

In [1]: `# MADRID 2013`

In [1]: `import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression, LogisticRegression, Lasso, Ridge
from sklearn.model_selection import train_test_split`

In [2]: `df=pd.read_csv("C:/Users/user/Downloads/FP1_air/csvs_per_year/csvs_per_year/madrid_2013_14.csv")
df`

Out[2]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL	
0	2013-11-01 01:00:00	NaN	0.6	NaN	NaN	135.0	74.0	NaN	NaN	NaN	7.0	NaN	NaN	2
1	2013-11-01 01:00:00	1.5	0.5	1.3	NaN	71.0	83.0	2.0	23.0	16.0	12.0	NaN	8.3	2
2	2013-11-01 01:00:00	3.9	NaN	2.8	NaN	49.0	70.0	NaN	NaN	NaN	NaN	NaN	9.0	2
3	2013-11-01 01:00:00	NaN	0.5	NaN	NaN	82.0	87.0	3.0	NaN	NaN	NaN	NaN	NaN	2
4	2013-11-01 01:00:00	NaN	NaN	NaN	NaN	242.0	111.0	2.0	NaN	NaN	12.0	NaN	NaN	2
...
209875	2013-03-01 00:00:00	NaN	0.4	NaN	NaN	8.0	39.0	52.0	NaN	NaN	NaN	NaN	NaN	2
209876	2013-03-01 00:00:00	NaN	0.4	NaN	NaN	1.0	11.0	NaN	6.0	NaN	2.0	NaN	NaN	2
209877	2013-03-01 00:00:00	NaN	NaN	NaN	NaN	2.0	4.0	75.0	NaN	NaN	NaN	NaN	NaN	2
209878	2013-03-01 00:00:00	NaN	NaN	NaN	NaN	2.0	11.0	52.0	NaN	NaN	NaN	NaN	NaN	2
209879	2013-03-01 00:00:00	NaN	NaN	NaN	NaN	1.0	10.0	75.0	3.0	NaN	NaN	NaN	NaN	2

209880 rows × 14 columns



```
In [3]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 209880 entries, 0 to 209879
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   date        209880 non-null object
1   BEN         50462 non-null float64
2   CO          87018 non-null float64
3   EBE         50463 non-null float64
4   NMHC        25935 non-null float64
5   NO          209108 non-null float64
6   NO_2        209108 non-null float64
7   O_3         121858 non-null float64
8   PM10        104339 non-null float64
9   PM25        51980 non-null float64
10  SO_2        86970 non-null float64
11  TCH         25935 non-null float64
12  TOL         50317 non-null float64
13  station     209880 non-null int64
dtypes: float64(12), int64(1), object(1)
memory usage: 22.4+ MB
```

```
In [4]: df1=df.dropna()
df1
```

```
Out[4]:
```

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL	st
17286	2013-08-01 01:00:00	0.4	0.2	0.8	0.28	1.0	24.0	79.0	35.0	8.0	3.0	1.49	1.3	2807
17310	2013-08-01 02:00:00	0.5	0.2	0.9	0.28	1.0	16.0	93.0	60.0	18.0	3.0	1.61	4.0	2807
17334	2013-08-01 03:00:00	0.5	0.2	1.1	0.29	1.0	14.0	90.0	38.0	12.0	3.0	1.71	2.8	2807
17358	2013-08-01 04:00:00	0.6	0.2	1.2	0.26	1.0	12.0	84.0	30.0	8.0	3.0	1.44	2.8	2807
17382	2013-08-01 05:00:00	0.3	0.2	0.8	0.25	1.0	15.0	72.0	25.0	7.0	3.0	1.40	1.7	2807
...
209622	2013-02-28 14:00:00	1.1	0.3	0.3	0.27	3.0	17.0	64.0	5.0	5.0	2.0	1.41	0.9	2807
209646	2013-02-28 15:00:00	1.3	0.4	0.3	0.27	2.0	16.0	66.0	6.0	5.0	1.0	1.40	0.9	2807
209670	2013-02-28 16:00:00	1.1	0.3	0.3	0.27	1.0	17.0	65.0	5.0	4.0	1.0	1.40	0.7	2807
209694	2013-02-28 17:00:00	1.0	0.3	0.4	0.27	1.0	18.0	64.0	5.0	5.0	1.0	1.39	0.7	2807
209718	2013-02-28 18:00:00	1.0	0.3	0.4	0.27	1.0	22.0	62.0	6.0	6.0	1.0	1.39	0.7	2807

