

# Importing Libraries

```
In [137]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

# Importing Datasets

```
In [138]: df=pd.read_csv("madrid_2011.csv")
df
```

Out[138]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL
0	2011-11-01 01:00:00	NaN	1.0	NaN	NaN	154.0	84.0	NaN	NaN	NaN	6.0	NaN	NaN 2
1	2011-11-01 01:00:00	2.5	0.4	3.5	0.26	68.0	92.0	3.0	40.0	24.0	9.0	1.54	8.7 2
2	2011-11-01 01:00:00	2.9	NaN	3.8	NaN	96.0	99.0	NaN	NaN	NaN	NaN	NaN	7.2 2
3	2011-11-01 01:00:00	NaN	0.6	NaN	NaN	60.0	83.0	2.0	NaN	NaN	NaN	NaN	NaN 2
4	2011-11-01 01:00:00	NaN	NaN	NaN	NaN	44.0	62.0	3.0	NaN	NaN	3.0	NaN	NaN 2
...	...	...	...	...	...	...	...	...	...	...	...	...	...
209923	2011-09-01 00:00:00	NaN	0.2	NaN	NaN	5.0	19.0	44.0	NaN	NaN	NaN	NaN	NaN 2
209924	2011-09-01 00:00:00	NaN	0.1	NaN	NaN	6.0	29.0	NaN	11.0	NaN	7.0	NaN	NaN 2
209925	2011-09-01 00:00:00	NaN	NaN	NaN	0.23	1.0	21.0	28.0	NaN	NaN	NaN	1.44	NaN 2
209926	2011-09-01 00:00:00	NaN	NaN	NaN	NaN	3.0	15.0	48.0	NaN	NaN	NaN	NaN	NaN 2
209927	2011-09-01 00:00:00	NaN	NaN	NaN	NaN	4.0	33.0	38.0	13.0	NaN	NaN	NaN	NaN 2

209928 rows × 14 columns



# Data Cleaning and Data Preprocessing

```
In [139]: df=df.dropna()
```

```
In [140]: df.columns
```

```
Out[140]: Index(['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',
      'SO_2', 'TCH', 'TOL', 'station'],
                 dtype='object')
```

```
In [141]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 16460 entries, 1 to 209910
Data columns (total 14 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   date      16460 non-null   object 
 1   BEN       16460 non-null   float64
 2   CO        16460 non-null   float64
 3   EBE       16460 non-null   float64
 4   NMHC      16460 non-null   float64
 5   NO        16460 non-null   float64
 6   NO_2      16460 non-null   float64
 7   O_3       16460 non-null   float64
 8   PM10      16460 non-null   float64
 9   PM25      16460 non-null   float64
 10  SO_2      16460 non-null   float64
 11  TCH       16460 non-null   float64
 12  TOL       16460 non-null   float64
 13  station    16460 non-null   int64  
dtypes: float64(12), int64(1), object(1)
memory usage: 1.9+ MB
```

In [142]: `data=df[['CO' , 'station']]  
data`

Out[142]:

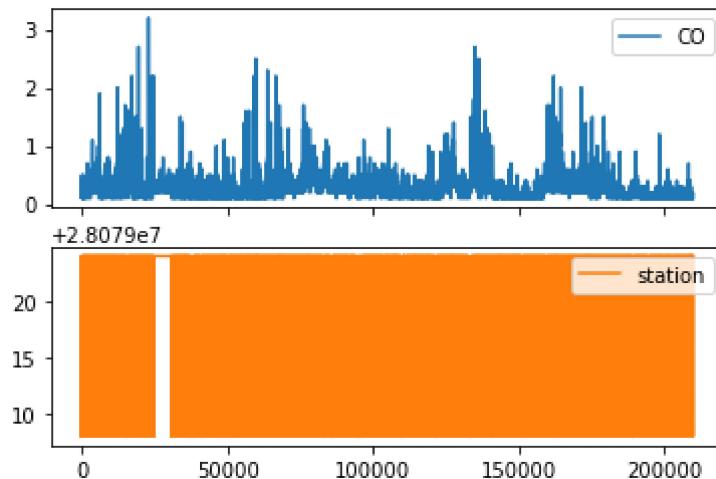
	CO	station
1	0.4	28079008
6	0.3	28079024
25	0.3	28079008
30	0.4	28079024
49	0.2	28079008
...	...	...
209862	0.1	28079024
209881	0.1	28079008
209886	0.1	28079024
209905	0.1	28079008
209910	0.1	28079024

16460 rows × 2 columns

## Line chart

In [143]: `data.plot.line(subplots=True)`

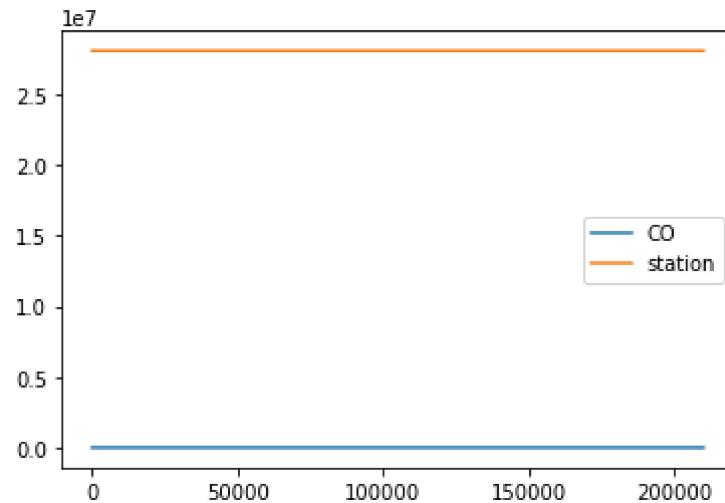
Out[143]: `array([<AxesSubplot:>, <AxesSubplot:>], dtype=object)`



## Line chart

