

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: df = pd.read_csv(r"C:\Users\user\Downloads\C10_air\csvs_per_year\csvs_per_year\
df
```

Out[2]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL	
0	2014-06-01 01:00:00	NaN	0.2	NaN	NaN	3.0	10.0	NaN	NaN	NaN	3.0	NaN	NaN	28
1	2014-06-01 01:00:00	0.2	0.2	0.1	0.11	3.0	17.0	68.0	10.0	5.0	5.0	1.36	1.3	28
2	2014-06-01 01:00:00	0.3	NaN	0.1	NaN	2.0	6.0	NaN	NaN	NaN	NaN	NaN	1.1	28
3	2014-06-01 01:00:00	NaN	0.2	NaN	NaN	1.0	6.0	79.0	NaN	NaN	NaN	NaN	NaN	28
4	2014-06-01 01:00:00	NaN	NaN	NaN	NaN	1.0	6.0	75.0	NaN	NaN	4.0	NaN	NaN	28
...
210019	2014-09-01 00:00:00	NaN	0.5	NaN	NaN	20.0	84.0	29.0	NaN	NaN	NaN	NaN	NaN	28
210020	2014-09-01 00:00:00	NaN	0.3	NaN	NaN	1.0	22.0	NaN	15.0	NaN	6.0	NaN	NaN	28
210021	2014-09-01 00:00:00	NaN	NaN	NaN	NaN	1.0	13.0	70.0	NaN	NaN	NaN	NaN	NaN	28
210022	2014-09-01 00:00:00	NaN	NaN	NaN	NaN	3.0	38.0	42.0	NaN	NaN	NaN	NaN	NaN	28
210023	2014-09-01 00:00:00	NaN	NaN	NaN	NaN	1.0	26.0	65.0	11.0	NaN	NaN	NaN	NaN	28

210024 rows × 14 columns



```
In [3]: df1 = df.fillna(0)
df1
```

Out[3]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL	s
0	2014-06-01 01:00:00	0.0	0.2	0.0	0.00	3.0	10.0	0.0	0.0	0.0	3.0	0.00	0.0	280
1	2014-06-01 01:00:00	0.2	0.2	0.1	0.11	3.0	17.0	68.0	10.0	5.0	5.0	1.36	1.3	280
2	2014-06-01 01:00:00	0.3	0.0	0.1	0.00	2.0	6.0	0.0	0.0	0.0	0.0	0.00	1.1	280
3	2014-06-01 01:00:00	0.0	0.2	0.0	0.00	1.0	6.0	79.0	0.0	0.0	0.0	0.00	0.0	280
4	2014-06-01 01:00:00	0.0	0.0	0.0	0.00	1.0	6.0	75.0	0.0	0.0	4.0	0.00	0.0	280
...
210019	2014-09-01 00:00:00	0.0	0.5	0.0	0.00	20.0	84.0	29.0	0.0	0.0	0.0	0.00	0.0	280
210020	2014-09-01 00:00:00	0.0	0.3	0.0	0.00	1.0	22.0	0.0	15.0	0.0	6.0	0.00	0.0	280
210021	2014-09-01 00:00:00	0.0	0.0	0.0	0.00	1.0	13.0	70.0	0.0	0.0	0.0	0.00	0.0	280
210022	2014-09-01 00:00:00	0.0	0.0	0.0	0.00	3.0	38.0	42.0	0.0	0.0	0.0	0.00	0.0	280
210023	2014-09-01 00:00:00	0.0	0.0	0.0	0.00	1.0	26.0	65.0	11.0	0.0	0.0	0.00	0.0	280

210024 rows × 14 columns