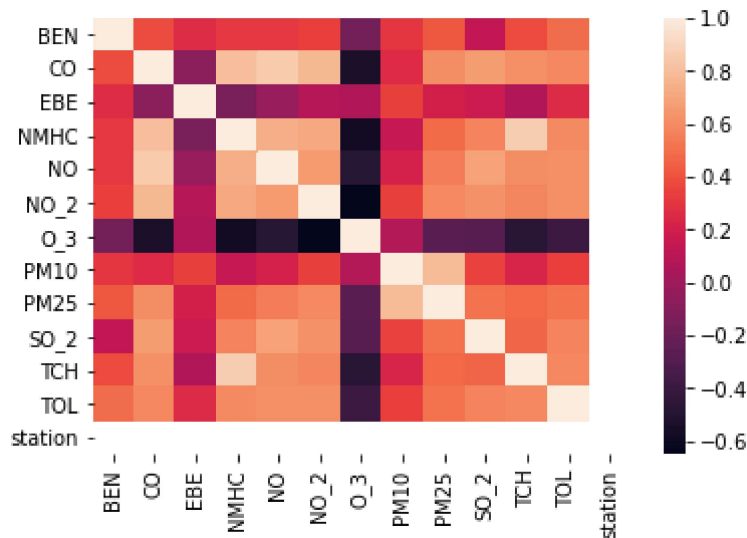
7315 rows × 14 columns



```
In [5]: df1=df1.drop(["date"],axis=1)
```

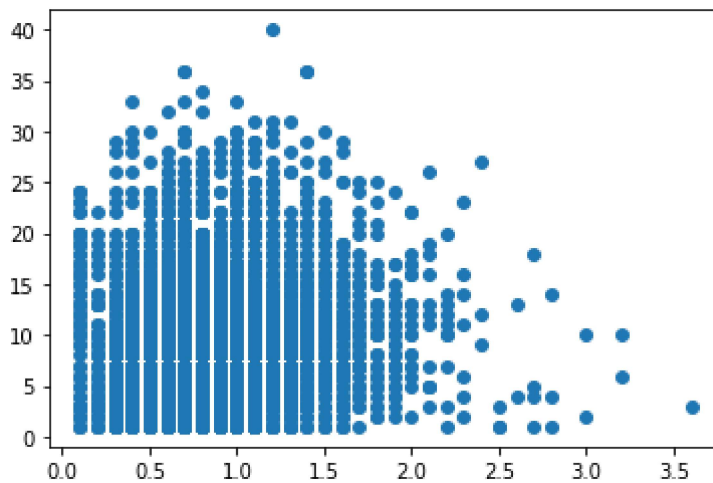
```
In [6]: sns.heatmap(df1.corr())
```

```
Out[6]: <AxesSubplot:>
```



```
In [7]: plt.plot(df1["EBE"],df1["PM25"],"o")
```

```
Out[7]: [<matplotlib.lines.Line2D at 0x25a0674f250>]
```



```
In [8]: x=df1.drop(["EBE"],axis=1)
y=df1["EBE"]
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

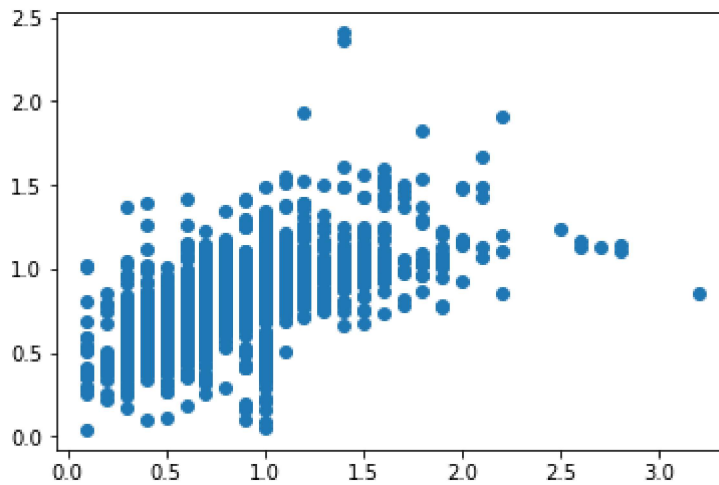
Linear