```
In [144]: data.plot.line()
```

```
Out[144]: <AxesSubplot:>
```

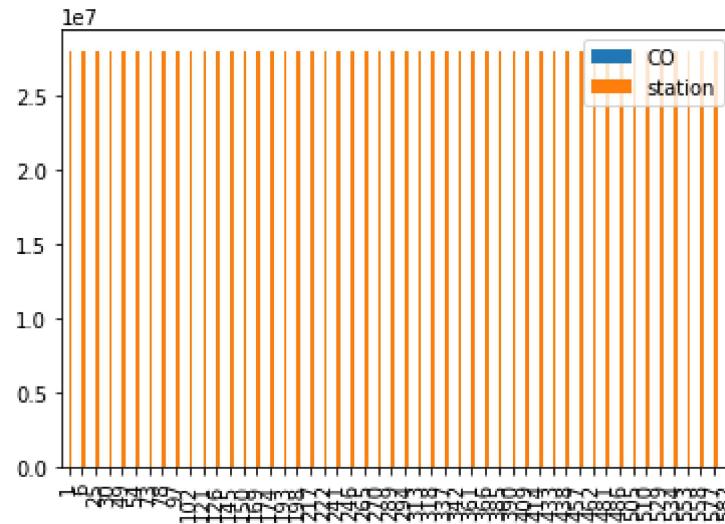


## Bar chart

```
In [145]: b=data[0:50]
```

```
In [146]: b.plot.bar()
```

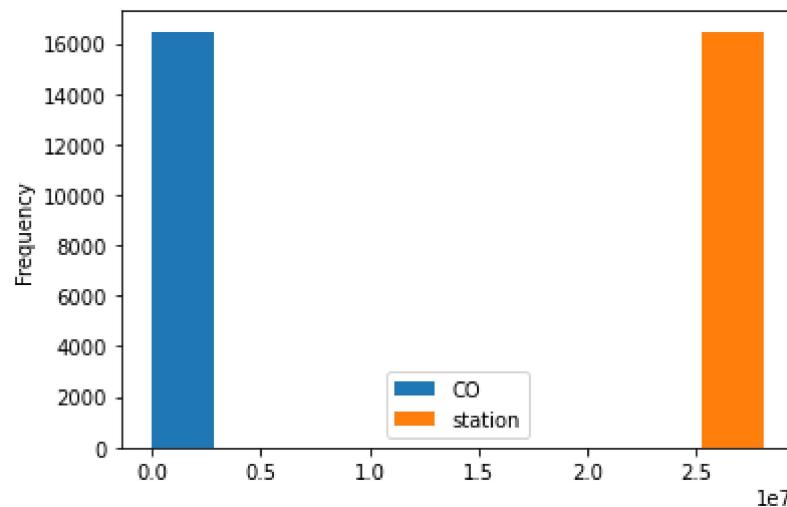
```
Out[146]: <AxesSubplot:>
```



## Histogram

```
In [147]: data.plot.hist()
```

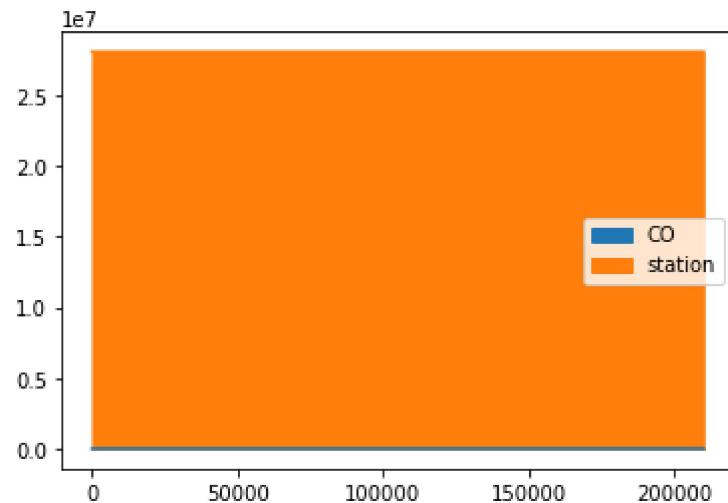
```
Out[147]: <AxesSubplot:ylabel='Frequency'>
```



## Area chart