```
In [4]: df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 210024 entries, 0 to 210023
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   date        210024 non-null  object
 1   BEN         210024 non-null  float64
 2   CO          210024 non-null  float64
 3   EBE         210024 non-null  float64
 4   NMHC        210024 non-null  float64
 5   NO          210024 non-null  float64
 6   NO_2        210024 non-null  float64
 7   O_3         210024 non-null  float64
 8   PM10        210024 non-null  float64
 9   PM25        210024 non-null  float64
10   SO_2        210024 non-null  float64
11   TCH         210024 non-null  float64
12   TOL         210024 non-null  float64
13   station     210024 non-null  int64
dtypes: float64(12), int64(1), object(1)
memory usage: 22.4+ MB
```

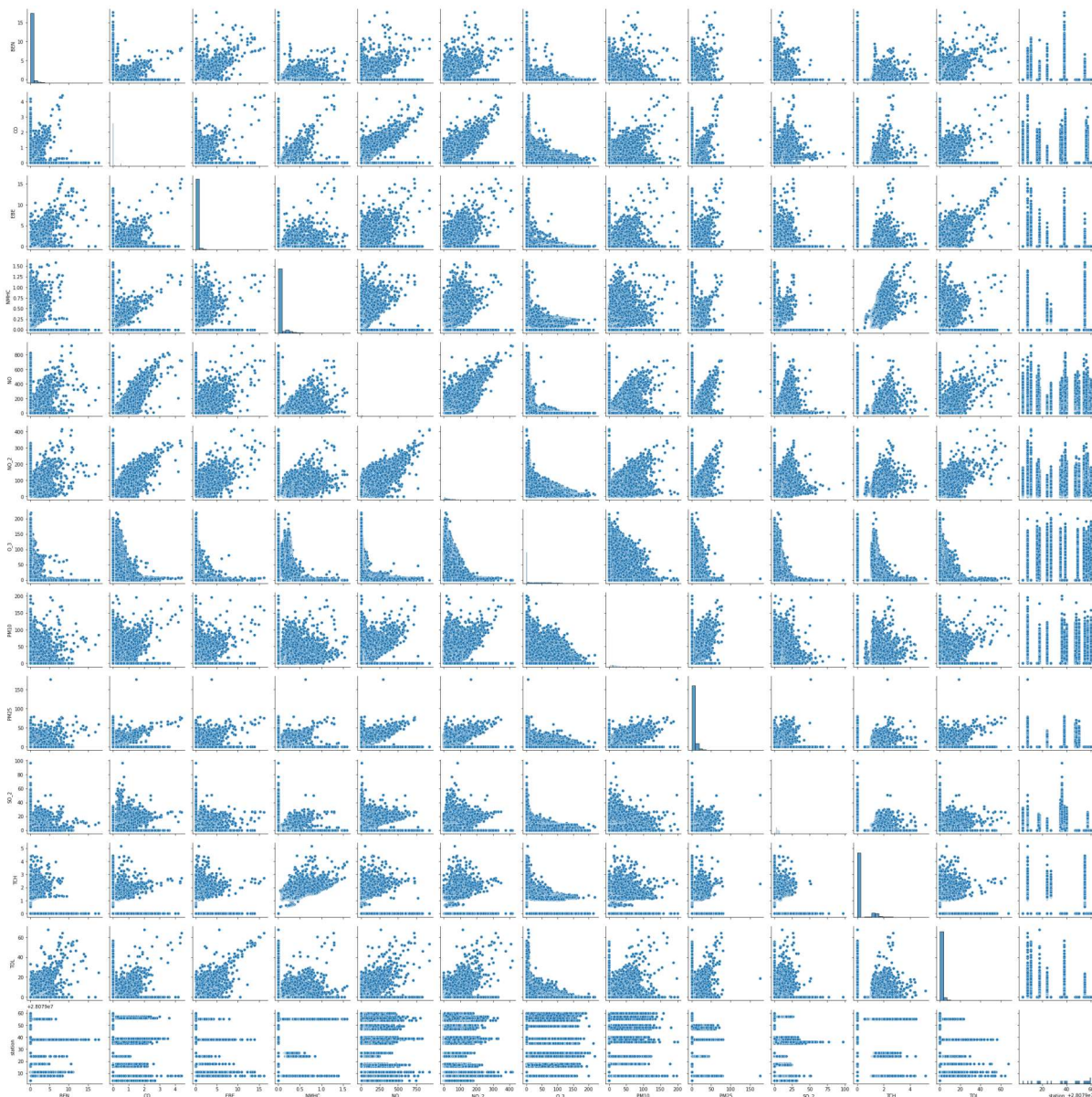
```
In [5]: df1.columns
```

```
Out[5]: Index(['date', 'BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',
              'SO_2', 'TCH', 'TOL', 'station'],
              dtype='object')
```

```
In [6]: df2 = df1[['BEN', 'CO', 'EBE', 'NMHC', 'NO', 'NO_2', 'O_3', 'PM10', 'PM25',
                  'SO_2', 'TCH', 'TOL', 'station']]
```

