

# 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("madrid_2005.csv")
df
```

Out[2]:

|        | date                | BEN  | CO   | EBE  | MXY  | NMHC | NO_2      | NOx        | OXY  | O_3       | PM10  |
|--------|---------------------|------|------|------|------|------|-----------|------------|------|-----------|-------|
| 0      | 2005-11-01 01:00:00 | NaN  | 0.77 | NaN  | NaN  | NaN  | 57.130001 | 128.699997 | NaN  | 14.720000 | 14.91 |
| 1      | 2005-11-01 01:00:00 | 1.52 | 0.65 | 1.49 | 4.57 | 0.25 | 86.559998 | 181.699997 | 1.27 | 11.680000 | 30.93 |
| 2      | 2005-11-01 01:00:00 | NaN  | 0.40 | NaN  | NaN  | NaN  | 46.119999 | 53.000000  | NaN  | 30.469999 | 14.60 |
| 3      | 2005-11-01 01:00:00 | NaN  | 0.42 | NaN  | NaN  | NaN  | 37.220001 | 52.009998  | NaN  | 21.379999 | 15.16 |
| 4      | 2005-11-01 01:00:00 | NaN  | 0.57 | NaN  | NaN  | NaN  | 32.160000 | 36.680000  | NaN  | 33.410000 | 5.00  |
| ...    | ...                 | ...  | ...  | ...  | ...  | ...  | ...       | ...        | ...  | ...       | ...   |
| 236995 | 2006-01-01 00:00:00 | 1.08 | 0.36 | 1.01 | NaN  | 0.11 | 21.990000 | 23.610001  | NaN  | 43.349998 | 5.00  |
| 236996 | 2006-01-01 00:00:00 | 0.39 | 0.54 | 1.00 | 1.00 | 0.11 | 2.200000  | 4.220000   | 1.00 | 69.639999 | 4.95  |
| 236997 | 2006-01-01 00:00:00 | 0.19 | NaN  | 0.26 | NaN  | 0.08 | 26.730000 | 30.809999  | NaN  | 43.840000 | 4.31  |
| 236998 | 2006-01-01 00:00:00 | 0.14 | NaN  | 1.00 | NaN  | 0.06 | 13.770000 | 17.770000  | NaN  | NaN       | 5.00  |
| 236999 | 2006-01-01 00:00:00 | 0.50 | 0.40 | 0.73 | 1.84 | 0.13 | 20.940001 | 26.950001  | 1.49 | 48.259998 | 5.67  |

237000 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: 20070 entries, 5 to 236999
Data columns (total 17 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   date      20070 non-null   object 
 1   BEN        20070 non-null   float64
 2   CO         20070 non-null   float64
 3   EBE        20070 non-null   float64
 4   MXY        20070 non-null   float64
 5   NMHC       20070 non-null   float64
 6   NO_2       20070 non-null   float64
 7   NOx        20070 non-null   float64
 8   OXY        20070 non-null   float64
 9   O_3         20070 non-null   float64
 10  PM10       20070 non-null   float64
 11  PM25       20070 non-null   float64
 12  PXY        20070 non-null   float64
 13  SO_2       20070 non-null   float64
 14  TCH         20070 non-null   float64
 15  TOL         20070 non-null   float64
 16  station    20070 non-null   int64  
dtypes: float64(15), int64(1), object(1)
memory usage: 2.8+ MB
```

```
In [6]: data=df[['CO' , 'station']]  
data
```

Out[6]:

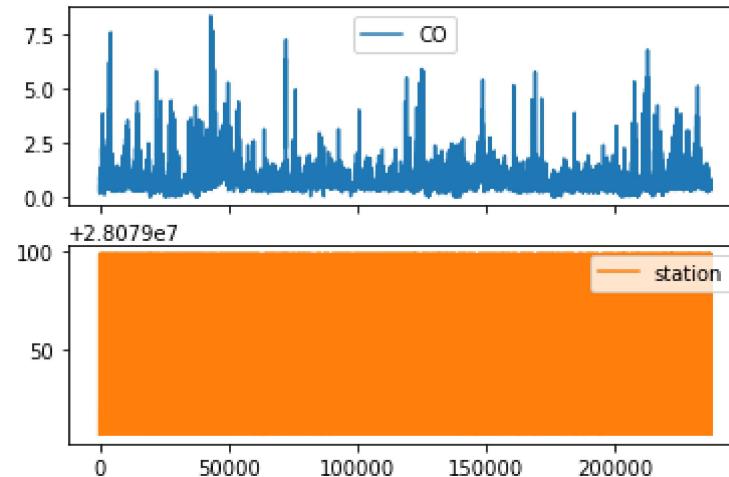
|        | CO   | station  |
|--------|------|----------|
| 5      | 0.88 | 28079006 |
| 22     | 0.22 | 28079024 |
| 25     | 0.49 | 28079099 |
| 31     | 0.84 | 28079006 |
| 48     | 0.20 | 28079024 |
| ...    | ...  | ...      |
| 236970 | 0.39 | 28079024 |
| 236973 | 0.45 | 28079099 |
| 236979 | 0.38 | 28079006 |
| 236996 | 0.54 | 28079024 |
| 236999 | 0.40 | 28079099 |