```
In [148]: data.plot.area()
```

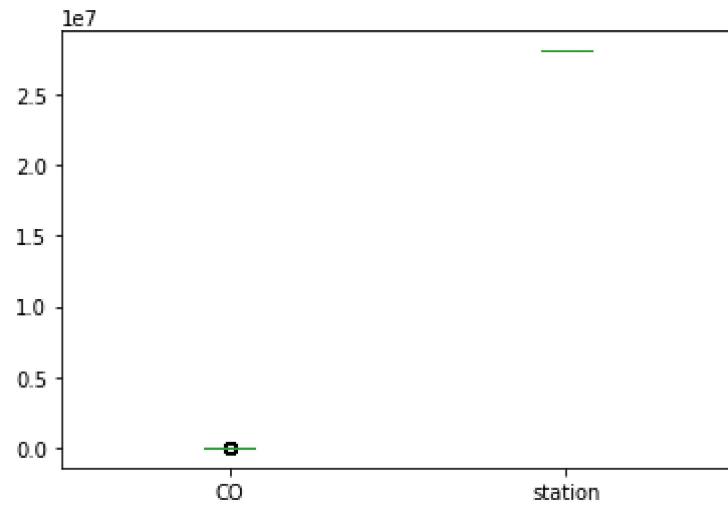
```
Out[148]: <AxesSubplot:>
```



## Box chart

In [149]: `data.plot.box()`

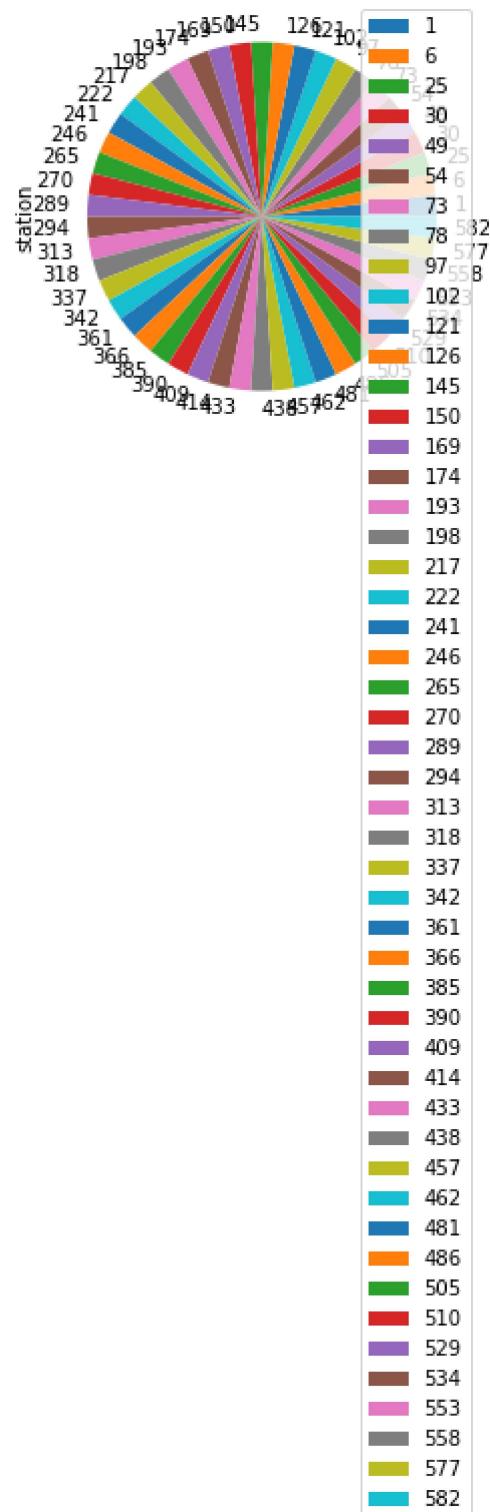
Out[149]: <AxesSubplot:>



## Pie chart

```
In [150]: b.plot.pie(y='station' )
```

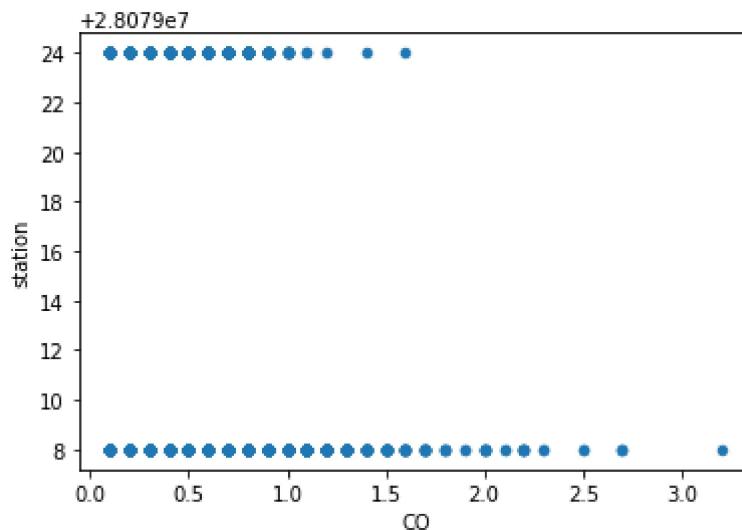
```
Out[150]: <AxesSubplot:ylabel='station'>
```



## Scatter chart