```
In [9]: li=LinearRegression()
li.fit(x_train,y_train)
```

```
Out[9]: LinearRegression()
```

```
In [10]: prediction=li.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[10]: <matplotlib.collections.PathCollection at 0x25a06816790>



```
In [11]: lis=li.score(x_test,y_test)
```

```
In [12]: df1["TCH"].value_counts()
```

Out[12]:

1.32	888
1.33	843
1.34	729
1.31	719
1.35	556
...	
1.23	1
2.09	1
1.84	1
2.25	1
2.29	1

Name: TCH, Length: 114, dtype: int64

```
In [13]: df1.loc[df1["TCH"]<1.40,"TCH"]=1
df1.loc[df1["TCH"]>1.40,"TCH"]=2
df1["TCH"].value_counts()
```

Out[13]:

1.0	5718
2.0	1597

Name: TCH, dtype: int64

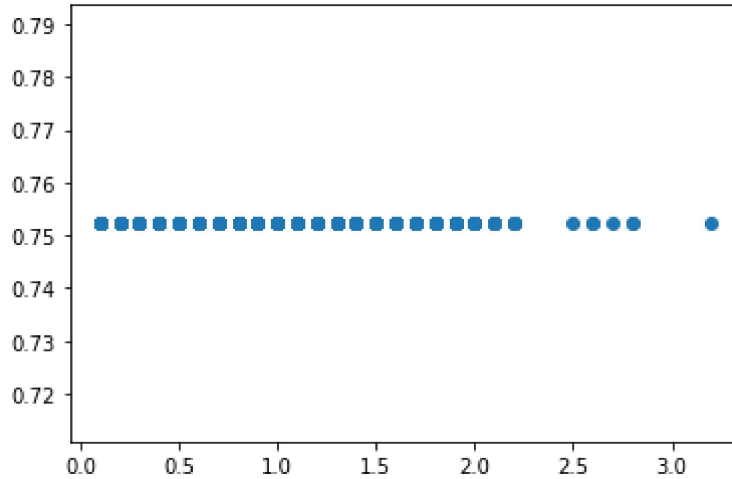
Lasso

```
In [14]: la=Lasso(alpha=5)
la.fit(x_train,y_train)
```

Out[14]: Lasso(alpha=5)

```
In [15]: prediction1=la.predict(x_test)
plt.scatter(y_test,prediction1)
```

Out[15]: <matplotlib.collections.PathCollection at 0x25a0687ed90>



```
In [16]: las=la.score(x_test,y_test)
```

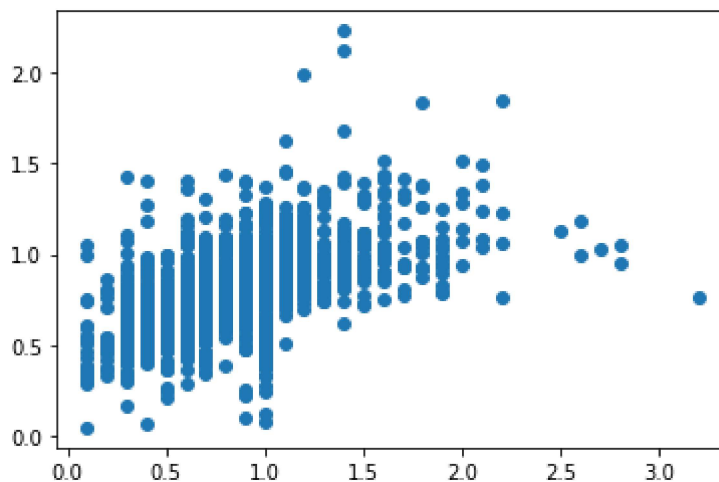
Ridge

```
In [17]: rr=Ridge(alpha=1)
rr.fit(x_train,y_train)
```

Out[17]: Ridge(alpha=1)

```
In [18]: prediction2=rr.predict(x_test)
plt.scatter(y_test,prediction2)
```

Out[18]: <matplotlib.collections.PathCollection at 0x25a066009d0>



```
In [19]: rrs=rr.score(x_test,y_test)
```