20070 rows × 2 columns

## Line chart

```
In [7]: data.plot.line(subplots=True)
```

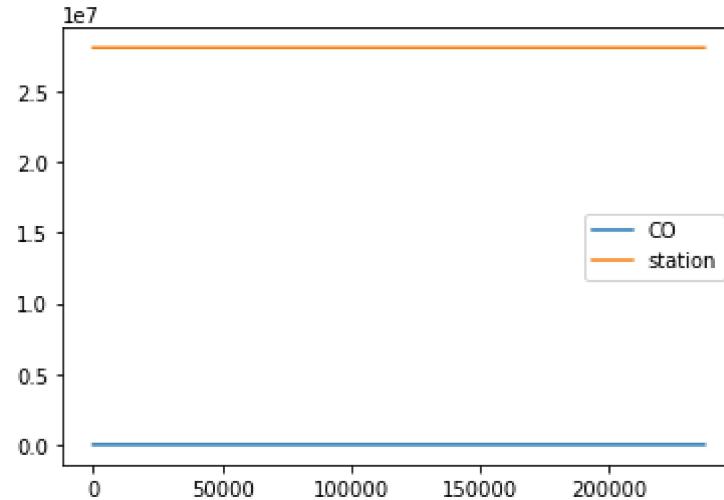
Out[7]: array([<AxesSubplot:>, <AxesSubplot:>], dtype=object)



## Line chart

```
In [8]: data.plot.line()
```

```
Out[8]: <AxesSubplot:>
```

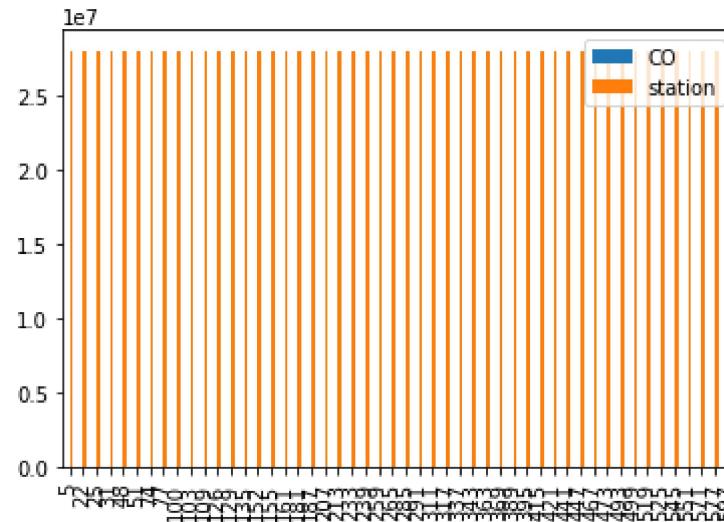


## Bar chart

```
In [9]: b=data[0:50]
```

```
In [10]: b.plot.bar()
```

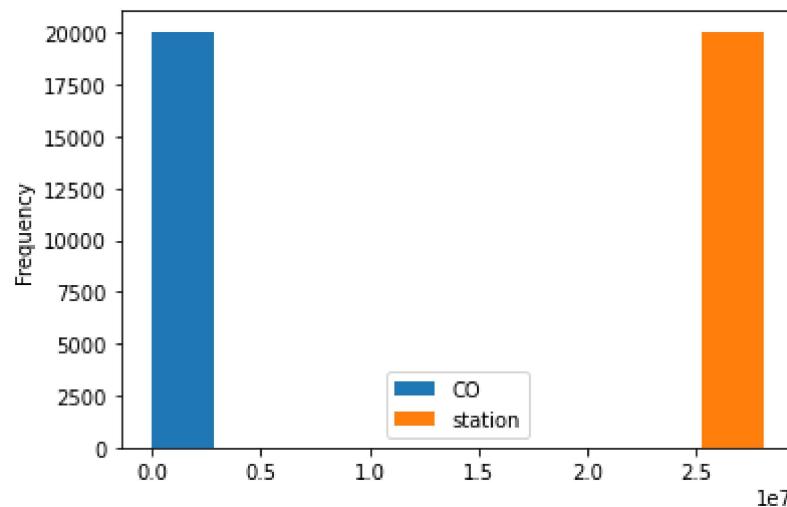
```
Out[10]: <AxesSubplot:>
```



## Histogram

```
In [11]: data.plot.hist()
```

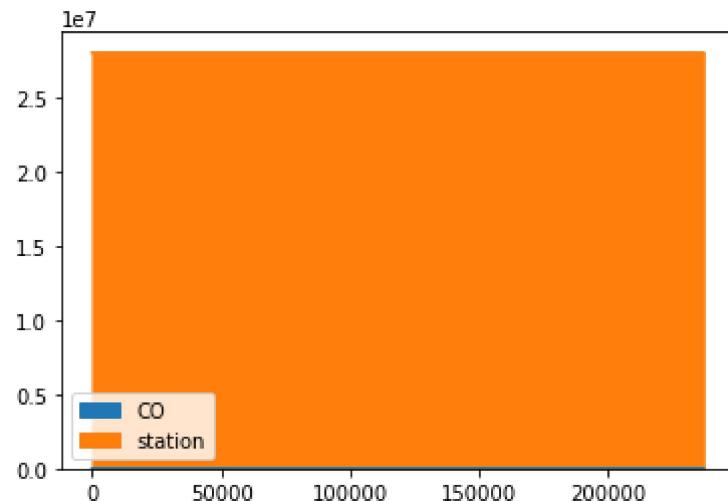
```
Out[11]: <AxesSubplot:ylabel='Frequency'>
```



## Area chart

```
In [12]: data.plot.area()
```

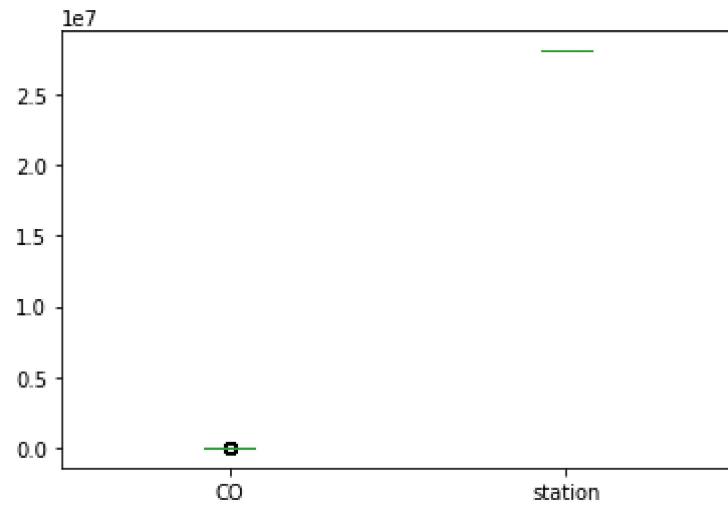
```
Out[12]: <AxesSubplot:>
```



## Box chart

```
In [13]: data.plot.box()
```

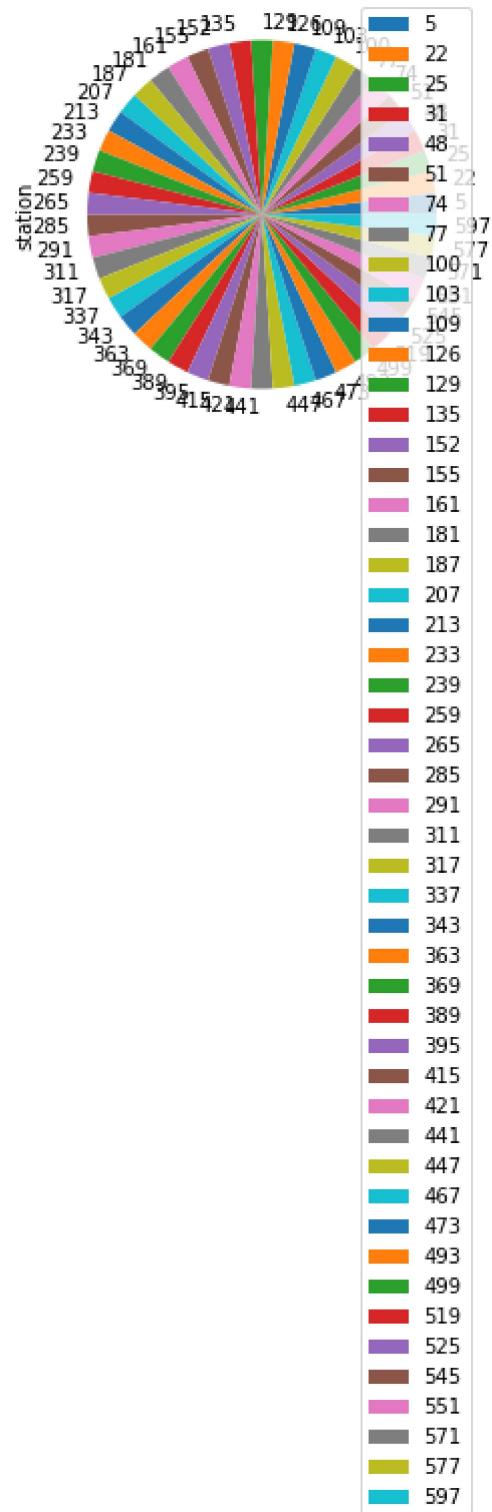
```
Out[13]: <AxesSubplot:>
```



## Pie chart

```
In [14]: b.plot.pie(y='station' )
```

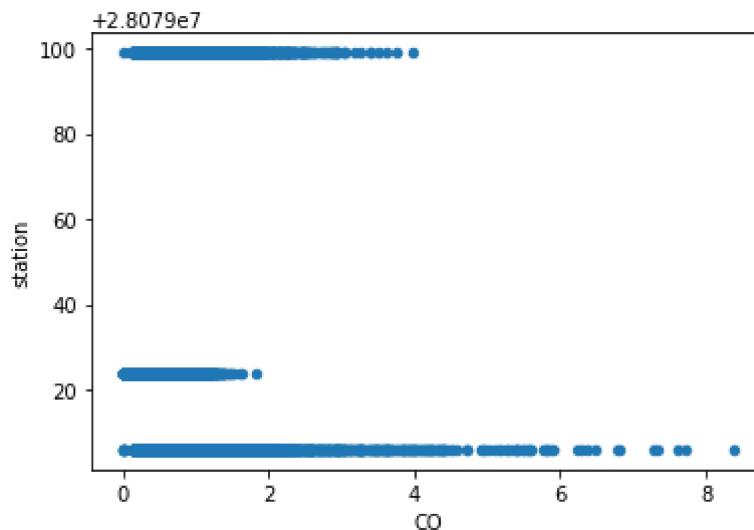
```
Out[14]: <AxesSubplot:ylabel='station'>
```



## Scatter chart

```
In [15]: data.plot.scatter(x='CO' ,y='station')
```

```
Out[15]: <AxesSubplot:xlabel='CO', ylabel='station'>
```



```
In [16]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 20070 entries, 5 to 236999
Data columns (total 17 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   date      20070 non-null   object 
 1   BEN       20070 non-null   float64
 2   CO        20070 non-null   float64
 3   EBE       20070 non-null   float64
 4   MXY       20070 non-null   float64
 5   NMHC      20070 non-null   float64
 6   NO_2      20070 non-null   float64
 7   NOx       20070 non-null   float64
 8   OXY       20070 non-null   float64
 9   O_3        20070 non-null   float64
 10  PM10      20070 non-null   float64
 11  PM25      20070 non-null   float64
 12  PXY       20070 non-null   float64
 13  SO_2      20070 non-null   float64
 14  TCU       20070 non-null   float64
```

```
In [17]: df.describe()
```

Out[17]:

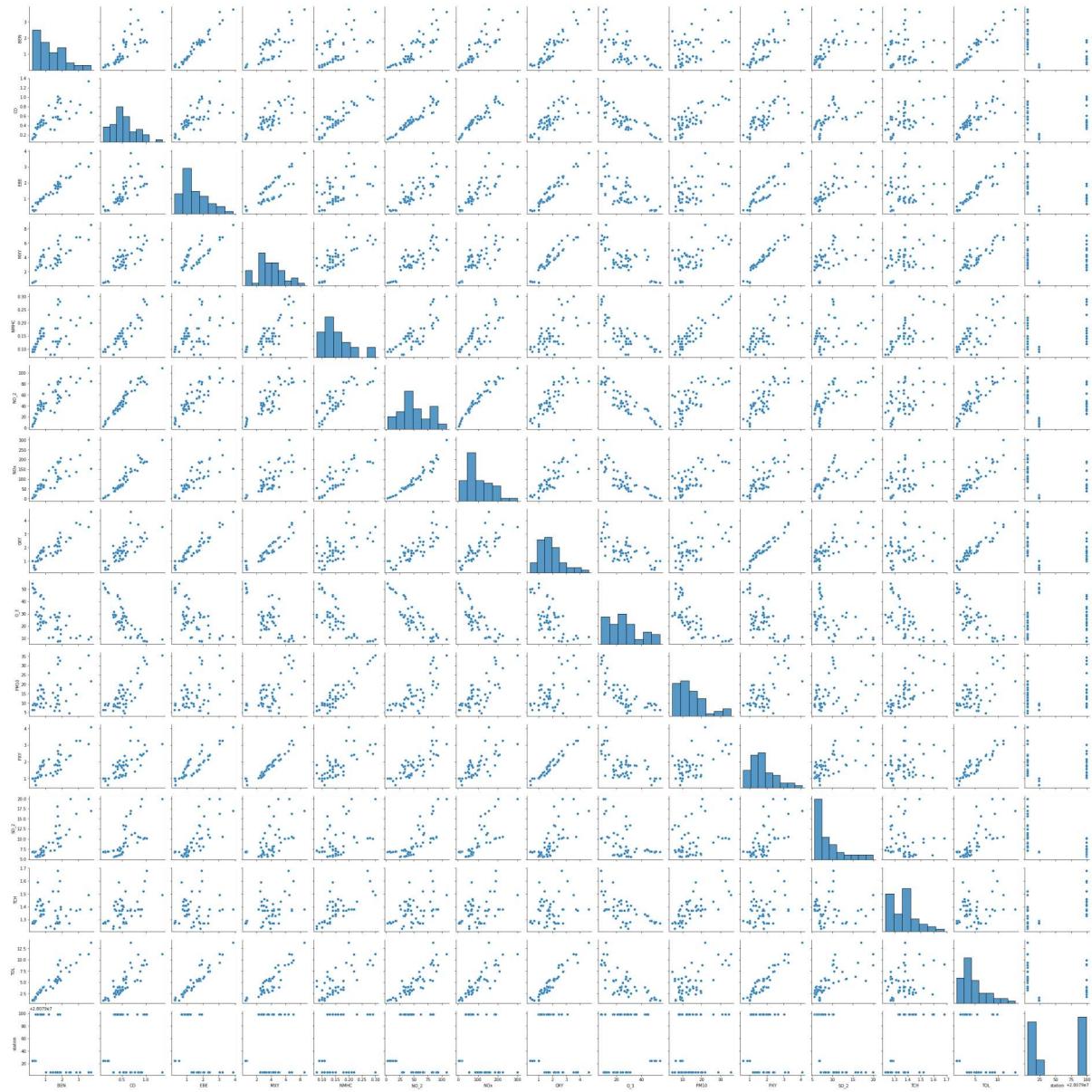
|       | BEN          | CO           | EBE          | MXY          | NMHC         | NO_2         | NOx          |
|-------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| count | 20070.000000 | 20070.000000 | 20070.000000 | 20070.000000 | 20070.000000 | 20070.000000 | 20070.000000 |
| mean  | 1.923656     | 0.720657     | 2.345423     | 5.457855     | 0.179282     | 66.226924    | 14.100000    |
| std   | 2.019061     | 0.549723     | 2.379219     | 5.495147     | 0.152783     | 40.568197    | 11.100000    |
| min   | 0.000000     | 0.000000     | 0.000000     | 0.000000     | 0.000000     | 0.000000     | 0.000000     |
| 25%   | 0.690000     | 0.400000     | 0.950000     | 1.930000     | 0.090000     | 36.602499    | 1.100000     |
| 50%   | 1.260000     | 0.580000     | 1.480000     | 3.800000     | 0.150000     | 60.525000    | 11.100000    |
| 75%   | 2.510000     | 0.880000     | 2.950000     | 7.210000     | 0.220000     | 89.317499    | 19.100000    |
| max   | 26.570000    | 8.380000     | 29.870001    | 71.050003    | 1.880000     | 419.500000   | 17.100000    |

```
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])
```

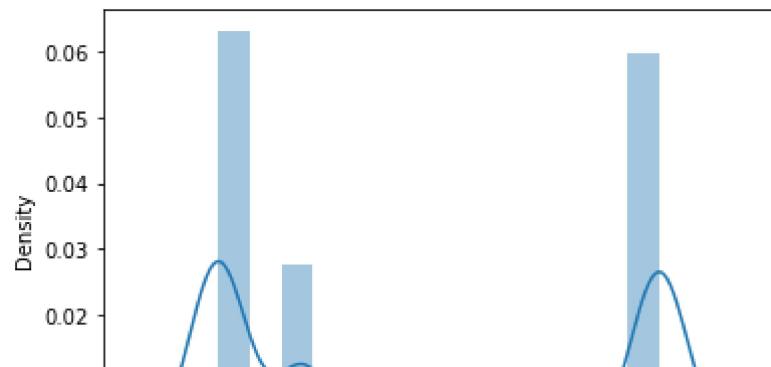
```
Out[19]: <seaborn.axisgrid.PairGrid at 0x175d2dc0670>
```



In [20]: `sns.distplot(df1['station'])`

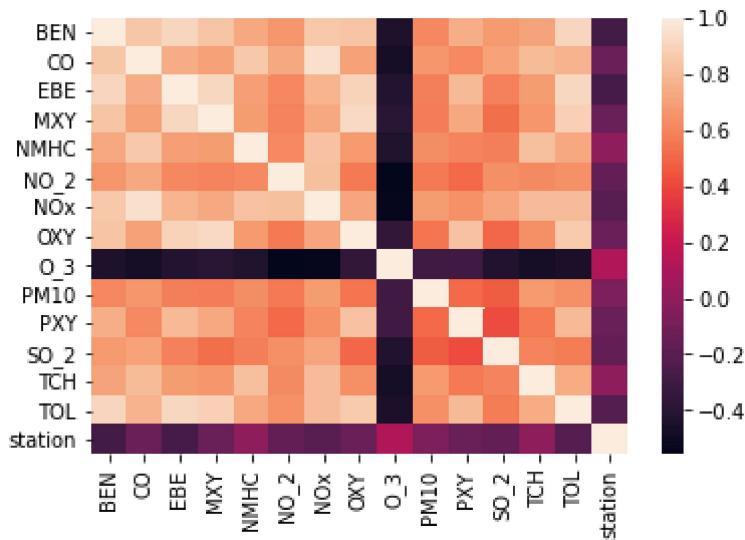
```
C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function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

Out[20]: <AxesSubplot:xlabel='station', 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]: 28078953.671107177
```

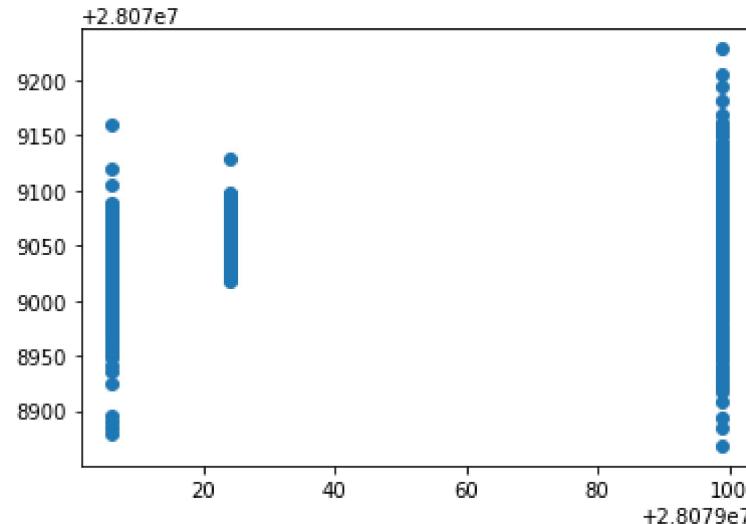
```
In [26]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])  
coeff
```

```
Out[26]:
```

|      | Co-efficient |
|------|--------------|
| BEN  | -9.277323    |
| CO   | 36.063969    |
| EBE  | -14.111555   |
| MXY  | 4.072225     |
| NMHC | 79.436208    |
| NO_2 | 0.117829     |
| NOx  | -0.261722    |
| OXY  | 2.848690     |
| O_3  | 0.036495     |
| PM10 | 0.045172     |
| PXY  | 2.894766     |
| SO_2 | 0.259341     |
| TCH  | 66.132314    |
| TOL  | -0.602787    |

```
In [27]: prediction = lr.predict(x_test)
plt.scatter(y_test,prediction)
```

```
Out[27]: <matplotlib.collections.PathCollection at 0x175e172e1c0>
```



## ACCURACY

```
In [28]: lr.score(x_test,y_test)
```

```
Out[28]: 0.29080411783731663
```

```
In [29]: lr.score(x_train,y_train)
```

```
Out[29]: 0.3095233169436309
```

## 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]: 0.2907627371075677
```

```
In [33]: rr.score(x_train,y_train)
```

```
Out[33]: 0.3092810921066417
```

## Accuracy(Lasso)

```
In [34]: la=Lasso(alpha=10)  
la.fit(x_train,y_train)
```

```
Out[34]: Lasso(alpha=10)
```

```
In [35]: la.score(x_train,y_train)
```

```
Out[35]: 0.06390078956605738
```

## ElasticNet

```
In [36]: la.score(x_test,y_test)
```

```
Out[36]: 0.06475489702601345
```

```
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([-5.66198376,  1.44191515, -7.68377765,  2.78409068,  0.90664961,  
   -0.05911858, -0.00954366,  1.79892077, -0.0099166 ,  0.23262699,  
   1.49376699,  0.18706343,  1.53340308, -0.76055315])
```

```
In [39]: en.intercept_
```

```
Out[39]: 28079049.123029206
```

```
In [40]: prediction=en.predict(x_test)
```

```
In [41]: en.score(x_test,y_test)
```

```
Out[41]: 0.17468526074111912
```

## 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)))
```

36.80772978720453  
1547.6869797757122  
39.34065301664059

## Logistic Regression

```
In [43]: from sklearn.linear_model import LogisticRegression
```

```
In [44]: feature_matrix=df[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_P10', 'PXY', 'SO_2', 'TCH', 'TOL']]
target_vector=df[ 'station']
```

```
In [45]: feature_matrix.shape
```

Out[45]: (20070, 14)

```
In [46]: target_vector.shape
```

Out[46]: (20070,)

```
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)
```

[28079006]

```
In [52]: logr.classes_
```

Out[52]: array([28079006, 28079024, 28079099], dtype=int64)

```
In [53]: logr.score(fs,target_vector)
```

Out[53]: 0.879023418036871

```
In [54]: logr.predict_proba(observation)[0][0]
```

```
Out[54]: 0.9998967601812779
```

```
In [55]: logr.predict_proba(observation)
```

```
Out[55]: array([[9.99896760e-01, 3.21124597e-30, 1.03239819e-04]])
```

## 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="acc")  
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]: 0.8644028810220575
```

```
In [61]: rfc_best=grid_search.best_estimator_
```

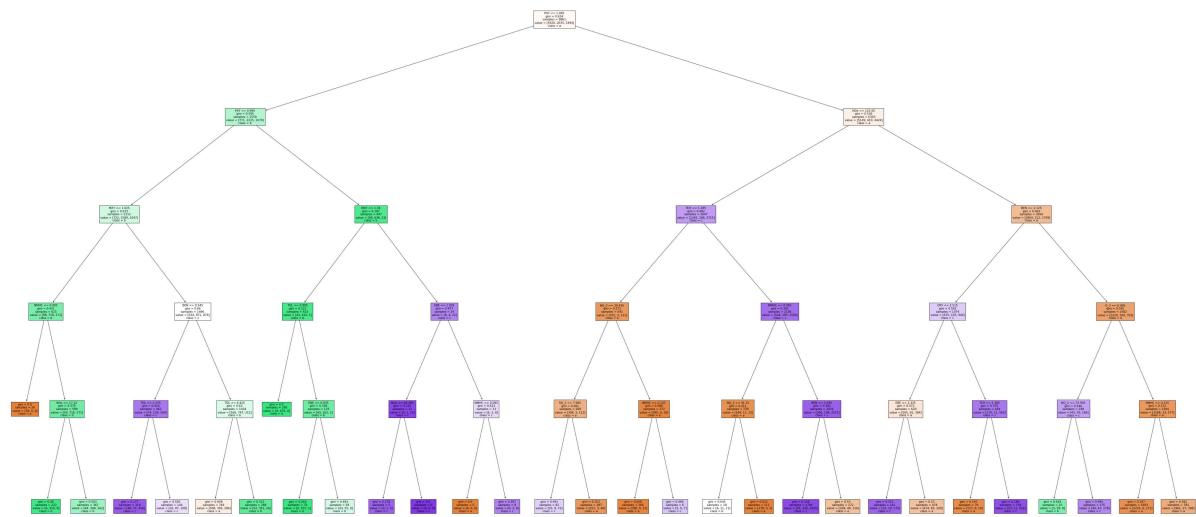
```
In [62]: from sklearn.tree import plot_tree
```

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

```
Out[62]: [Text(2055.7894736842104, 1993.2, 'PXY <= 1.005\ngini = 0.634\nsamples = 8861\nvalue = [5920, 2635, 5494]\nnclass = a'),  
 Text(900.6315789473683, 1630.8000000000002, 'PXY <= 0.995\ngini = 0.595\nsamples = 2558\nvalue = [771, 2225, 1070]\nnclass = b'),  
 Text(430.7368421052631, 1268.4, 'MXY <= 1.025\ngini = 0.633\nsamples = 2111\nvalue = [722, 1589, 1047]\nnclass = b'),  
 Text(156.6315789473684, 906.0, 'NMHC <= 0.005\ngini = 0.421\nsamples = 615\nvalue = [88, 718, 171]\nnclass = b'),  
 Text(78.3157894736842, 543.5999999999999, 'gini = 0.0\nsamples = 26\nvalue = [38, 0, 0]\nnclass = a'),  
 Text(234.9473684210526, 543.5999999999999, 'NOx <= 17.22\ngini = 0.379\nsamples = 589\nvalue = [50, 718, 171]\nnclass = b'),  
 Text(156.6315789473684, 181.1999999999982, 'gini = 0.08\nsamples = 222\nvalue = [6, 350, 9]\nnclass = b'),  
 Text(313.2631578947368, 181.1999999999982, 'gini = 0.503\nsamples = 367\nvalue = [44, 368, 162]\nnclass = b'),  
 Text(704.8421052631578, 906.0, 'BEN <= 0.545\ngini = 0.66\nsamples = 1496\nvalue = [634, 871, 876]\nnclass = c'),  
 Text(548.2105263157895, 543.5999999999999, 'TOL <= 2.375\ngini = 0.416\nsamples = 462\nvalue = [74, 124, 564]\nnclass = c'),  
 Text(469.8947368421052, 181.1999999999982, 'gini = 0.277\nsamples = 318\nvalue = [48, 37, 456]\nnclass = c'),  
 Text(626.5263157894736, 181.1999999999982, 'gini = 0.592\nsamples = 144\nvalue = [26, 87, 108]\nnclass = c'),  
 Text(861.4736842105262, 543.5999999999999, 'TOL <= 4.425\ngini = 0.63\nsamples = 1034\nvalue = [560, 747, 312]\nnclass = b'),  
 Text(783.1578947368421, 181.1999999999982, 'gini = 0.649\nsamples = 746\nvalue = [508, 395, 286]\nnclass = a'),  
 Text(939.7894736842104, 181.1999999999982, 'gini = 0.312\nsamples = 288\nvalue = [52, 352, 26]\nnclass = b'),  
 Text(1370.5263157894735, 1268.4, 'MXY <= 1.34\ngini = 0.187\nsamples = 447\nvalue = [49, 636, 23]\nnclass = b'),  
 Text(1096.421052631579, 906.0, 'TOL <= 0.995\ngini = 0.122\nsamples = 423\nvalue = [43, 632, 1]\nnclass = b'),  
 Text(1018.1052631578947, 543.5999999999999, 'gini = 0.0\nsamples = 290\nvalue = [0, 470, 0]\nnclass = b'),  
 Text(1174.7368421052631, 543.5999999999999, 'EBE <= 0.935\ngini = 0.338\nsamples = 133\nvalue = [43, 162, 1]\nnclass = b'),  
 Text(1096.421052631579, 181.1999999999982, 'gini = 0.018\nsamples = 75\nvalue = [0, 107, 1]\nnclass = b'),  
 Text(1253.0526315789473, 181.1999999999982, 'gini = 0.493\nsamples = 58\nvalue = [43, 55, 0]\nnclass = b'),  
 Text(1644.6315789473683, 906.0, 'EBE <= 1.035\ngini = 0.477\nsamples = 24\nvalue = [6, 4, 22]\nnclass = c'),  
 Text(1488.0, 543.5999999999999, 'NOx <= 44.385\ngini = 0.124\nsamples = 11\nvalue = [0, 1, 14]\nnclass = c'),  
 Text(1409.6842105263156, 181.1999999999982, 'gini = 0.278\nsamples = 5\nvalue = [0, 1, 5]\nnclass = c'),  
 Text(1566.3157894736842, 181.1999999999982, 'gini = 0.0\nsamples = 6\nvalue = [0, 0, 9]\nnclass = c'),  
 Text(1801.2631578947367, 543.5999999999999, 'NMHC <= 0.095\ngini = 0.623\nsamples = 13\nvalue = [6, 3, 8]\nnclass = c'),  
 Text(1722.9473684210525, 181.1999999999982, 'gini = 0.0\nsamples = 5\nvalue = [6, 0, 0]\nnclass = a'),  
 Text(1879.5789473684208, 181.1999999999982, 'gini = 0.397\nsamples = 8\nvalue = [0, 3, 8]\nnclass = c'),  
 Text(3210.9473684210525, 1630.8000000000002, 'NOx <= 122.05\ngini = 0.536\nsamples = 11\nvalue = [122, 0, 5]\nnclass = a')]
```

```
amples = 6303\nvalue = [5149, 410, 4424]\nclass = a'),  
    Text(2584.4210526315787, 1268.4, 'TCH <= 1.285\ngini = 0.482\nsamples = 2647  
\nvalue = [1245, 198, 2715]\nclass = c'),  
    Text(2271.157894736842, 906.0, 'NO_2 <= 39.435\ngini = 0.271\nsamples = 541  
\nvalue = [701, 3, 131]\nclass = a'),  
    Text(2114.5263157894738, 543.5999999999999, 'SO_2 <= 7.065\ngini = 0.401\nsa  
mple = 269\nvalue = [306, 3, 112]\nclass = a'),  
    Text(2036.2105263157894, 181.1999999999982, 'gini = 0.491\nsamples = 82\nva  
lue = [55, 0, 72]\nclass = c'),  
    Text(2192.842105263158, 181.1999999999982, 'gini = 0.253\nsamples = 187\nva  
lue = [251, 3, 40]\nclass = a'),  
    Text(2427.7894736842104, 543.5999999999999, 'NMHC <= 0.135\ngini = 0.088\nsa  
mple = 272\nvalue = [395, 0, 19]\nclass = a'),  
    Text(2349.4736842105262, 181.1999999999982, 'gini = 0.058\nsamples = 266\nv  
alue = [390, 0, 12]\nclass = a'),  
    Text(2506.1052631578946, 181.1999999999982, 'gini = 0.486\nsamples = 6\nval  
ue = [5, 0, 7]\nclass = c'),  
    Text(2897.6842105263154, 906.0, 'NMHC <= 0.065\ngini = 0.365\nsamples = 2106  
\nvalue = [544, 195, 2584]\nclass = c'),  
    Text(2741.052631578947, 543.5999999999999, 'NO_2 <= 41.31\ngini = 0.211\nsa  
mple = 130\nvalue = [184, 11, 13]\nclass = a'),  
    Text(2662.736842105263, 181.1999999999982, 'gini = 0.645\nsamples = 18\nval  
ue = [6, 11, 11]\nclass = b'),  
    Text(2819.368421052631, 181.1999999999982, 'gini = 0.022\nsamples = 112\nva  
lue = [178, 0, 2]\nclass = a'),  
    Text(3054.315789473684, 543.5999999999999, 'BEN <= 1.635\ngini = 0.302\nsa  
mple = 1976\nvalue = [360, 184, 2571]\nclass = c'),  
    Text(2976.0, 181.1999999999982, 'gini = 0.158\nsamples = 1704\nvalue = [91,  
136, 2455]\nclass = c'),  
    Text(3132.6315789473683, 181.1999999999982, 'gini = 0.53\nsamples = 272\nva  
lue = [269, 48, 116]\nclass = a'),  
    Text(3837.4736842105262, 1268.4, 'BEN <= 2.125\ngini = 0.463\nsamples = 3656  
\nvalue = [3904, 212, 1709]\nclass = a'),  
    Text(3524.210526315789, 906.0, 'OXY <= 2.515\ngini = 0.542\nsamples = 1074\nn  
value = [675, 103, 946]\nclass = c'),  
    Text(3367.578947368421, 543.5999999999999, 'EBE <= 1.325\ngini = 0.573\nsa  
mple = 620\nvalue = [505, 92, 384]\nclass = a'),  
    Text(3289.2631578947367, 181.1999999999982, 'gini = 0.321\nsamples = 142\nva  
lue = [31, 10, 175]\nclass = c'),  
    Text(3445.894736842105, 181.1999999999982, 'gini = 0.53\nsamples = 478\nval  
ue = [474, 82, 209]\nclass = a'),  
    Text(3680.8421052631575, 543.5999999999999, 'TCH <= 1.385\ngini = 0.375\nsa  
mple = 454\nvalue = [170, 11, 562]\nclass = c'),  
    Text(3602.5263157894733, 181.1999999999982, 'gini = 0.145\nsamples = 76\nva  
lue = [117, 0, 10]\nclass = a'),  
    Text(3759.1578947368416, 181.1999999999982, 'gini = 0.189\nsamples = 378\nva  
lue = [53, 11, 552]\nclass = c'),  
    Text(4150.736842105262, 906.0, 'O_3 <= 6.095\ngini = 0.345\nsamples = 2582\nn  
value = [3229, 109, 763]\nclass = a'),  
    Text(3994.1052631578946, 543.5999999999999, 'NO_2 <= 72.935\ngini = 0.546\nsa  
mple = 198\nvalue = [43, 76, 186]\nclass = c'),  
    Text(3915.7894736842104, 181.1999999999982, 'gini = 0.429\nsamples = 23\nva  
lue = [3, 29, 8]\nclass = b'),  
    Text(4072.4210526315787, 181.1999999999982, 'gini = 0.495\nsamples = 175\nva  
lue = [40, 47, 178]\nclass = c'),  
    Text(4307.368421052632, 543.5999999999999, 'NMHC <= 0.315\ngini = 0.272\nsa  
mple = 2384\nvalue = [3186, 33, 577]\nclass = a'),
```

```
Text(4229.0526315789475, 181.19999999999982, 'gini = 0.197\nsamples = 1603\nvalue = [2225, 6, 271]\nnclass = a'),  
Text(4385.684210526316, 181.19999999999982, 'gini = 0.392\nsamples = 781\nvalue = [961, 27, 306]\nnclass = a'))]
```



# Conclusion

## Accuracy

In [63]: lr.score(x\_train,y\_train)

Out[63]: 0.3095233169436309

In [64]: rr.score(x\_train,y\_train)

Out[64]: 0.3092810921066417

In [65]: la.score(x\_train,y\_train)

Out[65]: 0.06390078956605738

In [66]: en.score(x\_test,y\_test)

Out[66]: 0.17468526074111912

In [67]: logr.score(fs,target\_vector)

Out[67]: 0.879023418036871

In [68]: grid\_search.best\_score\_

Out[68]: 0.8644028810220575

**Logistic Regression is suitable for this dataset**