```
In [151]: data.plot.scatter(x='CO' ,y='station')
```

```
Out[151]: <AxesSubplot:xlabel='CO', ylabel='station'>
```



```
In [152]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 16460 entries, 1 to 209910
Data columns (total 14 columns):
 #   Column   Non-Null Count  Dtype  
--- 
 0   date      16460 non-null   object 
 1   BEN       16460 non-null   float64
 2   CO        16460 non-null   float64
 3   EBE       16460 non-null   float64
 4   NMHC      16460 non-null   float64
 5   NO        16460 non-null   float64
 6   NO_2      16460 non-null   float64
 7   O_3       16460 non-null   float64
 8   PM10      16460 non-null   float64
 9   PM25      16460 non-null   float64
 10  SO_2      16460 non-null   float64
 11  TCH       16460 non-null   float64
 12  TOL       16460 non-null   float64
 13  station    16460 non-null   int64
```

```
In [153]: df.describe()
```

Out[153]:

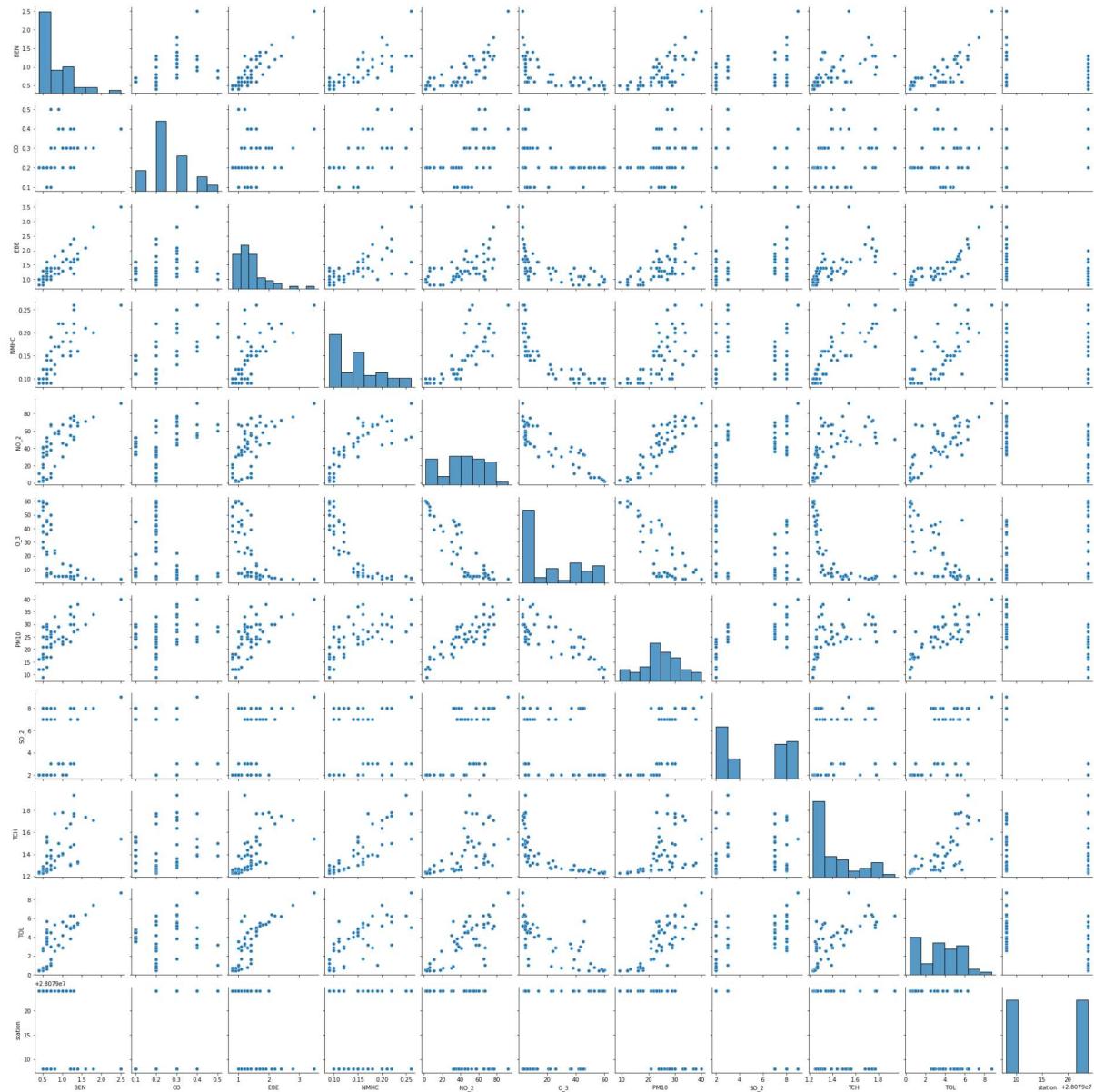
	BEN	CO	EBE	NMHC	NO	NO_2	NO_x
count	16460.000000	16460.000000	16460.000000	16460.000000	16460.000000	16460.000000	16460.000000
mean	0.900680	0.277758	1.471871	0.167043	23.671810	44.583961	1.44
std	0.768892	0.206143	1.051004	0.075068	44.362859	31.569185	1.44
min	0.100000	0.100000	0.200000	0.010000	1.000000	1.000000	0.00
25%	0.500000	0.200000	0.800000	0.120000	2.000000	19.000000	0.00
50%	0.700000	0.200000	1.200000	0.160000	7.000000	40.000000	0.00
75%	1.100000	0.300000	1.700000	0.200000	25.000000	63.000000	0.00
max	9.500000	3.200000	12.800000	0.840000	615.000000	289.000000	11.00

```
In [154]: df1=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',  
'PM10', 'SO_2', 'TCH', 'TOL', 'station']]
```

## EDA AND VISUALIZATION

```
In [155]: sns.pairplot(df1[0:50])
```

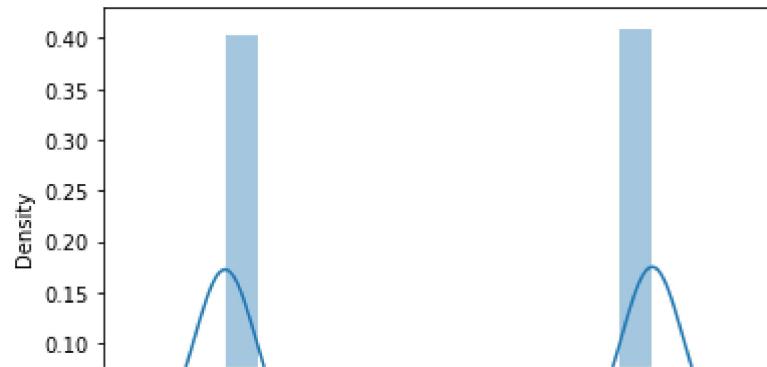
```
Out[155]: <seaborn.axisgrid.PairGrid at 0x2d01b5453d0>
```



```
In [156]: sns.distplot(df1['station'])
```

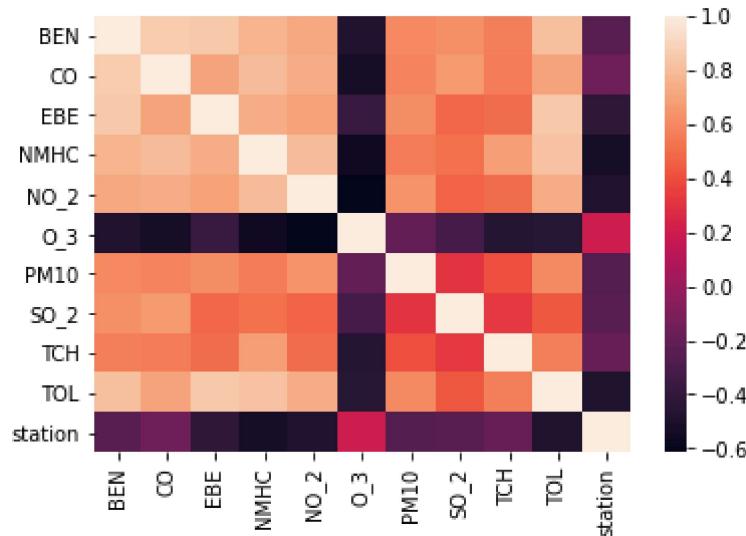
C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function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

```
Out[156]: <AxesSubplot:xlabel='station', ylabel='Density'>
```



```
In [157]: sns.heatmap(df1.corr())
```

```
Out[157]: <AxesSubplot:>
```



## TO TRAIN THE MODEL AND MODEL BUILDING

```
In [158]: x=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',
           'PM10', 'SO_2', 'TCH', 'TOL']]
y=df['station']
```

```
In [159]: 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 [160]: from sklearn.linear_model import LinearRegression  
lr=LinearRegression()  
lr.fit(x_train,y_train)
```

```
Out[160]: LinearRegression()
```

```
In [161]: lr.intercept_
```

```
Out[161]: 28079018.269854557
```

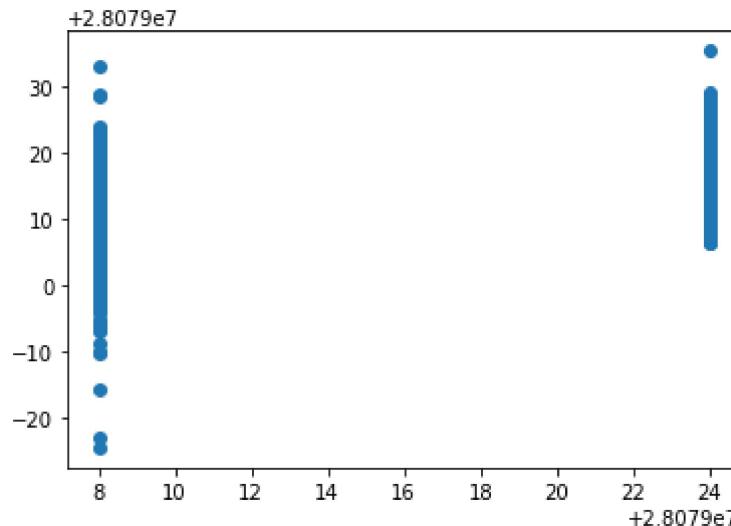
```
In [162]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])  
coeff
```