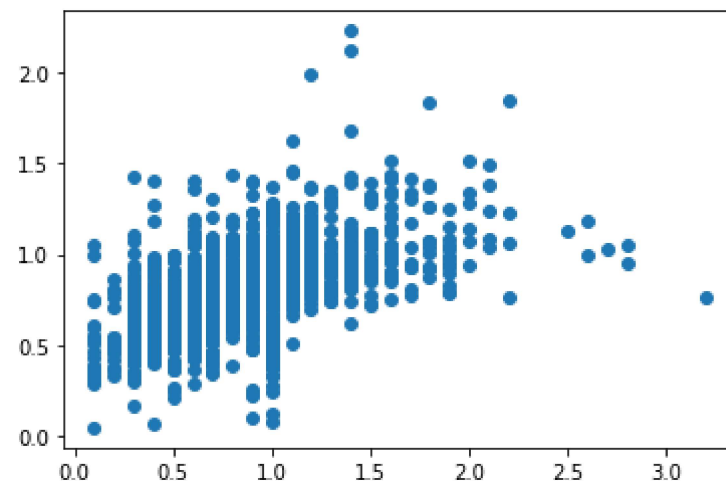
ElasticNet

```
In [20]: en=ElasticNet()
en.fit(x_train,y_train)
```

Out[20]: ElasticNet()

```
In [21]: prediction2=rr.predict(x_test)
plt.scatter(y_test,prediction2)
```

Out[21]: <matplotlib.collections.PathCollection at 0x25a07126cd0>



```
In [22]: ens=en.score(x_test,y_test)
```

```
In [23]: print(rr.score(x_test,y_test))
rr.score(x_train,y_train)
```

0.38794864124666883

Out[23]: 0.391803939842712

Logistic

```
In [24]: g={"TCH":{1.0:"Low",2.0:"High"}}
df1=df1.replace(g)
df1["TCH"].value_counts()
```

Out[24]: Low 5718
High 1597
Name: TCH, dtype: int64

```
In [25]: x=df1.drop(["TCH"],axis=1)
y=df1["TCH"]
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [26]: lo=LogisticRegression()
lo.fit(x_train,y_train)
```

Out[26]: LogisticRegression()

```
In [27]: prediction3=lo.predict(x_test)
plt.scatter(y_test,prediction3)
```

Out[27]: <matplotlib.collections.PathCollection at 0x25a07187280>



```
In [28]: los=lo.score(x_test,y_test)
```

Random Forest

```
In [29]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
```

```
In [30]: g1={"TCH":{"Low":1.0,"High":2.0}}
df1=df1.replace(g1)
```

```
In [31]: x=df1.drop(["TCH"],axis=1)
y=df1["TCH"]
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [32]: rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

Out[32]: RandomForestClassifier()


```
In [33]: parameter={
    'max_depth':[1,2,4,5,6],
    'min_samples_leaf':[5,10,15,20,25],
    'n_estimators':[10,20,30,40,50]
}
```

```
In [34]: grid_search=GridSearchCV(estimator=rfc,param_grid=parameter,cv=2,scoring="accuracy")
grid_search.fit(x_train,y_train)
```

```
Out[34]: GridSearchCV(cv=2, estimator=RandomForestClassifier(),
    param_grid={'max_depth': [1, 2, 4, 5, 6],
    'min_samples_leaf': [5, 10, 15, 20, 25],
    'n_estimators': [10, 20, 30, 40, 50]},
    scoring='accuracy')
```

```
In [35]: rfcs=grid_search.best_score_
```

```
In [36]: rfc_best=grid_search.best_estimator_
```

```
In [37]: from sklearn.tree import plot_tree

plt.figure(figsize=(80,40))
plot_tree(rfc_best.estimators_[5],feature_names=x.columns,class_names=['Yes','No'])
```

```
Out[37]: [Text(2305.4210526315787, 2019.0857142857144, 'TOL <= 1.75\ngini = 0.337\nsamples = 3228\nvalue = [4020, 1100]\nclass = Yes'),
    Text(1204.1052631578948, 1708.457142857143, 'O_3 <= 26.5\ngini = 0.195\nsamples = 2605\nvalue = [3677, 453]\nclass = Yes'),
    Text(469.89473684210526, 1397.8285714285716, 'NMHC <= 0.265\ngini = 0.436\nsamples = 262\nvalue = [124, 262]\nclass = No'),
    Text(234.94736842105263, 1087.2, 'PM10 <= 19.5\ngini = 0.266\nsamples = 59\nvalue = [80, 15]\nclass = Yes'),
    Text(176.21052631578948, 776.5714285714287, 'TOL <= 1.55\ngini = 0.217\nsamples = 54\nvalue = [78, 11]\nclass = Yes'),
    Text(117.47368421052632, 465.9428571428573, 'EBE <= 1.15\ngini = 0.147\nsamples = 46\nvalue = [69, 6]\nclass = Yes'),
    Text(58.73684210526316, 155.3142857142857, 'gini = 0.215\nsamples = 31\nvalue = [43, 6]\nclass = Yes'),
    Text(176.21052631578948, 155.3142857142857, 'gini = 0.0\nsamples = 15\nvalue = [26, 0]\nclass = Yes'),
    Text(234.94736842105263, 465.9428571428573, 'gini = 0.459\nsamples = 8\nvalue = [9, 5]\nclass = Yes'),
    Text(293.6842105263158, 776.5714285714287, 'gini = 0.444\nsamples = 5\nvalue = [1, 4]\nclass = No')]
```

```
In [38]: print("Linear:", lis)
print("Lasso:", las)
print("Ridge:", rrs)
print("ElasticNet:", ens)
print("Logistic:", los)
print("Random Forest:", rfcs)
```

```
Linear: 0.4114978071131482
Lasso: -6.234443005292967e-05
Ridge: 0.38794864124666883
ElasticNet: 0.09906134120330623
Logistic: 0.7831435079726652
Random Forest: 0.9505859375000001
```

Best Model is Random Forest

MADRID 2014

```
In [39]: df2=pd.read_csv("C:/Users/user/Downloads/FP1_air/csvs_per_year/csvs_per_year/madrid_data(2013_14).csv")
df2
```

```
Out[39]:
```

	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 [40]: df2.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         46703 non-null float64
 2   CO          87023 non-null float64
 3   EBE         46722 non-null float64
 4   NMHC        25021 non-null float64
 5   NO          209154 non-null float64
 6   NO_2        209154 non-null float64
 7   O_3         121681 non-null float64
 8   PM10        104311 non-null float64
 9   PM25        51954 non-null float64
10   SO_2        87141 non-null float64
11   TCH         25021 non-null float64
12   TOL         46570 non-null float64
13   station     210024 non-null int64
dtypes: float64(12), int64(1), object(1)
memory usage: 22.4+ MB
```

```
In [41]: df3=df2.dropna()  
df3
```

```
Out[41]:
```

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TOL	ε
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
6	2014-06-01 01:00:00	0.1	0.2	0.1	0.23	1.0	5.0	80.0	4.0	3.0	2.0	1.21	0.1	280
25	2014-06-01 02:00:00	0.2	0.2	0.1	0.11	4.0	21.0	63.0	9.0	6.0	5.0	1.36	0.8	280
30	2014-06-01 02:00:00	0.2	0.2	0.1	0.23	1.0	4.0	88.0	7.0	5.0	2.0	1.21	0.1	280
49	2014-06-01 03:00:00	0.1	0.2	0.1	0.11	4.0	18.0	66.0	9.0	7.0	6.0	1.36	0.9	280
...
209958	2014-08-31 22:00:00	0.2	0.2	0.1	0.22	1.0	28.0	96.0	61.0	15.0	3.0	1.28	0.1	280
209977	2014-08-31 23:00:00	1.1	0.7	0.7	0.19	36.0	118.0	23.0	60.0	25.0	9.0	1.27	6.5	280
209982	2014-08-31 23:00:00	0.2	0.2	0.1	0.21	1.0	17.0	90.0	28.0	14.0	3.0	1.27	0.2	280
210001	2014-09-01 00:00:00	0.6	0.4	0.4	0.12	6.0	63.0	41.0	26.0	15.0	8.0	1.19	4.1	280
210006	2014-09-01 00:00:00	0.2	0.2	0.1	0.23	1.0	30.0	69.0	18.0	13.0	3.0	1.30	0.1	280

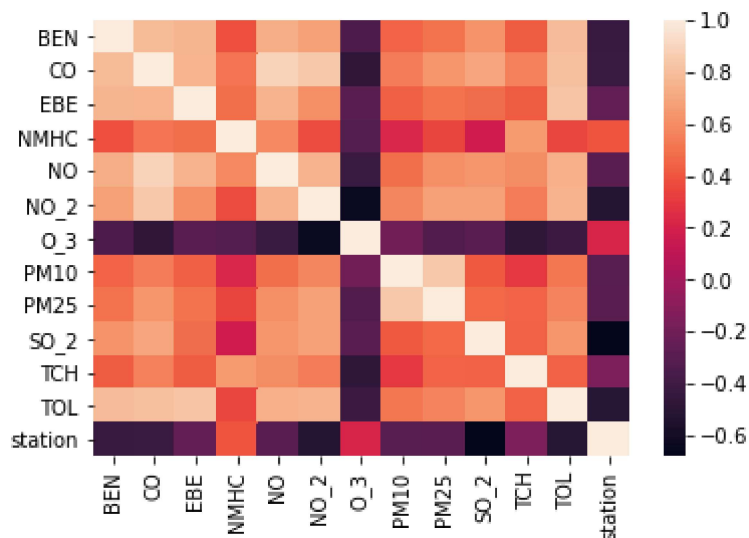
13946 rows × 14 columns



```
In [42]: df3=df3.drop(["date"],axis=1)
```

```
In [43]: sns.heatmap(df3.corr())
```

```
Out[43]: <AxesSubplot:>
```



```
In [44]: x=df3.drop(["TCH"],axis=1)
y=df3["TCH"]
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

Linear

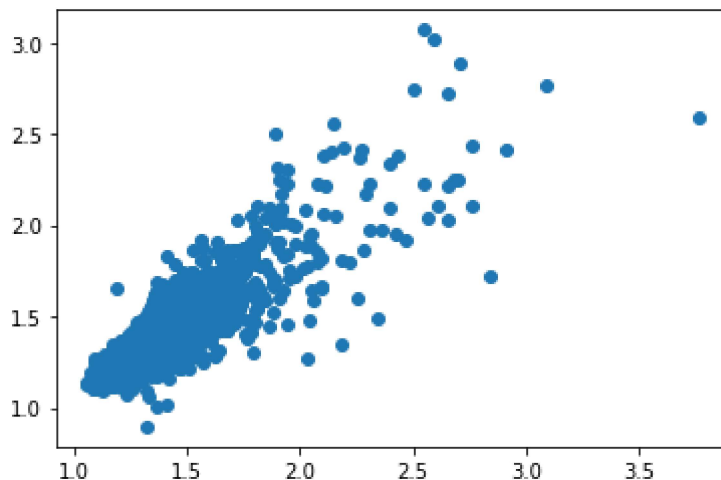
```
In [45]: li=LinearRegression()
li.fit(x_train,y_train)
```

```
Out[45]: LinearRegression()
```

```
In [ ]:
```

```
In [46]: prediction=li.predict(x_test)
plt.scatter(y_test,prediction)
```

```
Out[46]: <matplotlib.collections.PathCollection at 0x25a0b165d60>
```



```
In [47]: lis=li.score(x_test,y_test)
```

```
In [48]: df3["TCH"].value_counts()
```

```
Out[48]: 1.37    601
         1.36    598
         1.34    529
         1.35    528
         1.38    515
         ...
         2.50     1
         2.86     1
         2.70     1
         3.04     1
         4.37     1
         Name: TCH, Length: 184, dtype: int64
```

```
In [49]: df3.loc[df3["TCH"]<1.40,"TCH"]=1
         df3.loc[df3["TCH"]>1.40,"TCH"]=2
         df3["TCH"].value_counts()
```

```
Out[49]: 1.0    9997
         2.0    3949
         Name: TCH, dtype: int64
```

```
In [ ]:
```

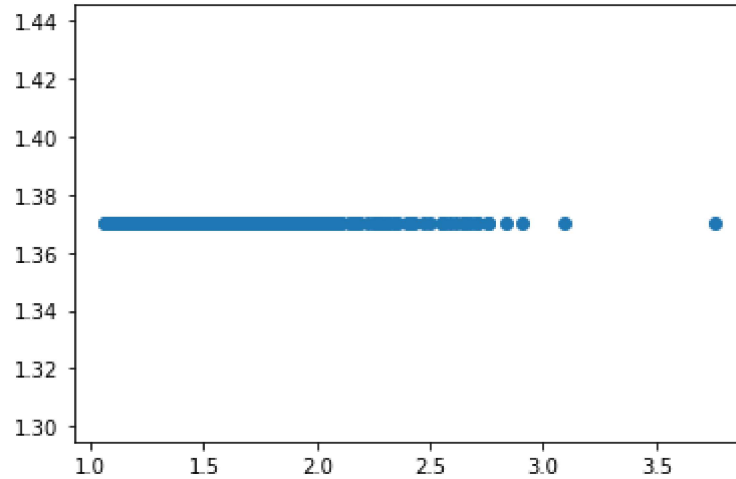
Lasso

```
In [50]: la=Lasso(alpha=5)
         la.fit(x_train,y_train)
```

```
Out[50]: Lasso(alpha=5)
```

```
In [51]: prediction1=la.predict(x_test)
plt.scatter(y_test,prediction1)
```

Out[51]: <matplotlib.collections.PathCollection at 0x25a0b1bca90>



```
In [52]: las=la.score(x_test,y_test)
```

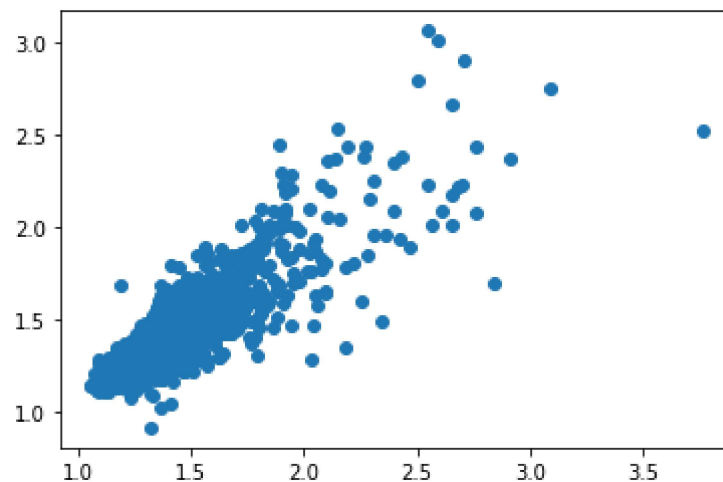
Ridge

```
In [53]: rr=Ridge(alpha=1)
rr.fit(x_train,y_train)
```

Out[53]: Ridge(alpha=1)

```
In [54]: prediction2=rr.predict(x_test)
plt.scatter(y_test,prediction2)
```

Out[54]: <matplotlib.collections.PathCollection at 0x25a0b21f100>



```
In [55]: rrs=rr.score(x_test,y_test)
```

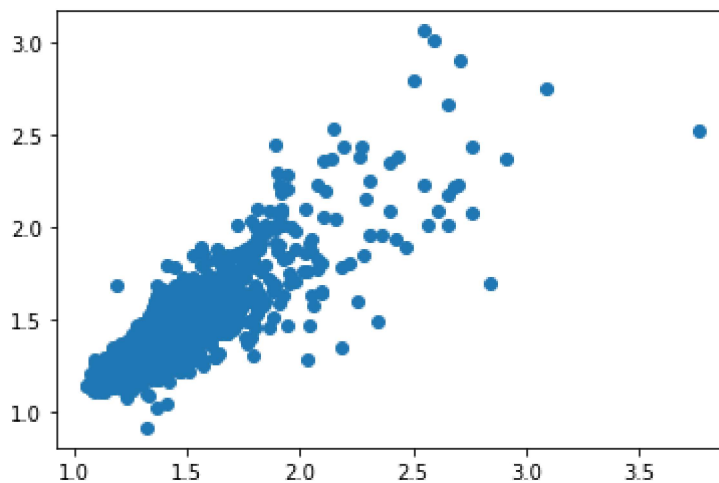

ElasticNet

```
In [56]: en=ElasticNet()  
en.fit(x_train,y_train)
```

Out[56]: ElasticNet()

```
In [57]: prediction2=rr.predict(x_test)  
plt.scatter(y_test,prediction2)
```

Out[57]: <matplotlib.collections.PathCollection at 0x25a0b26e6a0>



```
In [58]: ens=en.score(x_test,y_test)
```

```
In [59]: print(rr.score(x_test,y_test))  
rr.score(x_train,y_train)
```

0.7037488127556343

Out[59]: 0.7081013643502528

Logistic

```
In [60]: g={"TCH":{1.0:"Low",2.0:"High"}}  
df3=df3.replace(g)  
df3["TCH"].value_counts()
```

Out[60]: Low 9997
High 3949
Name: TCH, dtype: int64

```
In [61]: x=df3.drop(["TCH"],axis=1)
y=df3["TCH"]
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [62]: lo=LogisticRegression()
lo.fit(x_train,y_train)
```

Out[62]: LogisticRegression()

```
In [63]: prediction3=lo.predict(x_test)
plt.scatter(y_test,prediction3)
```

Out[63]: <matplotlib.collections.PathCollection at 0x25a0b2a06a0>



```
In [64]: los=lo.score(x_test,y_test)
```

Random Forest

```
In [65]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
```

```
In [66]: g1={"TCH":{"Low":1.0,"High":2.0}}
df3=df3.replace(g1)
```

```
In [67]: x=df3.drop(["TCH"],axis=1)
y=df3["TCH"]
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [68]: rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

Out[68]: RandomForestClassifier()

```
In [69]: parameter={
    'max_depth':[1,2,4,5,6],
    'min_samples_leaf':[5,10,15,20,25],
    'n_estimators':[10,20,30,40,50]
}
```

```
In [70]: grid_search=GridSearchCV(estimator=rfc,param_grid=parameter,cv=2,scoring="accuracy")
grid_search.fit(x_train,y_train)
```

```
Out[70]: GridSearchCV(cv=2, estimator=RandomForestClassifier(),
    param_grid={'max_depth': [1, 2, 4, 5, 6],
    'min_samples_leaf': [5, 10, 15, 20, 25],
    'n_estimators': [10, 20, 30, 40, 50]},
    scoring='accuracy')
```

```
In [71]: rfcs=grid_search.best_score_
```

```
In [72]: rfc_best=grid_search.best_estimator_
```

```
In [73]: from sklearn.tree import plot_tree
```

```
plt.figure(figsize=(80,40))
plot_tree(rfc_best.estimators_[5],feature_names=x.columns,class_names=['Yes','No'])
```

```
Out[73]: [Text(2730.6382978723404, 2019.0857142857144, 'TOL <= 3.65\ngini = 0.405\nsamples = 6192\nvalue = [7009, 2753]\nclass = Yes'),
    Text(1519.659574468085, 1708.457142857143, 'CO <= 0.35\ngini = 0.321\nsamples = 4987\nvalue = [6278, 1581]\nclass = Yes'),
    Text(759.8297872340426, 1397.8285714285716, 'O_3 <= 24.5\ngini = 0.237\nsamples = 4172\nvalue = [5672, 902]\nclass = Yes'),
    Text(379.9148936170213, 1087.2, 'PM25 <= 5.5\ngini = 0.472\nsamples = 398\nvalue = [233, 379]\nclass = No'),
    Text(189.95744680851064, 776.5714285714287, 'station <= 28079016.0\ngini = 0.426\nsamples = 31\nvalue = [36, 16]\nclass = Yes'),
    Text(94.97872340425532, 465.9428571428573, 'O_3 <= 19.0\ngini = 0.48\nsamples = 14\nvalue = [8, 12]\nclass = No'),
    Text(47.48936170212766, 155.3142857142857, 'gini = 0.0\nsamples = 8\nvalue = [0, 12]\nclass = No'),
    Text(142.46808510638297, 155.3142857142857, 'gini = 0.0\nsamples = 6\nvalue = [8, 0]\nclass = Yes'),
    Text(284.93617021276594, 465.9428571428573, 'NO_2 <= 21.0\ngini = 0.219\nsamples = 17\nvalue = [28, 4]\nclass = Yes'),
    Text(237.4468085106383, 155.3142857142857, 'gini = 0.0\nsamples = 5\nvalue = [12, 0]\nclass = Yes')]
```

```
In [74]: print("Linear:",lis)
          print("Lasso:",las)
          print("Ridge:",rrs)
          print("ElasticNet:",ens)
          print("Logistic:",los)
          print("Random Forest:",rfcs)
```

Linear: 0.6984961782299648
Lasso: -0.00010435529624808204
Ridge: 0.7037488127556343
ElasticNet: 0.45599972417723134
Logistic: 0.7194072657743786
Random Forest: 0.8882401147305881

Best model is Random Forest

In []: