

```
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	MXV	NMHC	NO_2	NOx	OXY	O_3	
0	2007-12-01 01:00:00	NaN	2.86	NaN	NaN	NaN	282.200012	1054.000000	NaN	4.030000	156.1
1	2007-12-01 01:00:00	NaN	1.82	NaN	NaN	NaN	86.419998	354.600006	NaN	3.260000	80.8
2	2007-12-01 01:00:00	NaN	1.47	NaN	NaN	NaN	94.639999	319.000000	NaN	5.310000	53.0
3	2007-12-01 01:00:00	NaN	1.64	NaN	NaN	NaN	127.900002	476.700012	NaN	4.500000	105.3
4	2007-12-01 01:00:00	4.64	1.86	4.26	7.98	0.57	145.100006	573.900024	3.49	52.689999	106.5
...
225115	2007-03-01 00:00:00	0.30	0.45	1.00	0.30	0.26	8.690000	11.690000	1.00	42.209999	6.7
225116	2007-03-01 00:00:00	NaN	0.16	NaN	NaN	NaN	46.820000	51.480000	NaN	22.150000	5.7
225117	2007-03-01 00:00:00	0.24	NaN	0.20	NaN	0.09	51.259998	66.809998	NaN	18.540001	13.0
225118	2007-03-01 00:00:00	0.11	NaN	1.00	NaN	0.05	24.240000	36.930000	NaN	NaN	6.6
225119	2007-03-01 00:00:00	0.53	0.40	1.00	1.70	0.12	32.360001	47.860001	1.37	24.150000	10.2

225120 rows × 17 columns



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

Out[3]:

	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3	
0	2007-12-01 01:00:00	0.00	2.86	0.00	0.00	0.00	282.200012	1054.000000	0.00	4.030000	156.1
1	2007-12-01 01:00:00	0.00	1.82	0.00	0.00	0.00	86.419998	354.600006	0.00	3.260000	80.8
2	2007-12-01 01:00:00	0.00	1.47	0.00	0.00	0.00	94.639999	319.000000	0.00	5.310000	53.0
3	2007-12-01 01:00:00	0.00	1.64	0.00	0.00	0.00	127.900002	476.700012	0.00	4.500000	105.3
4	2007-12-01 01:00:00	4.64	1.86	4.26	7.98	0.57	145.100006	573.900024	3.49	52.689999	106.5
...
225115	2007-03-01 00:00:00	0.30	0.45	1.00	0.30	0.26	8.690000	11.690000	1.00	42.209999	6.7
225116	2007-03-01 00:00:00	0.00	0.16	0.00	0.00	0.00	46.820000	51.480000	0.00	22.150000	5.7
225117	2007-03-01 00:00:00	0.24	0.00	0.20	0.00	0.09	51.259998	66.809998	0.00	18.540001	13.0
225118	2007-03-01 00:00:00	0.11	0.00	1.00	0.00	0.05	24.240000	36.930000	0.00	0.000000	6.6
225119	2007-03-01 00:00:00	0.53	0.40	1.00	1.70	0.12	32.360001	47.860001	1.37	24.150000	10.2

225120 rows × 17 columns

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 225120 entries, 0 to 225119
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        225120 non-null object
1   BEN         225120 non-null float64
2   CO          225120 non-null float64
3   EBE         225120 non-null float64
4   MXY         225120 non-null float64
5   NMHC        225120 non-null float64
6   NO_2        225120 non-null float64
7   NOx         225120 non-null float64
8   OXY         225120 non-null float64
9   O_3         225120 non-null float64
10  PM10        225120 non-null float64
11  PM25        225120 non-null float64
12  PXY         225120 non-null float64
13  SO_2        225120 non-null float64
14  TCH         225120 non-null float64
15  TOL         225120 non-null float64
16  station     225120 non-null int64
dtypes: float64(15), int64(1), object(1)
memory usage: 29.2+ MB
```

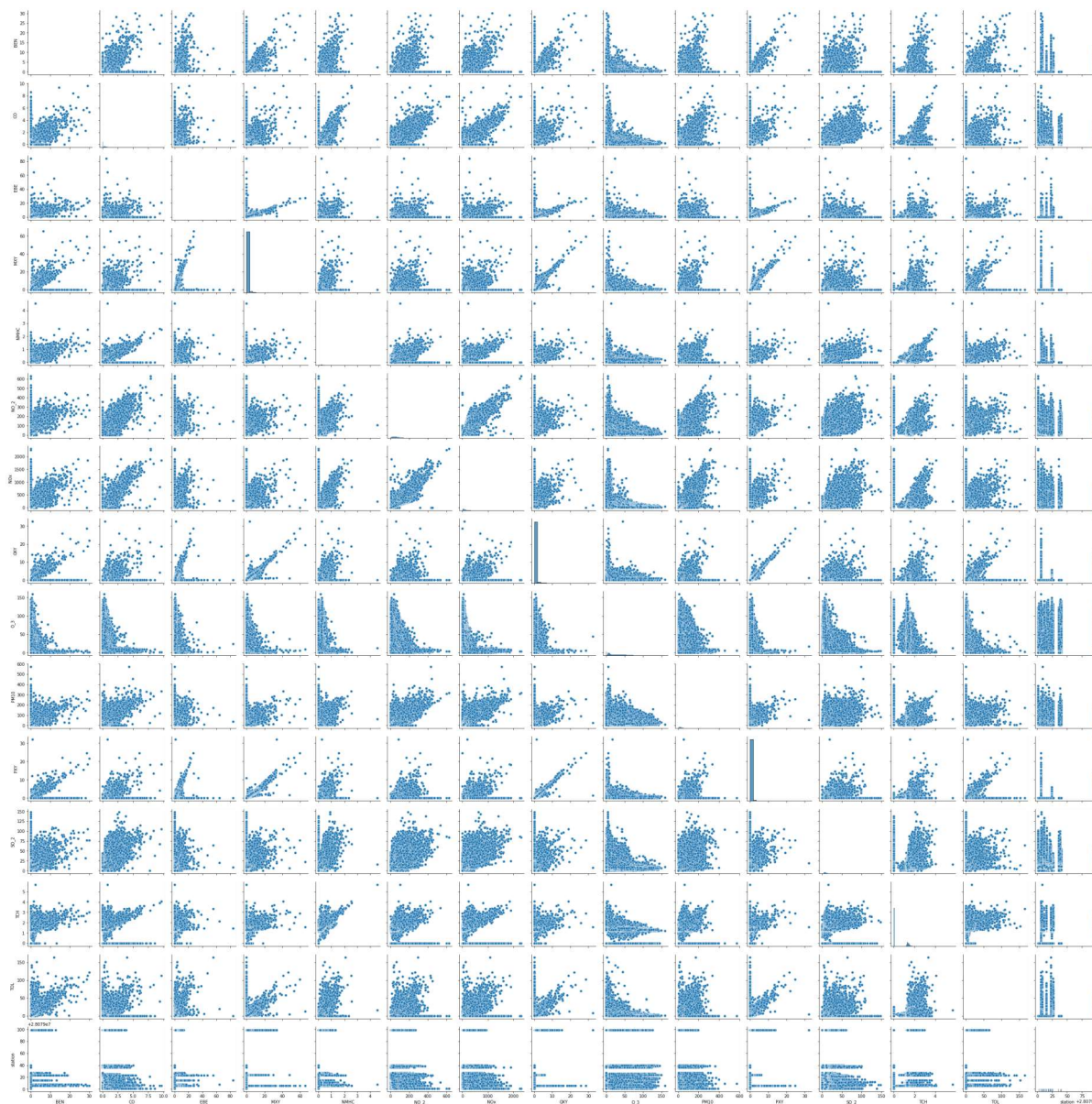
```
In [5]: df1.columns
```

```
Out[5]: Index(['date', 'BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
               'PM10', 'PM25', 'PXY', 'SO_2', 'TCH', 'TOL', 'station'],
              dtype='object')
```

```
In [6]: df2 = df1[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
                  'PM10', 'PXY', 'SO_2', 'TCH', 'TOL', 'station']]
```

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

```
Out[7]: <seaborn.axisgrid.PairGrid at 0x238032f7340>
```

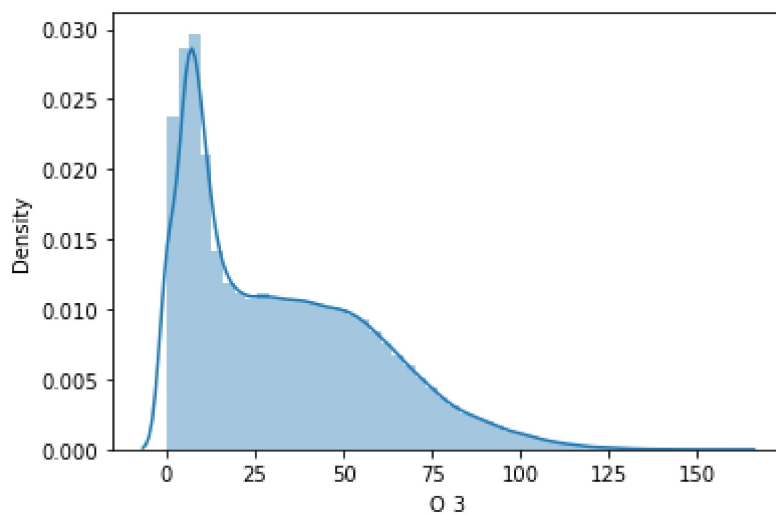


```
In [8]: sns.distplot(df2['O_3'])
```

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).

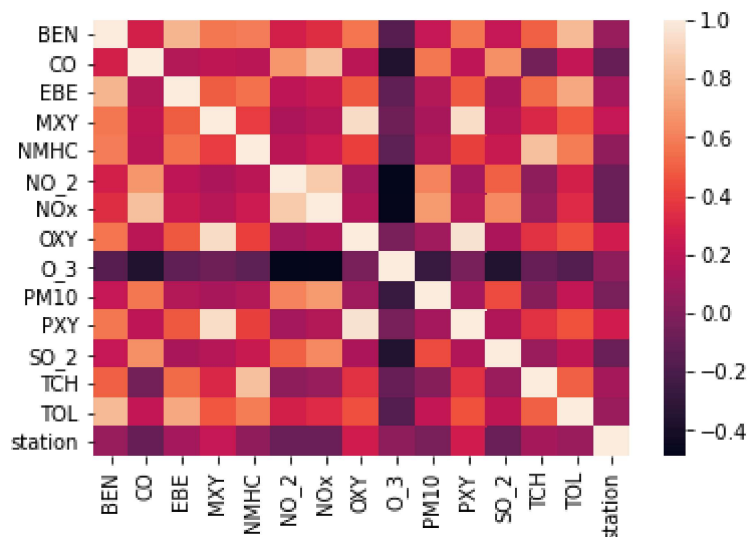
```
warnings.warn(msg, FutureWarning)
```

```
Out[8]: <AxesSubplot:xlabel='O_3', ylabel='Density'>
```



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

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



Linear Regression

```
In [10]: x = df2[['BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',  
               'PM10', 'SO_2', 'PXY', '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_)  
  
0.04001667203029302
```

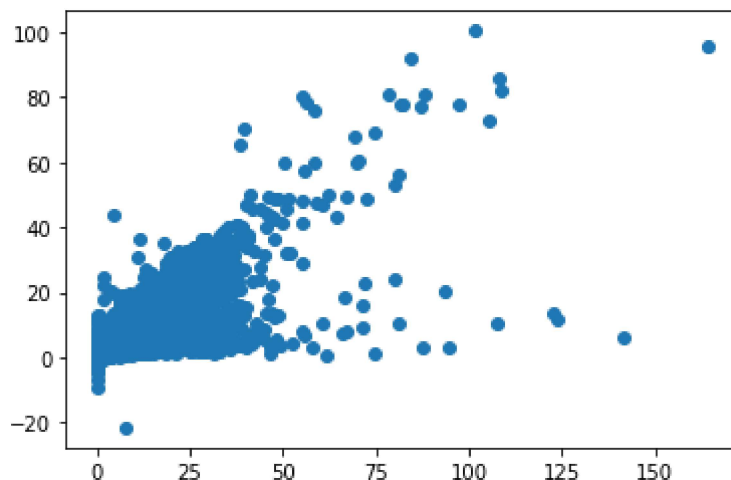
```
In [14]: coeff = pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])  
coeff
```

Out[14]:

	Co-efficient
BEN	2.774757
CO	-1.896856
EBE	0.914089
MXY	0.358328
NMHC	5.912057
NO_2	0.000512
NOx	0.006967
OXY	-0.865775
O_3	-0.001003
PM10	0.001280
SO_2	0.003489
PXY	0.088023
TCH	-0.382367

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

Out[15]: <matplotlib.collections.PathCollection at 0x2381fcdeaf0>



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

0.7029680841987095

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

Out[17]: 0.6997230429619776

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]: 0.6997214702742314

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

Out[20]: 0.7029800822222487

Lasso Regression

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

Out[21]: 0.11491515415927389

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

```
Out[22]: 0.09695612637041073
```

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_)
```

```
[ 9.96608222e-01 -0.00000000e+00  1.07555157e+00  2.56888091e-01  
 0.00000000e+00  0.00000000e+00  7.35919829e-03  0.00000000e+00  
 -3.43432023e-04  0.00000000e+00 -0.00000000e+00  0.00000000e+00  
 2.13947761e-01]
```

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

```
-0.2429278308475067
```

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

```
0.5897212219107215
```

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

```
0.5891546718974117
```

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]: (225120, 15)
```

```
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 [34]: observation = [[1.4,1.5,1.6,2.7,2.3,3.3,2.3,4.1,2.3,4.2,1.2,2.1,4.3,6,2.2]]
prediction = logs.predict(observation)
```

```
In [35]: print(prediction)

[28079099]
```

```
In [36]: logs.classes_
```

```
Out[36]: array([28079001, 28079003, 28079004, 28079006, 28079007, 28079008,
                28079009, 28079011, 28079012, 28079014, 28079015, 28079016,
                28079018, 28079019, 28079021, 28079022, 28079023, 28079024,
                28079025, 28079026, 28079027, 28079036, 28079038, 28079039,
                28079040, 28079099], dtype=int64)
```

```
In [37]: 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 [38]: print(logs.score(x_test,y_test))

0.03879412461502014
```

```
In [39]: print(logs.score(x_train,y_train))

0.038963346532642905
```

Conclusion

Ridge regression is bestfit model

Ridge regression is best fit model for dataset madrid_2001. The score of x_train,y_train is 0.6997214702742314 and x_test and y_test score is 0.7029800822222487.