```
Out[162]:
```

	Co-efficient
BEN	4.017502
CO	31.986266
EBE	-1.905570
NMHC	-94.561288
NO_2	-0.088893
O_3	-0.016792
PM10	0.004594
SO_2	-0.520397
TCH	9.675602
TOL	-0.476223

```
In [163]: prediction = lr.predict(x_test)
plt.scatter(y_test,prediction)
```

```
Out[163]: <matplotlib.collections.PathCollection at 0x2d0267d90d0>
```



## ACCURACY

```
In [164]: lr.score(x_test,y_test)
```

```
Out[164]: 0.611007100086842
```

```
In [165]: lr.score(x_train,y_train)
```

```
Out[165]: 0.6264125273757513
```

## Ridge and Lasso

```
In [166]: from sklearn.linear_model import Ridge,Lasso
```

```
In [167]: rr=Ridge(alpha=10)
rr.fit(x_train,y_train)
```

```
Out[167]: Ridge(alpha=10)
```

## Accuracy(Ridge)

```
In [168]: rr.score(x_test,y_test)
```

```
Out[168]: 0.5759505750015184
```

```
In [169]: rr.score(x_train,y_train)
```

```
Out[169]: 0.5932042205659791
```

## Accuracy(Lasso)

```
In [170]: la=Lasso(alpha=10)  
la.fit(x_train,y_train)
```

```
Out[170]: Lasso(alpha=10)
```

```
In [171]: la.score(x_train,y_train)
```

```
Out[171]: 0.23744684486461143
```

## ElasticNet

```
In [172]: la.score(x_test,y_test)
```

```
Out[172]: 0.2299586069038282
```

```
In [173]: from sklearn.linear_model import ElasticNet  
en=ElasticNet()  
en.fit(x_train,y_train)
```

```
Out[173]: ElasticNet()
```

```
In [174]: en.coef_
```

```
Out[174]: array([ 0.63905176,  0.          , -0.          , -0.          , -0.13475689,  
                  -0.05397622,  0.07612869, -0.          ,  0.          , -0.71451597])
```

```
In [175]: en.intercept_
```

```
Out[175]: 28079024.351645764
```

```
In [176]: prediction=en.predict(x_test)
```

```
In [177]: en.score(x_test,y_test)
```

```
Out[177]: 0.3172206432771122
```

## Evaluation Metrics

```
In [178]: 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)))
```

```
5.849244923141863  
43.69729819528614  
6.610393195210564
```

## Logistic Regression

```
In [179]: from sklearn.linear_model import LogisticRegression
```

```
In [180]: feature_matrix=df[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',  
    'PM10', 'SO_2', 'TCH', 'TOL']]  
target_vector=df['station']
```

```
In [181]: feature_matrix.shape
```

```
Out[181]: (16460, 10)
```

```
In [182]: target_vector.shape
```

```
Out[182]: (16460,)
```

```
In [183]: from sklearn.preprocessing import StandardScaler
```

```
In [184]: fs=StandardScaler().fit_transform(feature_matrix)
```

```
In [185]: logr=LogisticRegression(max_iter=10000)  
logr.fit(fs,target_vector)
```

```
Out[185]: LogisticRegression(max_iter=10000)
```

```
In [186]: observation=[[1,2,3,4,5,6,7,8,9,10]]
```

```
In [187]: prediction=logr.predict(observation)  
print(prediction)
```

```
[28079008]
```

```
In [188]: logr.classes_
```

```
Out[188]: array([28079008, 28079024], dtype=int64)
```

```
In [189]: logr.score(fs,target_vector)
```

```
Out[189]: 0.9237545565006076
```

```
In [190]: logr.predict_proba(observation)[0][0]
```

```
Out[190]: 0.9999999999999966
```

```
In [191]: logr.predict_proba(observation)
```

```
Out[191]: array([[1.0000000e+00, 3.47334507e-15]])
```

## Random Forest

```
In [192]: from sklearn.ensemble import RandomForestClassifier
```

```
In [193]: rfc=RandomForestClassifier()  
rfc.fit(x_train,y_train)
```

```
Out[193]: RandomForestClassifier()
```

```
In [194]: parameters={'max_depth':[1,2,3,4,5],  
                   'min_samples_leaf':[5,10,15,20,25],  
                   'n_estimators':[10,20,30,40,50]  
}
```

```
In [195]: 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[195]: 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 [196]: grid_search.best_score_
```

```
Out[196]: 0.9339524388127061
```

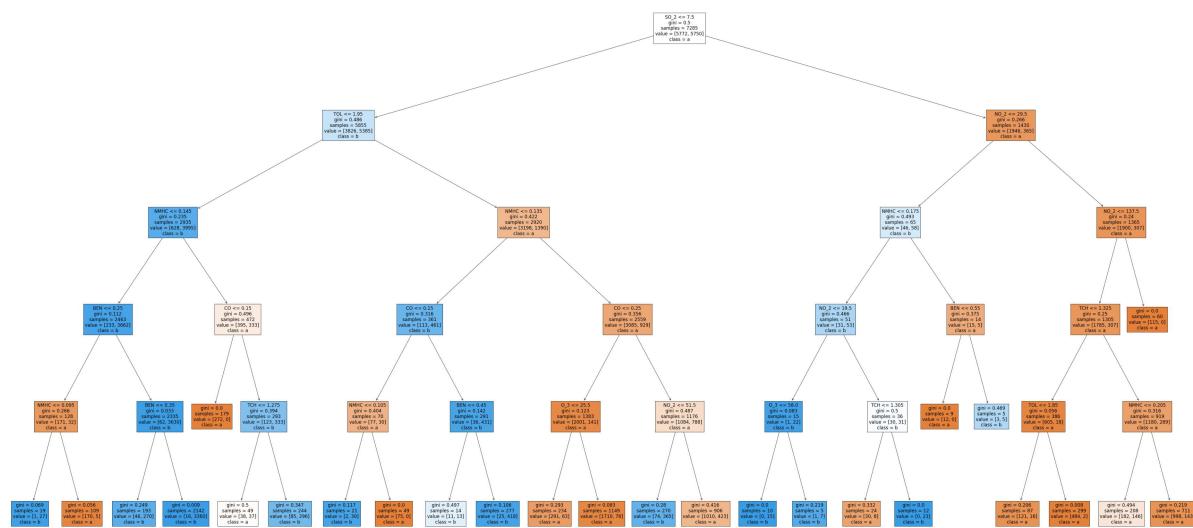
```
In [197]: rfc_best=grid_search.best_estimator_
```

```
In [198]: 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[198]: [Text(2523.1304347826085, 1993.2, 'SO_2 <= 7.5\ngini = 0.5\nsamples = 7285\nvalue = [5772, 5750]\nclass = a'),
Text(1285.8260869565217, 1630.8000000000002, 'TOL <= 1.95\ngini = 0.486\nsamples = 5855\nvalue = [3826, 5385]\nclass = b'),
Text(630.7826086956521, 1268.4, 'NMHC <= 0.145\ngini = 0.235\nsamples = 2935\nvalue = [628, 3995]\nclass = b'),
Text(388.17391304347825, 906.0, 'BEN <= 0.25\ngini = 0.112\nsamples = 2463\nvalue = [233, 3662]\nclass = b'),
Text(194.08695652173913, 543.5999999999999, 'NMHC <= 0.095\ngini = 0.266\nsamples = 128\nvalue = [171, 32]\nclass = a'),
Text(97.04347826086956, 181.1999999999982, 'gini = 0.069\nsamples = 19\nvalue = [1, 27]\nclass = b'),
Text(291.1304347826087, 181.1999999999982, 'gini = 0.056\nsamples = 109\nvalue = [170, 5]\nclass = a'),
Text(582.2608695652174, 543.5999999999999, 'BEN <= 0.35\ngini = 0.033\nsamples = 2335\nvalue = [62, 3630]\nclass = b'),
Text(485.2173913043478, 181.1999999999982, 'gini = 0.249\nsamples = 193\nvalue = [46, 270]\nclass = b'),
Text(679.304347826087, 181.1999999999982, 'gini = 0.009\nsamples = 2142\nvalue = [16, 3360]\nclass = b'),
Text(873.391304347826, 906.0, 'CO <= 0.15\ngini = 0.496\nsamples = 472\nvalue = [395, 333]\nclass = a'),
Text(776.3478260869565, 543.5999999999999, 'gini = 0.0\nsamples = 179\nvalue = [272, 0]\nclass = a'),
Text(970.4347826086956, 543.5999999999999, 'TCH <= 1.275\ngini = 0.394\nsamples = 293\nvalue = [123, 333]\nclass = b'),
Text(873.391304347826, 181.1999999999982, 'gini = 0.5\nsamples = 49\nvalue = [38, 37]\nclass = a'),
Text(1067.4782608695652, 181.1999999999982, 'gini = 0.347\nsamples = 244\nvalue = [85, 296]\nclass = b'),
Text(1940.8695652173913, 1268.4, 'NMHC <= 0.135\ngini = 0.422\nsamples = 2920\nvalue = [3198, 1390]\nclass = a'),
Text(1552.695652173913, 906.0, 'CO <= 0.15\ngini = 0.316\nsamples = 361\nvalue = [113, 461]\nclass = b'),
Text(1358.608695652174, 543.5999999999999, 'NMHC <= 0.105\ngini = 0.404\nsamples = 70\nvalue = [77, 30]\nclass = a'),
Text(1261.5652173913043, 181.1999999999982, 'gini = 0.117\nsamples = 21\nvalue = [2, 30]\nclass = b'),
Text(1455.6521739130435, 181.1999999999982, 'gini = 0.0\nsamples = 49\nvalue = [75, 0]\nclass = a'),
Text(1746.782608695652, 543.5999999999999, 'BEN <= 0.45\ngini = 0.142\nsamples = 291\nvalue = [36, 431]\nclass = b'),
Text(1649.7391304347825, 181.1999999999982, 'gini = 0.497\nsamples = 14\nvalue = [11, 13]\nclass = b'),
Text(1843.8260869565217, 181.1999999999982, 'gini = 0.106\nsamples = 277\nvalue = [25, 418]\nclass = b'),
Text(2329.0434782608695, 906.0, 'CO <= 0.25\ngini = 0.356\nsamples = 2559\nvalue = [3085, 929]\nclass = a'),
Text(2134.9565217391305, 543.5999999999999, 'O_3 <= 25.5\ngini = 0.123\nsamples = 1383\nvalue = [2001, 141]\nclass = a'),
Text(2037.9130434782608, 181.1999999999982, 'gini = 0.293\nsamples = 234\nvalue = [291, 63]\nclass = a'),
Text(2232.0, 181.1999999999982, 'gini = 0.083\nsamples = 1149\nvalue = [1710, 78]\nclass = a'),
Text(2523.1304347826085, 543.5999999999999, 'NO_2 <= 51.5\ngini = 0.487\nsamples = 1176\nvalue = [1084, 788]\nclass = a'),
Text(2426.086956521739, 181.1999999999982, 'gini = 0.28\nsamples = 270\nvalue = [1084, 788]\nclass = a')]
```

```
ue = [74, 365]\nclass = b'),  
    Text(2620.173913043478, 181.19999999999982, 'gini = 0.416\ncount = 906\nvalue = [1010, 423]\nclass = a'),  
    Text(3760.4347826086955, 1630.8000000000002, 'NO_2 <= 29.5\ngini = 0.266\nsamples = 1430\nvalue = [1946, 365]\nclass = a'),  
    Text(3348.0, 1268.4, 'NMHC <= 0.175\ngini = 0.493\nsamples = 65\nvalue = [46, 58]\nclass = b'),  
    Text(3105.391304347826, 906.0, 'NO_2 <= 19.5\ngini = 0.466\nsamples = 51\nvalue = [31, 53]\nclass = b'),  
    Text(2911.304347826087, 543.5999999999999, 'O_3 <= 58.0\ngini = 0.083\nsamples = 15\nvalue = [1, 22]\nclass = b'),  
    Text(2814.2608695652175, 181.19999999999982, 'gini = 0.0\nsamples = 10\nvalue = [0, 15]\nclass = b'),  
    Text(3008.3478260869565, 181.19999999999982, 'gini = 0.219\nsamples = 5\nvalue = [1, 7]\nclass = b'),  
    Text(3299.478260869565, 543.5999999999999, 'TCH <= 1.305\ngini = 0.5\nsamples = 36\nvalue = [30, 31]\nclass = b'),  
    Text(3202.4347826086955, 181.19999999999982, 'gini = 0.332\nsamples = 24\nvalue = [30, 8]\nclass = a'),  
    Text(3396.5217391304345, 181.19999999999982, 'gini = 0.0\nsamples = 12\nvalue = [0, 23]\nclass = b'),  
    Text(3590.608695652174, 906.0, 'BEN <= 0.55\ngini = 0.375\nsamples = 14\nvalue = [15, 5]\nclass = a'),  
    Text(3493.565217391304, 543.5999999999999, 'gini = 0.0\nsamples = 9\nvalue = [12, 0]\nclass = a'),  
    Text(3687.6521739130435, 543.5999999999999, 'gini = 0.469\nsamples = 5\nvalue = [3, 5]\nclass = b'),  
    Text(4172.869565217391, 1268.4, 'NO_2 <= 137.5\ngini = 0.24\nsamples = 1365\nvalue = [1900, 307]\nclass = a'),  
    Text(4075.8260869565215, 906.0, 'TCH <= 1.325\ngini = 0.25\nsamples = 1305\nvalue = [1785, 307]\nclass = a'),  
    Text(3881.7391304347825, 543.5999999999999, 'TOL <= 1.85\ngini = 0.056\nsamples = 386\nvalue = [605, 18]\nclass = a'),  
    Text(3784.695652173913, 181.19999999999982, 'gini = 0.206\nsamples = 87\nvalue = [121, 16]\nclass = a'),  
    Text(3978.782608695652, 181.19999999999982, 'gini = 0.008\nsamples = 299\nvalue = [484, 2]\nclass = a'),  
    Text(4269.913043478261, 543.5999999999999, 'NMHC <= 0.205\ngini = 0.316\nsamples = 919\nvalue = [1180, 289]\nclass = a'),  
    Text(4172.869565217391, 181.19999999999982, 'gini = 0.494\nsamples = 208\nvalue = [182, 146]\nclass = a'),  
    Text(4366.95652173913, 181.19999999999982, 'gini = 0.219\nsamples = 711\nvalue = [998, 143]\nclass = a'),  
    Text(4269.913043478261, 906.0, 'gini = 0.0\nsamples = 60\nvalue = [115, 0]\nclass = a')]
```



## Conclusion

### Accuracy

```
In [199]: lr.score(x_train,y_train)
```

```
Out[199]: 0.6264125273757513
```

```
In [200]: rr.score(x_train,y_train)
```

```
Out[200]: 0.5932042205659791
```

```
In [201]: la.score(x_train,y_train)
```

```
Out[201]: 0.23744684486461143
```

```
In [202]: en.score(x_test,y_test)
```

```
Out[202]: 0.3172206432771122
```

```
In [203]: logr.score(fs,target_vector)
```

```
Out[203]: 0.9237545565006076
```

```
In [204]: grid_search.best_score_
```

```
Out[204]: 0.9339524388127061
```

**Random is suitable for this dataset**