```
In [7]: sns.pairplot(df2)
```

```
Out[7]: <seaborn.axisgrid.PairGrid at 0x1e25dd68070>
```

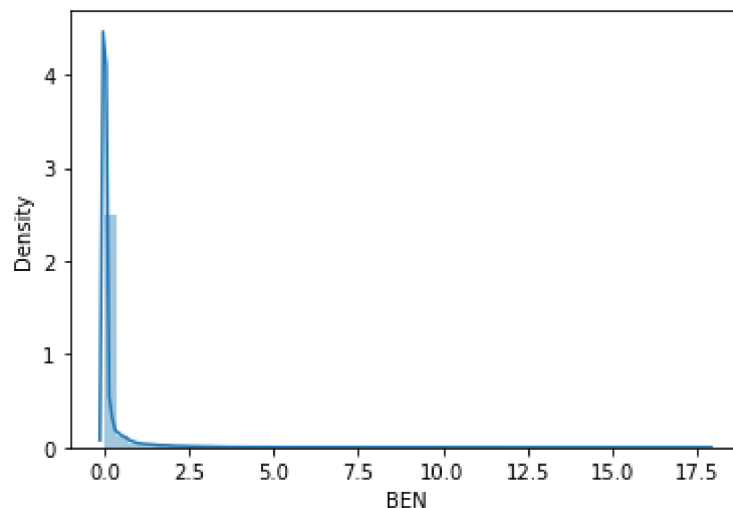


```
In [8]: sns.distplot(df2['BEN'])
```

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).

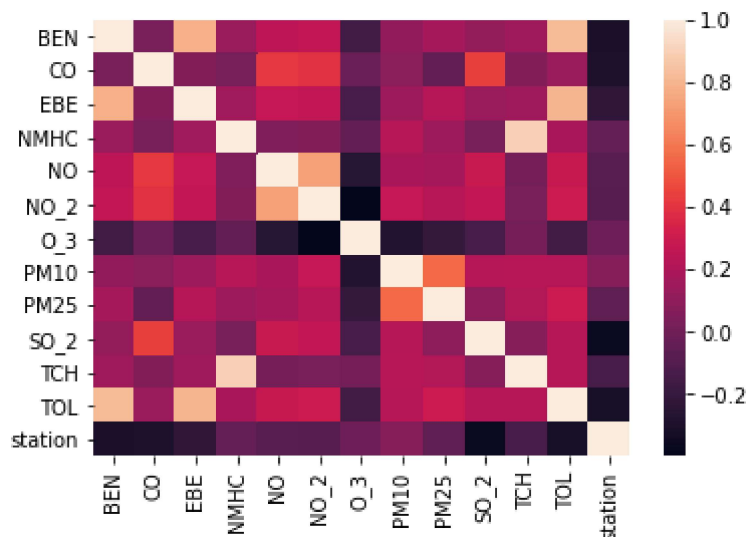
```
warnings.warn(msg, FutureWarning)
```

```
Out[8]: <AxesSubplot:xlabel='BEN', ylabel='Density'>
```



```
In [9]: sns.heatmap(df2.corr())
```

```
Out[9]: <AxesSubplot:>
```



Linear Regression

```
In [10]: x = df2[['BEN', 'CO', 'EBE', 'NMHC', 'NO_2', 'O_3',  
               'PM10', 'SO_2', 'TCH']]  
y = df2['TOL']
```

```
In [11]: from sklearn.model_selection import train_test_split  
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.30)
```

```
In [12]: from sklearn.linear_model import LinearRegression  
  
lr = LinearRegression()  
lr.fit(x_train,y_train)
```

Out[12]: LinearRegression()

```
In [13]: print(lr.intercept_)  
  
-0.2079872868055742
```

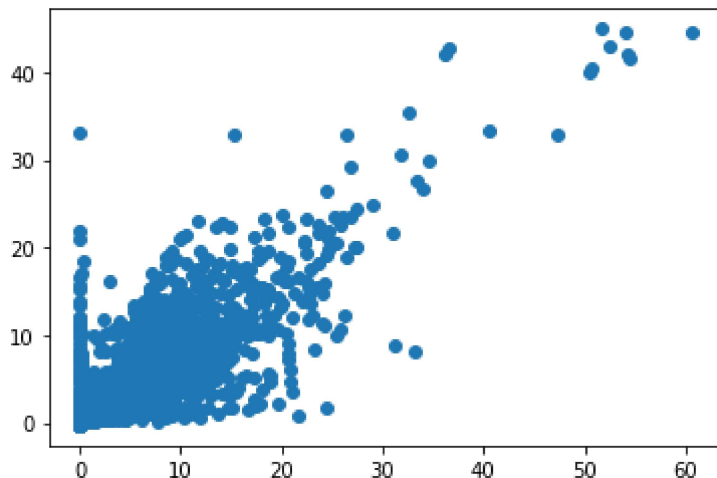
```
In [14]: coeff = pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])  
coeff
```

Out[14]:

	Co-efficient
BEN	1.882933
CO	0.471901
EBE	1.792157
NMHC	-1.031286
NO_2	0.000457
O_3	0.000544
PM10	0.011472
SO_2	0.037908
TCH	0.470564

```
In [15]: prediction = lr.predict(x_test)
plt.scatter(y_test, prediction)
```

Out[15]: <matplotlib.collections.PathCollection at 0x1e2046f8550>



```
In [16]: print(lr.score(x_test,y_test))
```

0.7778147341001751

```
In [17]: lr.score(x_train,y_train)
```

Out[17]: 0.770048129954465

Ridge and Lasso

```
In [18]: from sklearn.linear_model import Ridge,Lasso
```

```
In [19]: rr = Ridge(alpha=10)
rr.fit(x_train,y_train)
rr.score(x_train,y_train)
```

Out[19]: 0.7700473851362977

```
In [20]: rr.score(x_test,y_test)
```

Out[20]: 0.7778131325203637

Lasso Regression

```
In [21]: ls = Lasso(alpha=10)
ls.fit(x_train,y_train)
ls.score(x_train,y_train)
```

Out[21]: 0.05630626705848085

```
In [22]: ls.score(x_test,y_test)
```

```
Out[22]: 0.05870826224256287
```

ElasticNET regression

```
In [23]: from sklearn.linear_model import ElasticNet  
es = ElasticNet()  
es.fit(x_train,y_train)
```

```
Out[23]: ElasticNet()
```

```
In [24]: print(es.coef_)
```

```
[ 3.09871124e-01  0.00000000e+00  7.81127807e-02  0.00000000e+00  
 1.48616930e-02 -1.37716937e-04  1.67379460e-02  3.13069169e-02  
 0.00000000e+00]
```

```
In [25]: print(es.intercept_)
```

```
-0.19450666150573948
```

```
In [26]: print(es.score(x_test,y_test))
```

```
0.2611329784999282
```

```
In [27]: print(es.score(x_train,y_train))
```

```
0.2548407570864988
```

LogisticRegression

```
In [28]: from sklearn.linear_model import LogisticRegression
```

```
In [29]: feature_matrix = df2.iloc[:,0:15]  
target_vector = df2.iloc[:,-1]
```

```
In [30]: feature_matrix.shape
```

```
Out[30]: (210024, 13)
```

```
In [31]: from sklearn.preprocessing import StandardScaler
```

```
In [32]: fs = StandardScaler().fit_transform(feature_matrix)
```



```
In [33]: logs = LogisticRegression()
logs.fit(fs,target_vector)
```

C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model_logistic.py: 763: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)
n_iter_i = _check_optimize_result(

```
Out[33]: LogisticRegression()
```

```
In [35]: observation = [[1.4,1.5,1.6,2.7,2.3,3.3,2.3,4.1,2.3,4.2,1.2,12,2]]
prediction = logs.predict(observation)
```

```
In [36]: print(prediction)

[28079059]
```

```
In [37]: logs.classes_
```

```
Out[37]: array([28079004, 28079008, 28079011, 28079016, 28079017, 28079018,
                28079024, 28079027, 28079035, 28079036, 28079038, 28079039,
                28079040, 28079047, 28079048, 28079049, 28079050, 28079054,
                28079055, 28079056, 28079057, 28079058, 28079059, 28079060],
                dtype=int64)
```

```
In [38]: from sklearn.model_selection import train_test_split

x_train,x_test,y_train,y_test = train_test_split(feature_matrix,target_vector,t
```

```
In [39]: print(logs.score(x_test,y_test))

0.04161376333164043
```

```
In [40]: print(logs.score(x_train,y_train))

0.04175055776242042
```

Conclusion

linear regression is bestfit model

linear regression is best fit model for dataset madrid_2001. The score of x_train,y_train is 0.770048129954465 and x_test and y_test score is 0.7778147341001751.

In []: