      19397 non-null   float64
 2    CO       19397 non-null   float64
 3    EBE      19397 non-null   float64
 4    MXY      19397 non-null   float64
 5    NMHC     19397 non-null   float64
 6    NO_2     19397 non-null   float64
 7    NOx      19397 non-null   float64
 8    OXY      19397 non-null   float64
 9    O_3      19397 non-null   float64
 10   PM10     19397 non-null   float64
 11   PM25     19397 non-null   float64
 12   PXY      19397 non-null   float64
 13   SO_2     19397 non-null   float64
 14   TCH      19397 non-null   float64
 15   TOL      19397 non-null   float64
 16   station  19397 non-null   int64
dtypes: float64(15), int64(1), object(1)
memory usage: 2.7+ MB
```

In [114]: 
```python
data=df[['CO' ,'station']]
data
```
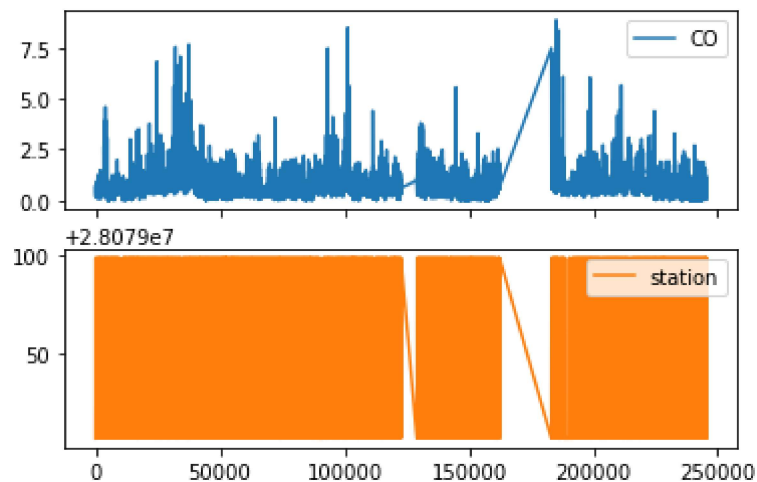
Out[114]:

|        | CO   | station  |
|--------|------|----------|
| 5      | 0.63 | 28079006 |
| 22     | 0.36 | 28079024 |
| 26     | 0.46 | 28079099 |
| 32     | 0.67 | 28079006 |
| 49     | 0.30 | 28079024 |
| ...    | ...  | ...      |
| 245463 | 0.08 | 28079024 |
| 245467 | 0.67 | 28079099 |
| 245473 | 1.12 | 28079006 |
| 245491 | 0.21 | 28079024 |
| 245495 | 0.67 | 28079099 |

19397 rows × 2 columns

# Line chart

In [115]: 
```python
data.plot.line(subplots=True)
```
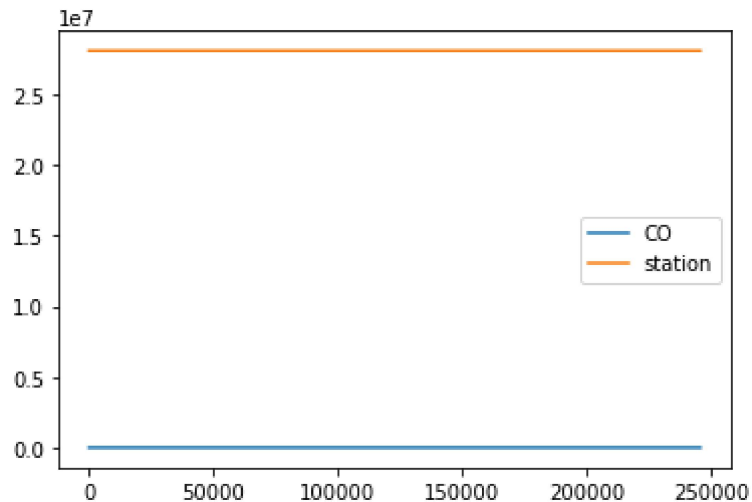
Out[115]: array([<AxesSubplot:>, <AxesSubplot:>], dtype=object)



# Line chart

In [116]:
```python
data.plot.line()
```
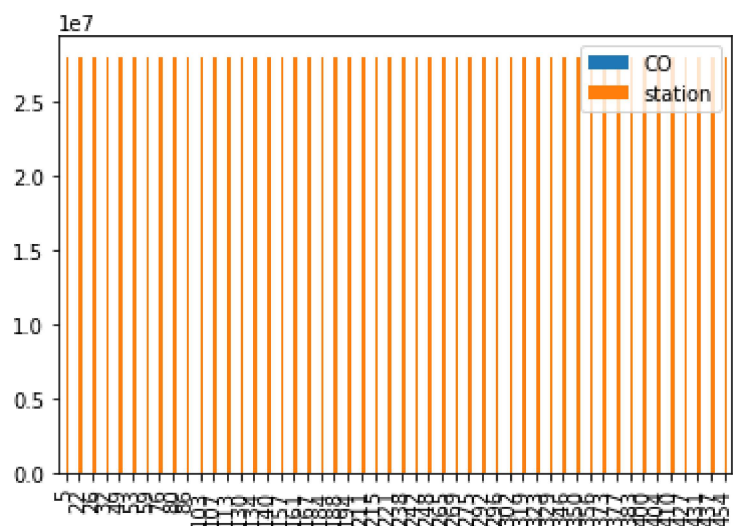
Out[116]:  <AxesSubplot:>



# Bar chart

In [117]:
```python
b=data[0:50]
```

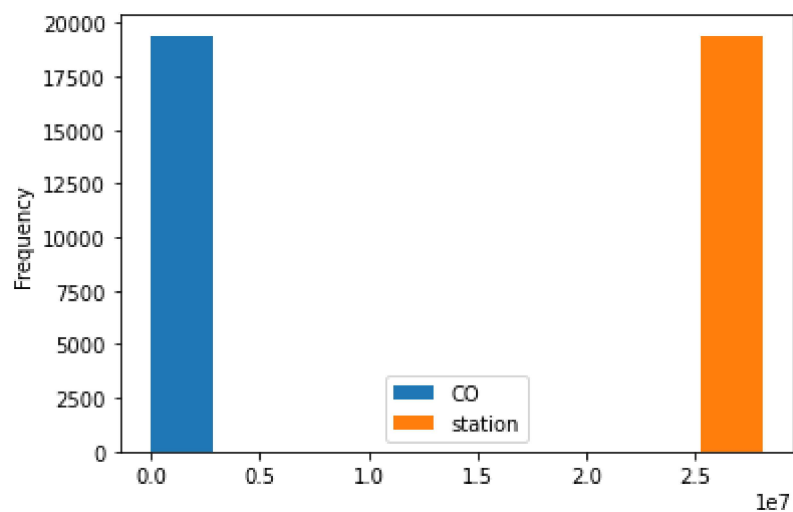In [118]:
```python
b.plot.bar()
```

Out[118]:  <AxesSubplot:>



# Histogram

In [119]: `data.plot.hist()`

Out[119]: `<AxesSubplot:ylabel='Frequency'>`



# Area chart

In [120]: `data.plot.area()`
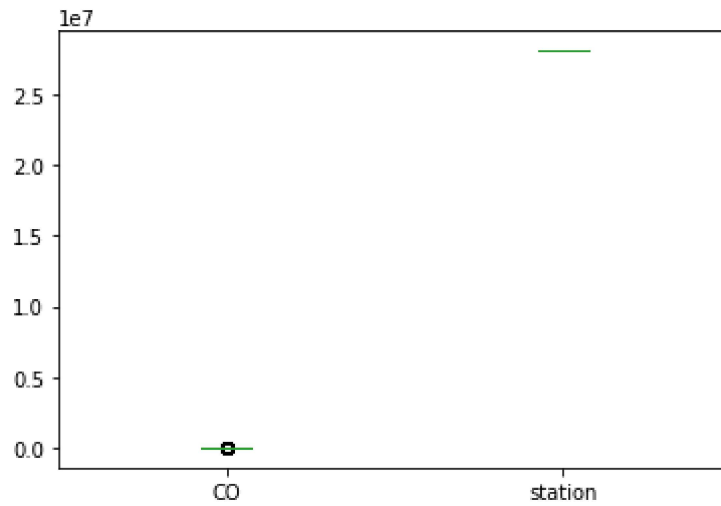
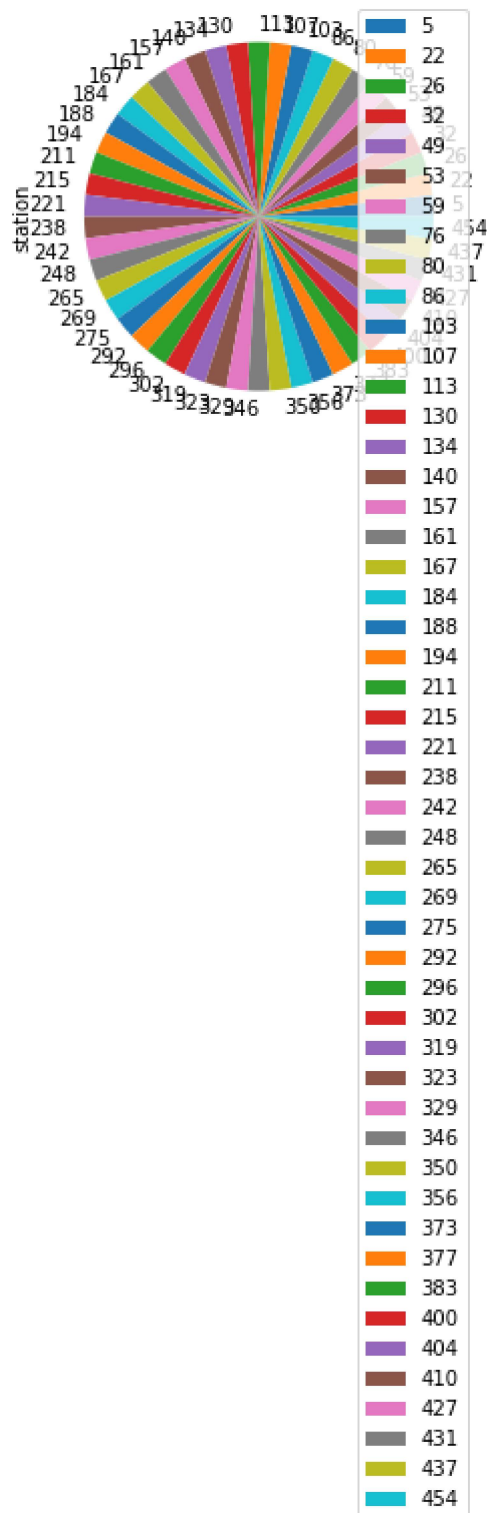Out[120]: `<AxesSubplot:>`



# Box chart

In [121]: `data.plot.box()`

Out[121]: `<AxesSubplot:>`



# Pie chart

In [122]:
```python
b.plot.pie(y='station' )
```
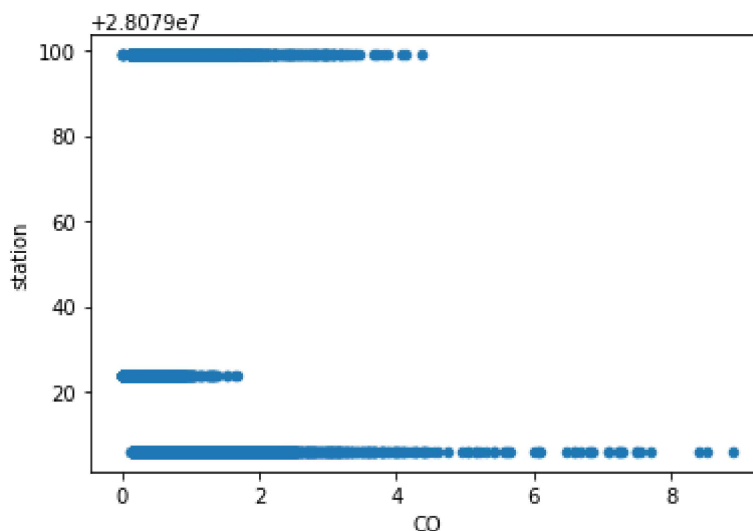
Out[122]:   <AxesSubplot:ylabel='station'>



## Scatter chart

In [123]: ```python
data.plot.scatter(x='CO' ,y='station')
```

Out[123]: `<AxesSubplot:xlabel='CO', ylabel='station'>`



In [124]: ```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 19397 entries, 5 to 245495
Data columns (total 17 columns):
 #   Column   Non-Null Count  Dtype
---  ------   --------------  -----
 0   date     19397 non-null  object
 1   BEN      19397 non-null  float64
 2   CO       19397 non-null  float64
 3   EBE      19397 non-null  float64
 4   MXY      19397 non-null  float64
 5   NMHC     19397 non-null  float64
 6   NO_2     19397 non-null  float64
 7   NOx      19397 non-null  float64
 8   OXY      19397 non-null  float64
 9   O_3      19397 non-null  float64
 10  PM10     19397 non-null  float64
 11  PM25     19397 non-null  float64
 12  PXY      19397 non-null  float64
 13  SO_2     19397 non-null  float64
 14  TCH      19397 non-null  float64
 15  TOL      19397 non-null  float64
 16  station  19397 non-null  int64
dtypes: float64(15), int64(1), object(1)
memory usage: 2.7+ MB
```

In [125]: `df.describe()`

Out[125]:

|         | BEN          | CO           | EBE          | MXY          | NMHC         | NO_2         | 193  |
|---------|--------------|--------------|--------------|--------------|--------------|--------------|------|
| count   | 19397.000000 | 19397.000000 | 19397.000000 | 19397.000000 | 19397.000000 | 19397.000000 | 193  |
| mean    | 2.250781     | 0.675347     | 2.775913     | 5.424809     | 0.151024     | 62.887023    | 1    |
| std     | 2.184724     | 0.591026     | 2.729622     | 5.554358     | 0.158603     | 37.952255    | 1    |
| min     | 0.000000     | 0.000000     | 0.000000     | 0.000000     | 0.000000     | 0.090000     |      |
| 25%     | 0.870000     | 0.320000     | 1.020000     | 1.780000     | 0.060000     | 35.150002    |      |
| 50%     | 1.620000     | 0.520000     | 1.970000     | 3.800000     | 0.110000     | 58.310001    |      |
| 75%     | 2.910000     | 0.860000     | 3.580000     | 7.260000     | 0.200000     | 85.730003    | 1    |
| max     | 34.180000    | 8.900000     | 41.880001    | 91.599998    | 4.810000     | 355.100006   | 17   |

In [126]: 
```python
df1=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
 'PM10', 'PXY', 'SO_2', 'TCH', 'TOL', 'station']]
```
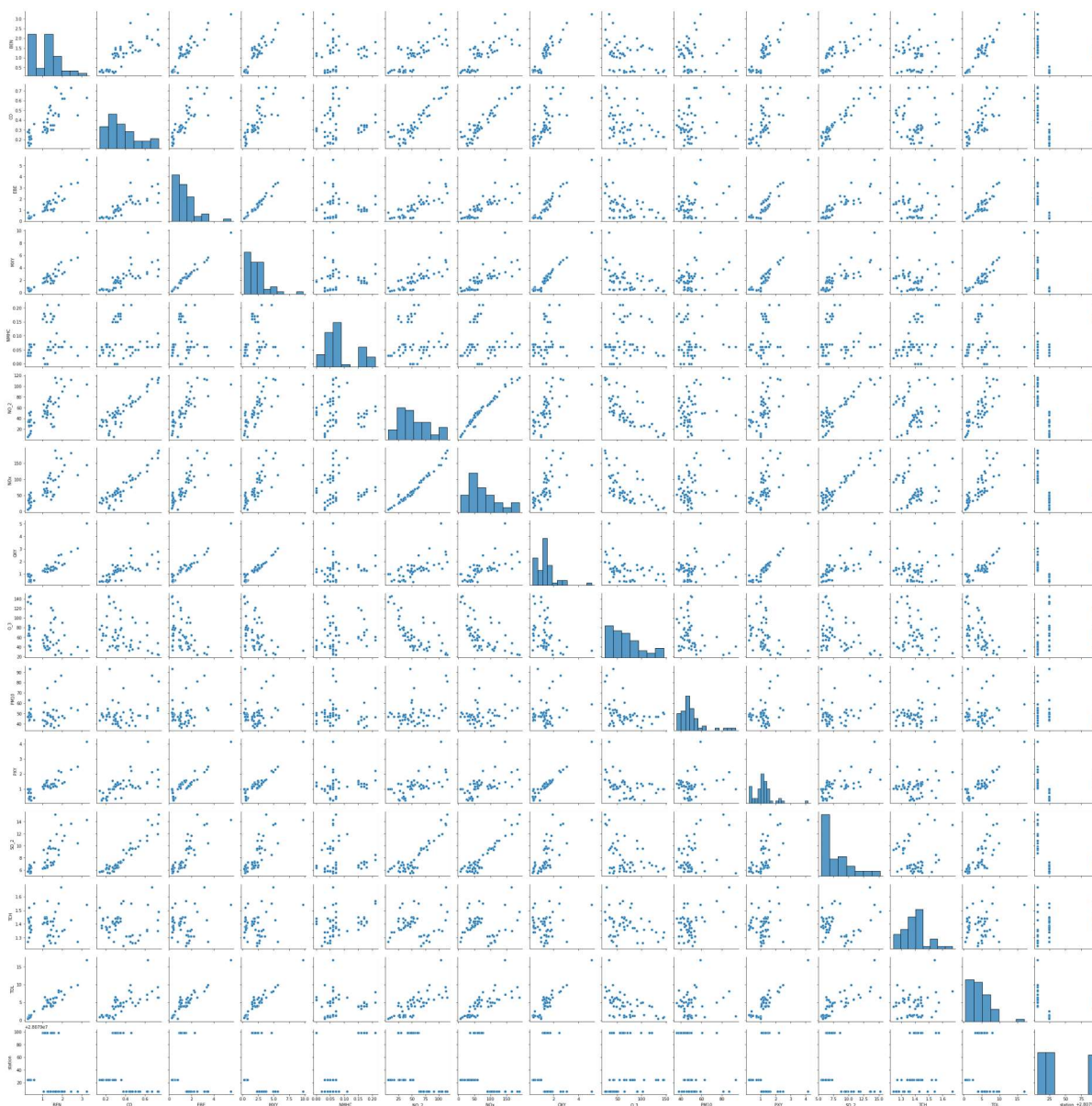
# EDA AND VISUALIZATION

In [138]: `sns.pairplot(df1[0:50])`
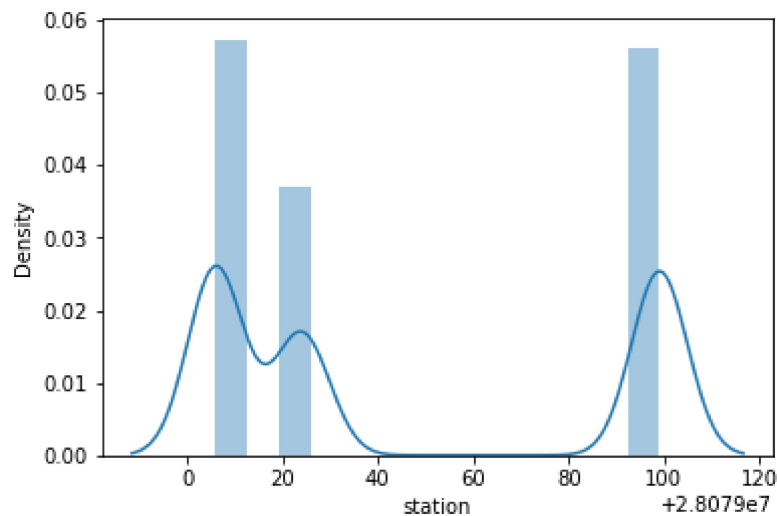
Out[138]: `<seaborn.axisgrid.PairGrid at 0x1cd2081d5e0>`

In [139]: `sns.distplot(df1['station'])`

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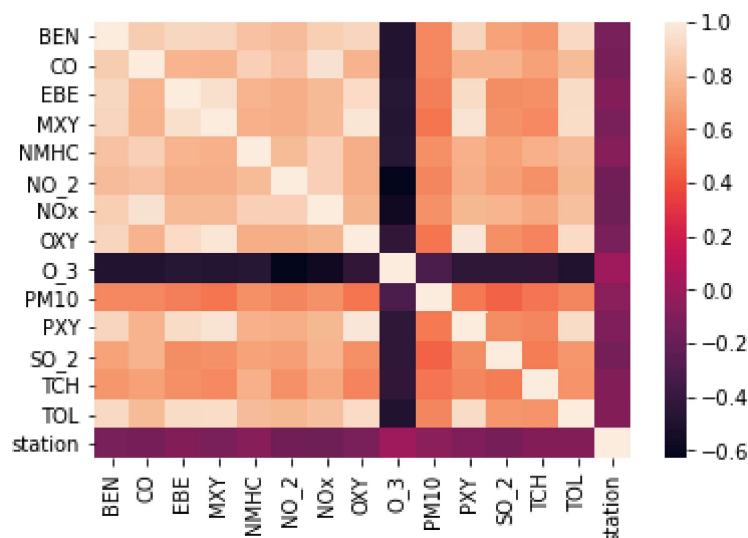stograms).
    warnings.warn(msg, FutureWarning)

Out[139]: `<AxesSubplot:xlabel='station', ylabel='Density'>`



In [140]: `sns.heatmap(df1.corr())`

Out[140]: `<AxesSubplot:>`



# TO TRAIN THE MODEL AND MODEL BULDING

In [141]:
```python
x=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
 'PM10', 'PXY', 'SO_2', 'TCH', 'TOL']]
y=df['station']
```

In [142]:
```python
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 [143]:
```python
from sklearn.linear_model import LinearRegression
lr=LinearRegression()
lr.fit(x_train,y_train)
```

Out[143]: LinearRegression()

In [144]:
```python
lr.intercept_
```

Out[144]: 28079074.904809926

In [145]:
```python
coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

Out[145]:

|      | Co-efficient |
| ---- | ------------ |
| BEN  | -4.149320    |
| CO   | 28.974299    |
| EBE  | 4.093538     |
| MXY  | -3.564947    |
| NMHC | 76.214166    |
| NO_2 | -0.152042    |
| NOx  | -0.263485    |
| OXY  | -2.700986    |
| O_3  | -0.299274    |
| PM10 | 0.091041     |
| PXY  | 6.509818     |
| SO_2 | -0.207503    |
| TCH  | -6.417410    |
| TOL  | 1.237076     |

```
In [146]: prediction =lr.predict(x_test)
          plt.scatter(y_test,prediction)
```

Out[146]: <matplotlib.collections.PathCollection at 0x1cd000e0c40>



# ACCURACY

```
In [147]: lr.score(x_test,y_test)
```

Out[147]: 0.10368510261153874

```
In [148]: lr.score(x_train,y_train)
```

Out[148]: 0.10720468459527421

# Ridge and Lasso

```
In [149]: from sklearn.linear_model import Ridge,Lasso
```

```
In [150]: rr=Ridge(alpha=10)
          rr.fit(x_train,y_train)
```

Out[150]: Ridge(alpha=10)

# Accuracy(Ridge)

In [151]: `rr.score(x_test,y_test)`

Out[151]: 0.1013285735884294

In [152]: `rr.score(x_train,y_train)`

Out[152]: 0.10691065764218544

In [153]:
```python
la=Lasso(alpha=10)
la.fit(x_train,y_train)
```

Out[153]: Lasso(alpha=10)

In [154]: `la.score(x_train,y_train)`

Out[154]: 0.0562981891575286

# Accuracy(Lasso)

In [155]: `la.score(x_test,y_test)`

Out[155]: 0.04747300651618447

In [156]:
```python
from sklearn.linear_model import ElasticNet
en=ElasticNet()
en.fit(x_train,y_train)
```

Out[156]: ElasticNet()

In [157]: `en.coef_`

Out[157]: array([-0.        ,  0.42069669,  1.41170788, -1.88395165,  0.        ,
              -0.16846589, -0.09328818, -0.        , -0.23225802,  0.11512537,
               0.38645125, -0.13953972,  0.        ,  1.21346959])

In [158]: `en.intercept_`

Out[158]: 28079067.29497689

In [159]: `prediction=en.predict(x_test)`

In [160]: `en.score(x_test,y_test)`

Out[160]: 0.06055311810140007

# Evaluation Metrics

In [161]:
```python
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)))
```

```
38.63234705078924
1666.5565468600973
40.823480337424655
```

# Logistic Regression

In [162]:
```python
from sklearn.linear_model import LogisticRegression
```

In [163]:
```python
feature_matrix=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_
 'PM10', 'PXY', 'SO_2', 'TCH', 'TOL']]
target_vector=df[ 'station']
```

In [164]:
```python
feature_matrix.shape
```

Out[164]: (19397, 14)

In [165]:
```python
target_vector.shape
```

Out[165]: (19397,)

In [166]:
```python
from sklearn.preprocessing import StandardScaler
```

In [167]:
```python
fs=StandardScaler().fit_transform(feature_matrix)
```

In [168]:
```python
logr=LogisticRegression(max_iter=10000)
logr.fit(fs,target_vector)
```

Out[168]: LogisticRegression(max_iter=10000)

In [169]:
```python
observation=[[1,2,3,4,5,6,7,8,9,10,11,12,13,14]]
```

In [170]:
```python
prediction=logr.predict(observation)
print(prediction)
```

```
[28079006]
```

In [171]:
```python
logr.classes_
```

Out[171]: array([28079006, 28079024, 28079099], dtype=int64)

In [172]:
```python
logr.score(fs,target_vector)
```

Out[172]: 0.7360416559261741

In [173]:
```python
logr.predict_proba(observation)[0][0]
```

Out[173]: 0.9999978255573396

In [174]:
```python
logr.predict_proba(observation)
```

Out[174]: array([[9.99997826e-01, 7.75018107e-20, 2.17444266e-06]])

# Random Forest

In [175]:
```python
from sklearn.ensemble import RandomForestClassifier
```

In [176]:
```python
rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

Out[176]: RandomForestClassifier()

In [177]:
```python
parameters={'max_depth':[1,2,3,4,5],
 'min_samples_leaf':[5,10,15,20,25],
 'n_estimators':[10,20,30,40,50]
}
```

```python
In [178]: from sklearn.model_selection import GridSearchCV
          grid_search =GridSearchCV(estimator=rfc,param_grid=parameters,cv=2,scoring="acc
          grid_search.fit(x_train,y_train)
```

```
Out[178]: 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')
```

```python
In [179]: grid_search.best_score_
```

```
Out[179]: 0.7760189647834945
```

```python
In [180]: rfc_best=grid_search.best_estimator_
```

```python
In [181]: 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',
```

```
 Text(2539.8620689655177, 181.19999999999982, 'gini = 0.349\nsamples = 1250
\nvalue = [207, 209, 1601]\nclass = c'),
 Text(2924.689655172414, 906.0, 'EBE <= 4.325\ngini = 0.545\nsamples = 2431
\nvalue = [2003, 215, 1552]\nclass = a'),
 Text(2770.7586206896553, 543.5999999999999, 'NO_2 <= 89.8\ngini = 0.541\ns
amples = 2293\nvalue = [1953, 207, 1389]\nclass = a'),
 Text(2693.7931034482763, 181.19999999999982, 'gini = 0.551\nsamples = 1700
\nvalue = [1298, 149, 1184]\nclass = a'),
 Text(2847.724137931035, 181.19999999999982, 'gini = 0.437\nsamples = 593\n
value = [655, 58, 205]\nclass = a'),
 Text(3078.6206896551726, 543.5999999999999, 'O_3 <= 6.98\ngini = 0.404\nsa
mples = 138\nvalue = [50, 8, 163]\nclass = c'),
 Text(3001.6551724137935, 181.19999999999982, 'gini = 0.127\nsamples = 55\n
value = [6, 0, 82]\nclass = c'),
 Text(3155.586206896552, 181.19999999999982, 'gini = 0.516\nsamples = 83\nv
alue = [44, 8, 81]\nclass = c'),
 Text(3848.275862068966, 1268.4, 'NOx <= 143.5\ngini = 0.39\nsamples = 1752
\nvalue = [2091, 25, 714]\nclass = a'),
 Text(3540.4137931034484, 906.0, 'TCH <= 1.365\ngini = 0.497\nsamples = 170
\nvalue = [113, 5, 171]\nclass = c'),
```

# Conclusion

# Accuracy

Linear Regression:0.11045993310581825

Ridge Regression:0.11011845641033136

Lasso Regression:0.044541624988104433

ElasticNet Regression:0.055622313951294466

Logistic Regression:0.7360416559261741

Random Forest:0.7754293098484298

# From the above data, we can conclude that logistic regression and random forest is preferrable to other regression types

In [ ]: