

```
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
from sklearn.model_selection import train_test_split
```

```
In [2]: df=pd.read_csv("/Users/bob/Downloads/FP1_air/csvs_per_year/csvs_per_
df
```

Out[2]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	TI
0	2015-10-01 01:00:00	NaN	0.8	NaN	NaN	90.0	82.0	NaN	NaN	NaN	10.0	NaN	N
1	2015-10-01 01:00:00	2.0	0.8	1.6	0.33	40.0	95.0	4.0	37.0	24.0	12.0	1.83	8
2	2015-10-01 01:00:00	3.1	NaN	1.8	NaN	29.0	97.0	NaN	NaN	NaN	NaN	NaN	5
3	2015-10-01 01:00:00	NaN	0.6	NaN	NaN	30.0	103.0	2.0	NaN	NaN	NaN	NaN	N
4	2015-10-01 01:00:00	NaN	NaN	NaN	NaN	95.0	96.0	2.0	NaN	NaN	9.0	NaN	N
...
210091	2015-08-01 00:00:00	NaN	0.2	NaN	NaN	11.0	33.0	53.0	NaN	NaN	NaN	NaN	N
210092	2015-08-01 00:00:00	NaN	0.2	NaN	NaN	1.0	5.0	NaN	26.0	NaN	10.0	NaN	N
210093	2015-08-01 00:00:00	NaN	NaN	NaN	NaN	1.0	7.0	74.0	NaN	NaN	NaN	NaN	N
210094	2015-08-01 00:00:00	NaN	NaN	NaN	NaN	3.0	7.0	65.0	NaN	NaN	NaN	NaN	N
210095	2015-08-01 00:00:00	NaN	NaN	NaN	NaN	1.0	9.0	54.0	29.0	NaN	NaN	NaN	N

210096 rows × 14 columns

In [3]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 210096 entries, 0 to 210095
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        210096 non-null  object
1   BEN         51039 non-null   float64
2   CO          86827 non-null   float64
3   EBE         50962 non-null   float64
4   NMHC        25756 non-null   float64
5   NO          208805 non-null   float64
6   NO_2        208805 non-null   float64
7   O_3         121574 non-null   float64
8   PM10        102745 non-null   float64
9   PM25        48798 non-null   float64
10  SO_2        86898 non-null   float64
11  TCH         25756 non-null   float64
12  TOL         50626 non-null   float64
13  station     210096 non-null   int64
dtypes: float64(12), int64(1), object(1)
memory usage: 22.4+ MB
```

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

Out[4]:

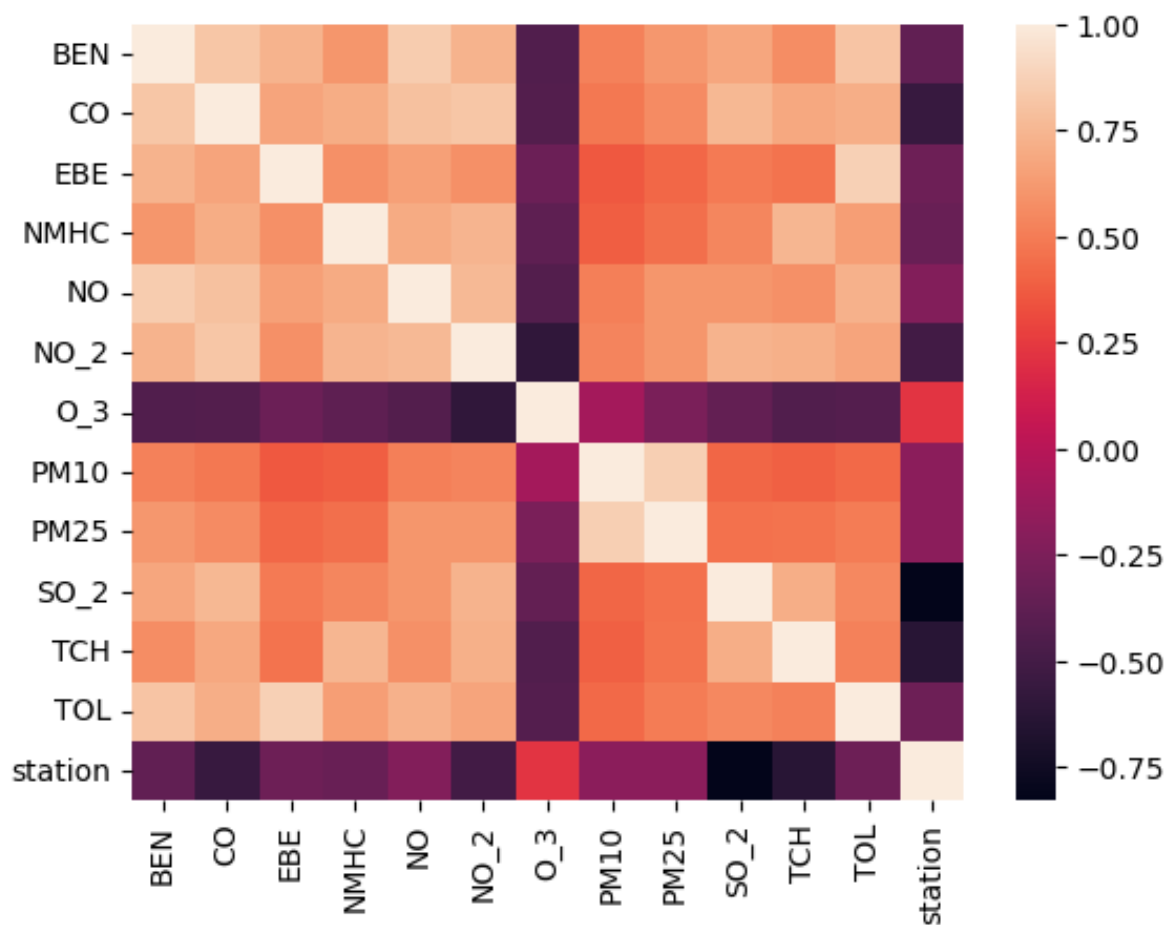
	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	T
1	2015-10-01 01:00:00	2.0	0.8	1.6	0.33	40.0	95.0	4.0	37.0	24.0	12.0	1.83	
6	2015-10-01 01:00:00	0.5	0.3	0.3	0.12	6.0	83.0	1.0	19.0	12.0	3.0	1.29	
25	2015-10-01 02:00:00	1.6	0.7	1.3	0.38	81.0	105.0	4.0	36.0	19.0	13.0	1.93	
30	2015-10-01 02:00:00	0.4	0.3	0.3	0.11	5.0	72.0	2.0	16.0	10.0	2.0	1.27	
49	2015-10-01 03:00:00	2.2	0.8	1.8	0.41	111.0	104.0	4.0	35.0	20.0	14.0	2.05	1
...
210030	2015-07-31 22:00:00	0.1	0.1	0.1	0.06	1.0	10.0	69.0	10.0	3.0	2.0	1.18	
210049	2015-07-31 23:00:00	0.4	0.3	0.1	0.12	3.0	28.0	56.0	15.0	7.0	12.0	1.45	
210054	2015-07-31 23:00:00	0.1	0.1	0.1	0.06	1.0	10.0	63.0	5.0	1.0	2.0	1.18	
210073	2015-08-01 00:00:00	0.1	0.3	0.1	0.11	2.0	23.0	59.0	5.0	2.0	11.0	1.44	
210078	2015-08-01 00:00:00	0.1	0.1	0.1	0.06	1.0	8.0	65.0	7.0	1.0	2.0	1.18	

16026 rows × 14 columns

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

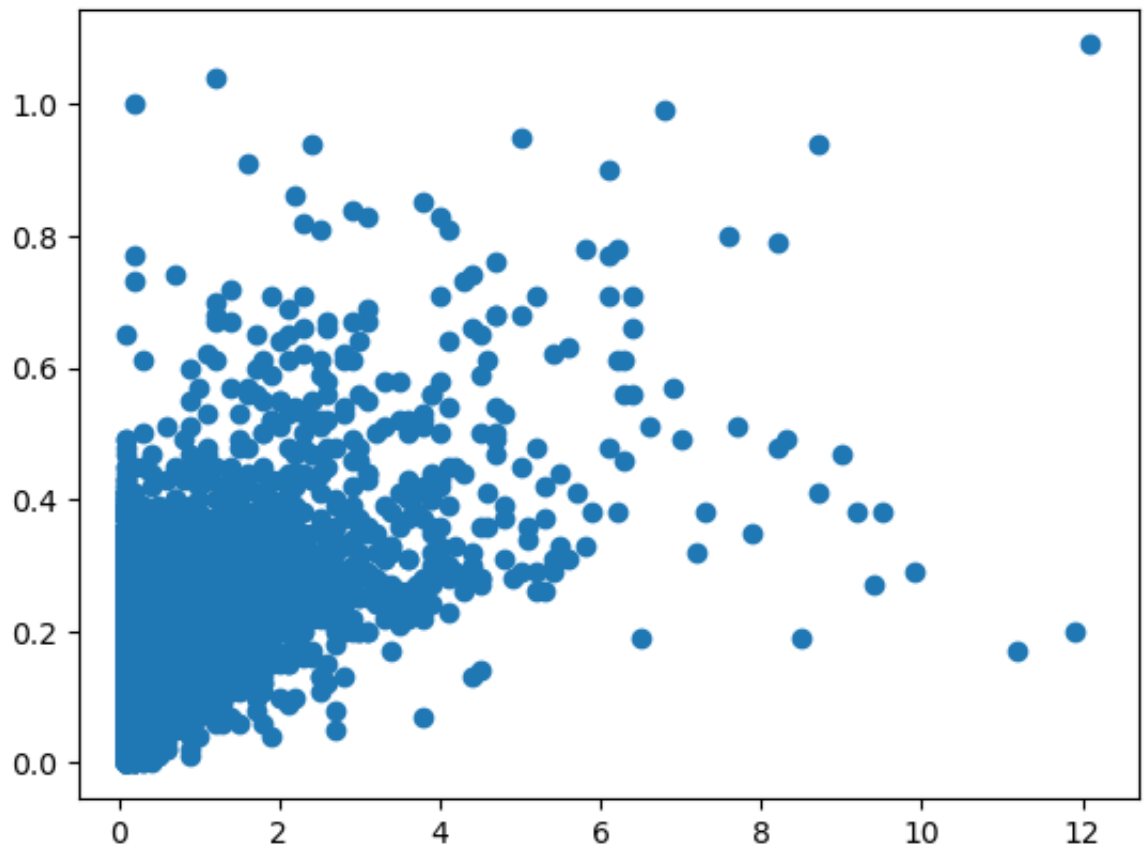
```
In [6]: sns.heatmap(df1.corr())
```

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



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

```
Out[7]: [ <matplotlib.lines.Line2D at 0x7fa5001f0880>]
```



```
In [8]: data=df[["EBE","NMHC"]]
```

```
In [9]: 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)
```

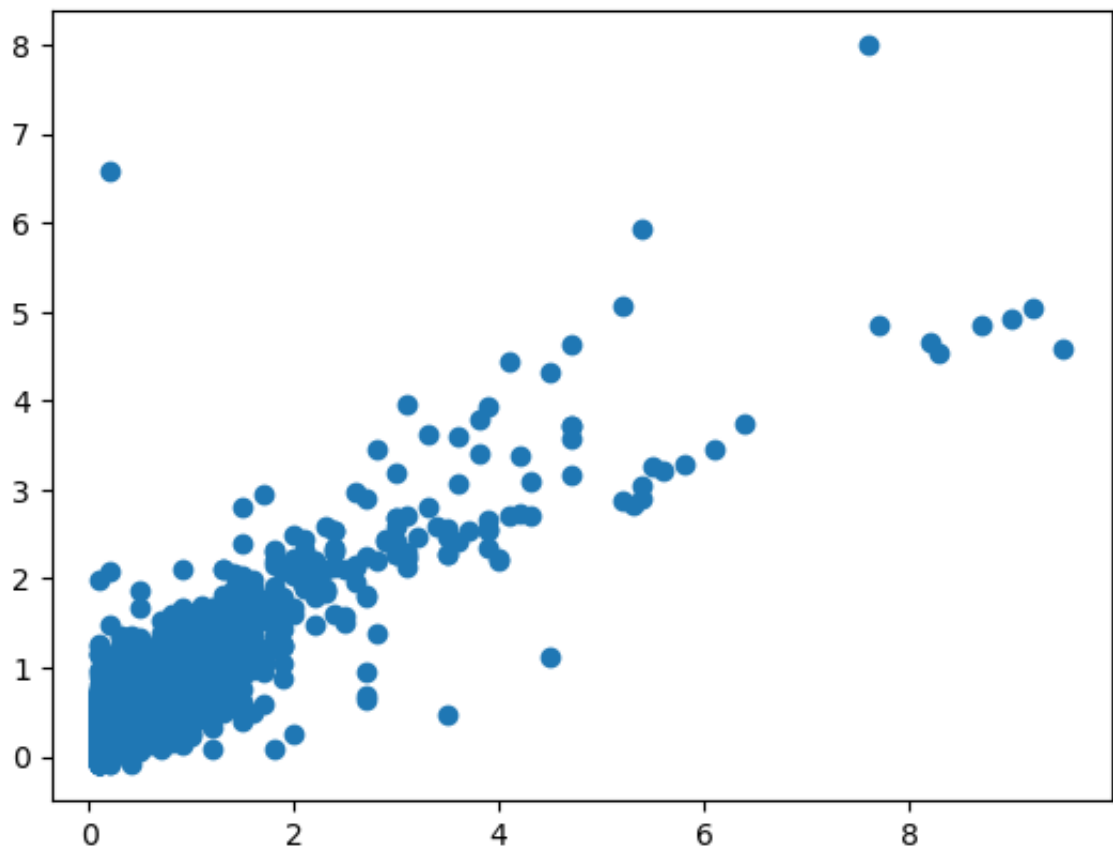
Linear

```
In [10]: li=LinearRegression()
li.fit(x_train,y_train)
```

```
Out[10]: ▼ LinearRegression
LinearRegression()
```

```
In [11]: prediction=li.predict(x_test)
plt.scatter(y_test,prediction)
```

```
Out[11]: <matplotlib.collections.PathCollection at 0x7fa50018b8b0>
```



```
In [12]: lis=li.score(x_test,y_test)
```

```
In [13]: df1["TCH"].value_counts()
```

```
Out[13]: 1.20    905
         1.19    873
         1.21    793
         1.22    638
         1.18    465
         ...
         2.57     1
         2.46     1
         2.93     1
         3.00     1
         3.02     1
         Name: TCH, Length: 184, dtype: int64
```

```
In [14]: df1.loc[df1["TCH"]<1.40,"TCH"]=1  
df1.loc[df1["TCH"]>1.40,"TCH"]=2  
df1["TCH"].value_counts()
```

```
Out[14]: 2.0    8290  
1.0    7736  
Name: TCH, dtype: int64
```

Lasso

```
In [15]: la=Lasso(alpha=5)  
la.fit(x_train,y_train)
```

```
Out[15]: 

▼



Lasso

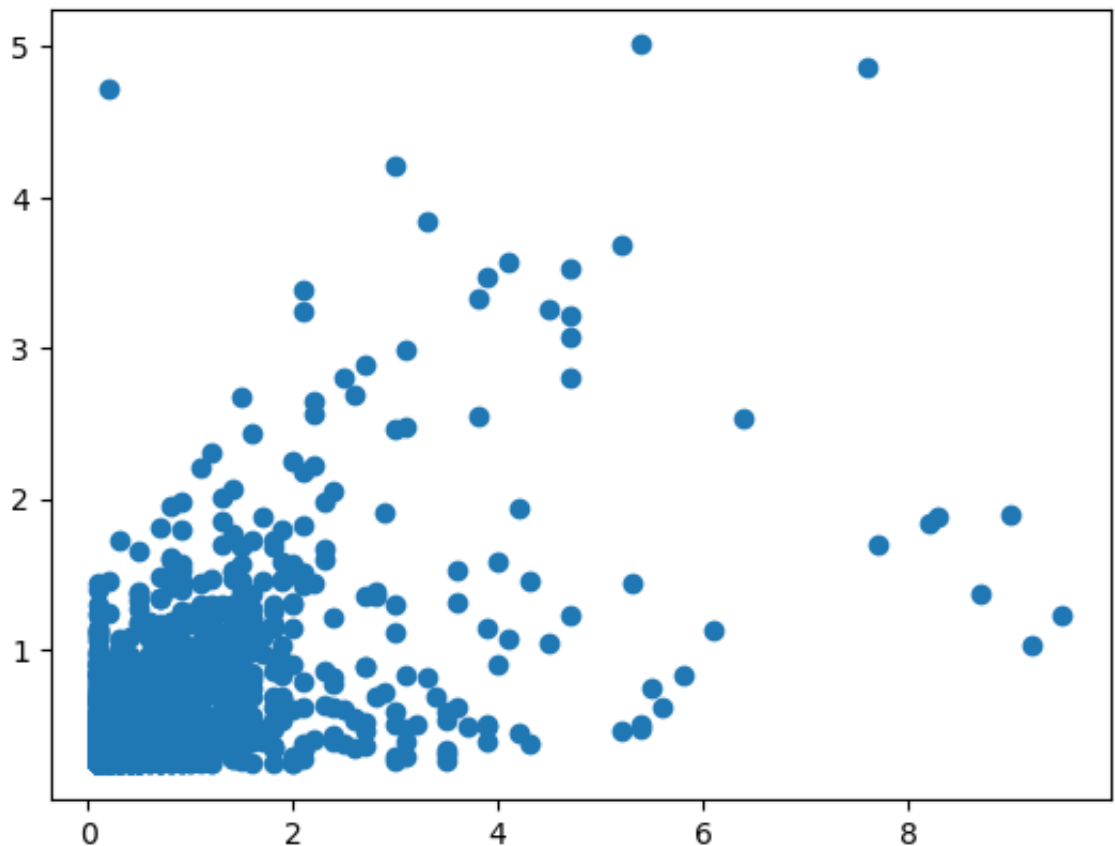


Lasso(alpha=5)


```

```
In [16]: prediction1=la.predict(x_test)  
plt.scatter(y_test,prediction1)
```

```
Out[16]: <matplotlib.collections.PathCollection at 0x7fa521f1e1a0>
```



```
In [17]: las=la.score(x_test,y_test)
```

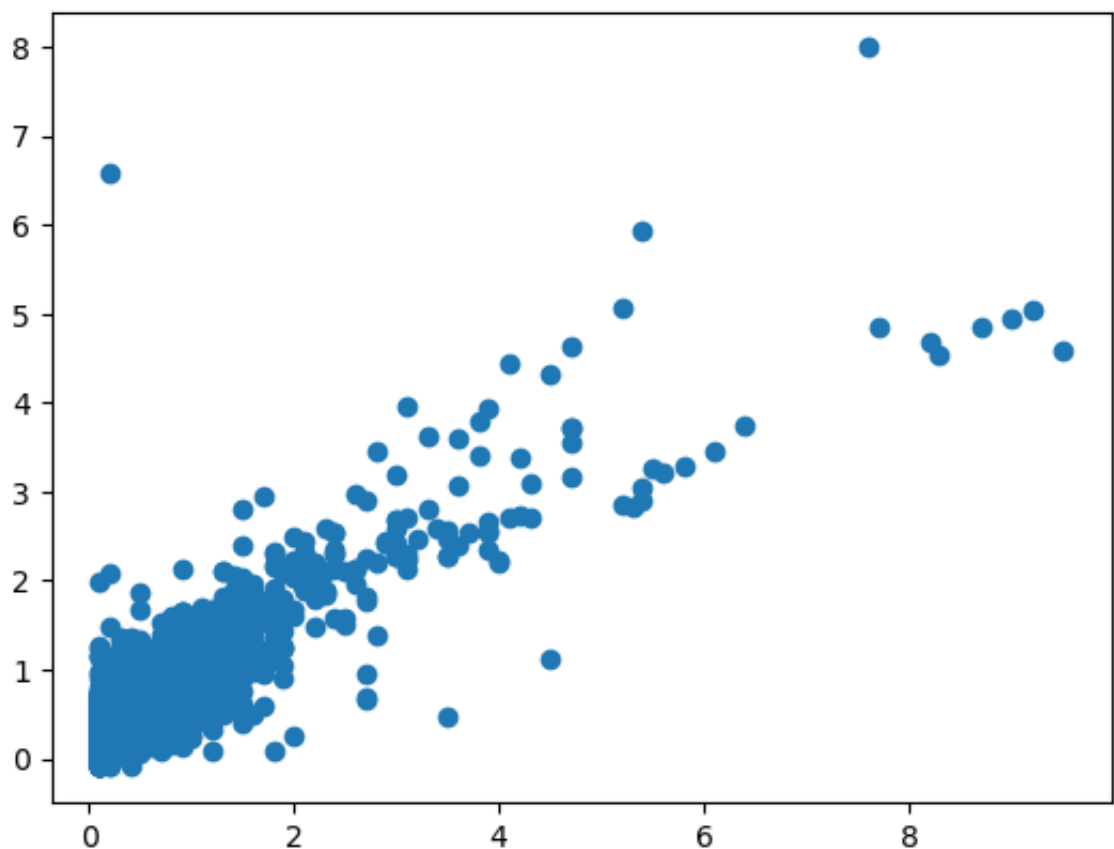
Ridge

```
In [18]: rr=Ridge(alpha=1)  
rr.fit(x_train,y_train)
```

```
Out [18]:  
▼      Ridge  
Ridge(alpha=1)
```

```
In [19]: prediction2=rr.predict(x_test)  
plt.scatter(y_test,prediction2)
```

```
Out [19]: <matplotlib.collections.PathCollection at 0x7fa521ea9780>
```



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

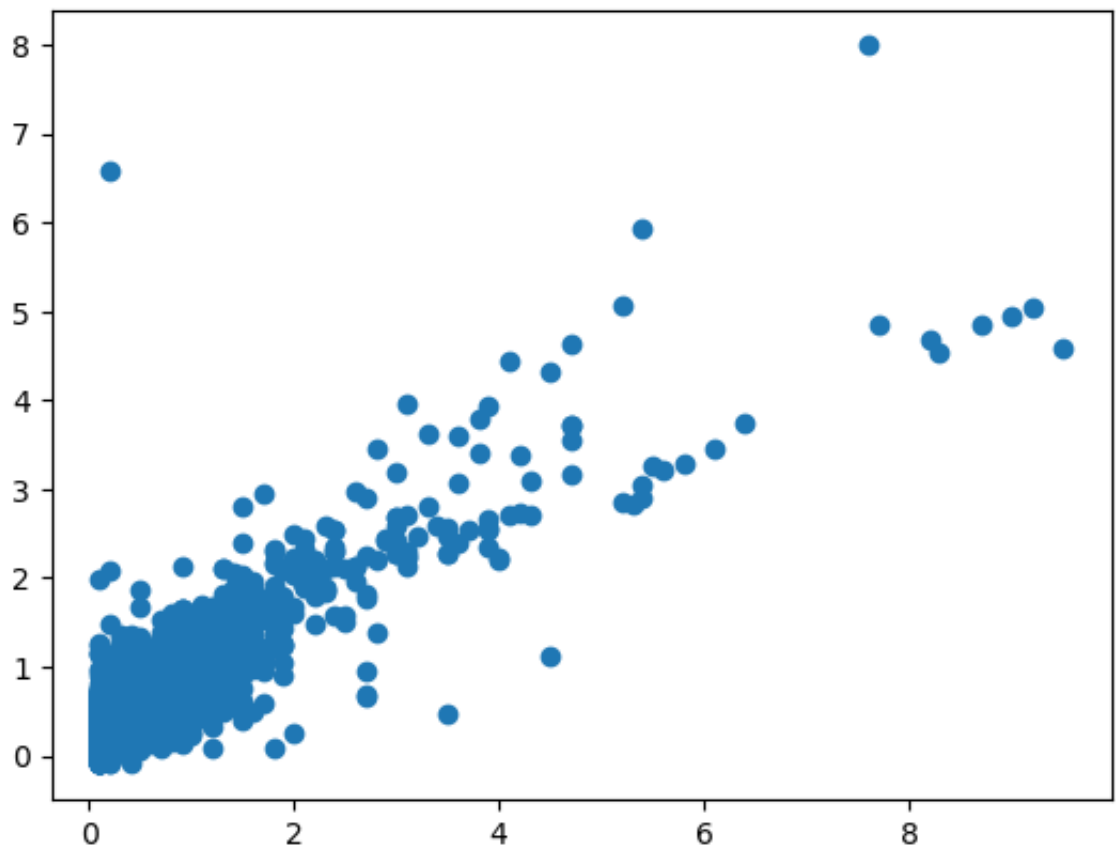
ElasticNet


```
In [21]: en=ElasticNet()  
en.fit(x_train,y_train)
```

```
Out [21]: ▾ ElasticNet  
ElasticNet()
```

```
In [22]: prediction2=rr.predict(x_test)  
plt.scatter(y_test,prediction2)
```

```
Out [22]: <matplotlib.collections.PathCollection at 0x7fa521f5c9d0>
```



```
In [23]: ens=en.score(x_test,y_test)
```

```
In [24]: print(rr.score(x_test,y_test))  
rr.score(x_train,y_train)
```

```
0.7718808700233138
```

```
Out [24]: 0.7690394789280244
```

Logistic

```
In [25]: g={"TCH":{1.0:"Low",2.0:"High"}}
df1=df1.replace(g)
df1["TCH"].value_counts()
```

```
Out[25]: High      8290
Low       7736
Name: TCH, dtype: int64
```

```
In [26]: 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 [27]: lo=LogisticRegression()
lo.fit(x_train,y_train)
```

```
Out[27]: ▾ LogisticRegression
LogisticRegression()
```

```
In [28]: prediction3=lo.predict(x_test)
plt.scatter(y_test,prediction3)
```

```
Out[28]: <matplotlib.collections.PathCollection at 0x7fa5318dba60>
```



```
In [29]: los=lo.score(x_test,y_test)
```

Random Forest

```
In [30]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
```

```
In [31]: g1={"TCH":{"Low":1.0,"High":2.0}}
df1=df1.replace(g1)
```

```
In [32]: 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 [33]: rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

```
Out [33]: ▼ RandomForestClassifier
RandomForestClassifier()
```

```
In [34]: parameter={
    'max_depth': [1,2,4,5,6],
    'min_samples_leaf': [5,10,15,20,25],
    'n_estimators': [10,20,30,40,50]
}
```

```
In [35]: grid_search=GridSearchCV(estimator=rfc,param_grid=parameter,cv=2,sc
grid_search.fit(x_train,y_train)
```

```
Out [35]: ▶ GridSearchCV
▶ estimator: RandomForestClassifier
    ▶ RandomForestClassifier
```

```
In [36]: rfcs=grid_search.best_score_
```

```
In [37]: rfc_best=grid_search.best_estimator_
```

```
In [38]: from sklearn.tree import plot_tree

plt.figure(figsize=(80,40))
plot_tree(rfc_best.estimators_[5],feature_names=x.columns,class_name=
Text(0.13513513513513514, 0.21428571428571427, 'PM10 <= 45.5\ngini
i = 0.095\nsamples = 737\nvalue = [1106, 58]\nclass = Yes'),
Text(0.12162162162162163, 0.07142857142857142, 'gini = 0.089\nsam
ples = 720\nvalue = [1087, 53]\nclass = Yes'),
Text(0.14864864864864866, 0.07142857142857142, 'gini = 0.33\nsamp
les = 17\nvalue = [19, 5]\nclass = Yes'),
Text(0.1891891891891892, 0.21428571428571427, 'TOL <= 2.95\ngini
= 0.451\nsamples = 21\nvalue = [21, 11]\nclass = Yes'),
Text(0.17567567567567569, 0.07142857142857142, 'gini = 0.105\nsam
ples = 13\nvalue = [17, 1]\nclass = Yes'),
Text(0.20270270270270271, 0.07142857142857142, 'gini = 0.408\nsam
ples = 8\nvalue = [4, 10]\nclass = No'),
Text(0.30405405405405406, 0.5, 'PM25 <= 37.5\ngini = 0.33\nsampl
es = 451\nvalue = [566, 149]\nclass = Yes'),
Text(0.2702702702702703, 0.35714285714285715, 'TOL <= 3.85\ngini
= 0.315\nsamples = 440\nvalue = [562, 137]\nclass = Yes'),
Text(0.24324324324324326, 0.21428571428571427, 'NO_2 <= 41.5\ngin
i = 0.174\nsamples = 143\nvalue = [206, 22]\nclass = Yes'),
Text(0.22972972972972974, 0.07142857142857142, 'gini = 0.076\nsam
ples = 79\nvalue = [122, 5]\nclass = Yes'),
```

```
In [39]: print("Linear:",lis)
print("Lasso:",las)
print("Ridge:",rrs)
print("ElasticNet:",ens)
print("Logistic:",los)
print("Random Forest:",rfcs)
```

```
Linear: 0.7718463798012134
Lasso: 0.3625346561801054
Ridge: 0.7718808700233138
ElasticNet: 0.651846033294754
Logistic: 0.512063227953411
Random Forest: 0.9568550543768943
```

Best Model is Random Forest

```
In [40]: df2=pd.read_csv("/Users/bob/Downloads/FP1_air/csvs_per_year/csvs_pe
df2
```

Out[40]:

	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	
0	2016-11-01 01:00:00	NaN	0.7	NaN	NaN	153.0	77.0	NaN	NaN	NaN	7.0	NaN	I
1	2016-11-01 01:00:00	3.1	1.1	2.0	0.53	260.0	144.0	4.0	46.0	24.0	18.0	2.44	
2	2016-11-01 01:00:00	5.9	NaN	7.5	NaN	297.0	139.0	NaN	NaN	NaN	NaN	NaN	:
3	2016-11-01 01:00:00	NaN	1.0	NaN	NaN	154.0	113.0	2.0	NaN	NaN	NaN	NaN	I
4	2016-11-01 01:00:00	NaN	NaN	NaN	NaN	275.0	127.0	2.0	NaN	NaN	18.0	NaN	I
...	
209491	2016-07-01 00:00:00	NaN	0.2	NaN	NaN	2.0	29.0	73.0	NaN	NaN	NaN	NaN	I
209492	2016-07-01 00:00:00	NaN	0.3	NaN	NaN	1.0	29.0	NaN	36.0	NaN	5.0	NaN	I
209493	2016-07-01 00:00:00	NaN	NaN	NaN	NaN	1.0	19.0	71.0	NaN	NaN	NaN	NaN	I
209494	2016-07-01 00:00:00	NaN	NaN	NaN	NaN	6.0	17.0	85.0	NaN	NaN	NaN	NaN	I
209495	2016-07-01 00:00:00	NaN	NaN	NaN	NaN	2.0	46.0	61.0	34.0	NaN	NaN	NaN	I

209496 rows × 14 columns

In [41]: df2.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 209496 entries, 0 to 209495
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype  
---  --
 0   date        209496 non-null object  
 1   BEN         50755 non-null  float64
 2   CO          85999 non-null  float64
 3   EBE         50335 non-null  float64
 4   NMHC        25970 non-null  float64
 5   NO          208614 non-null float64
 6   NO_2        208614 non-null float64
 7   O_3         121197 non-null float64
 8   PM10        102892 non-null float64
 9   PM25        52165 non-null  float64
10   SO_2        86023 non-null  float64
11   TCH         25970 non-null  float64
12   TOL         50662 non-null  float64
13   station     209496 non-null int64  
dtypes: float64(12), int64(1), object(1)
memory usage: 22.4+ MB
```

```
In [42]: df3=df2.dropna()  
df3
```

Out[42]:

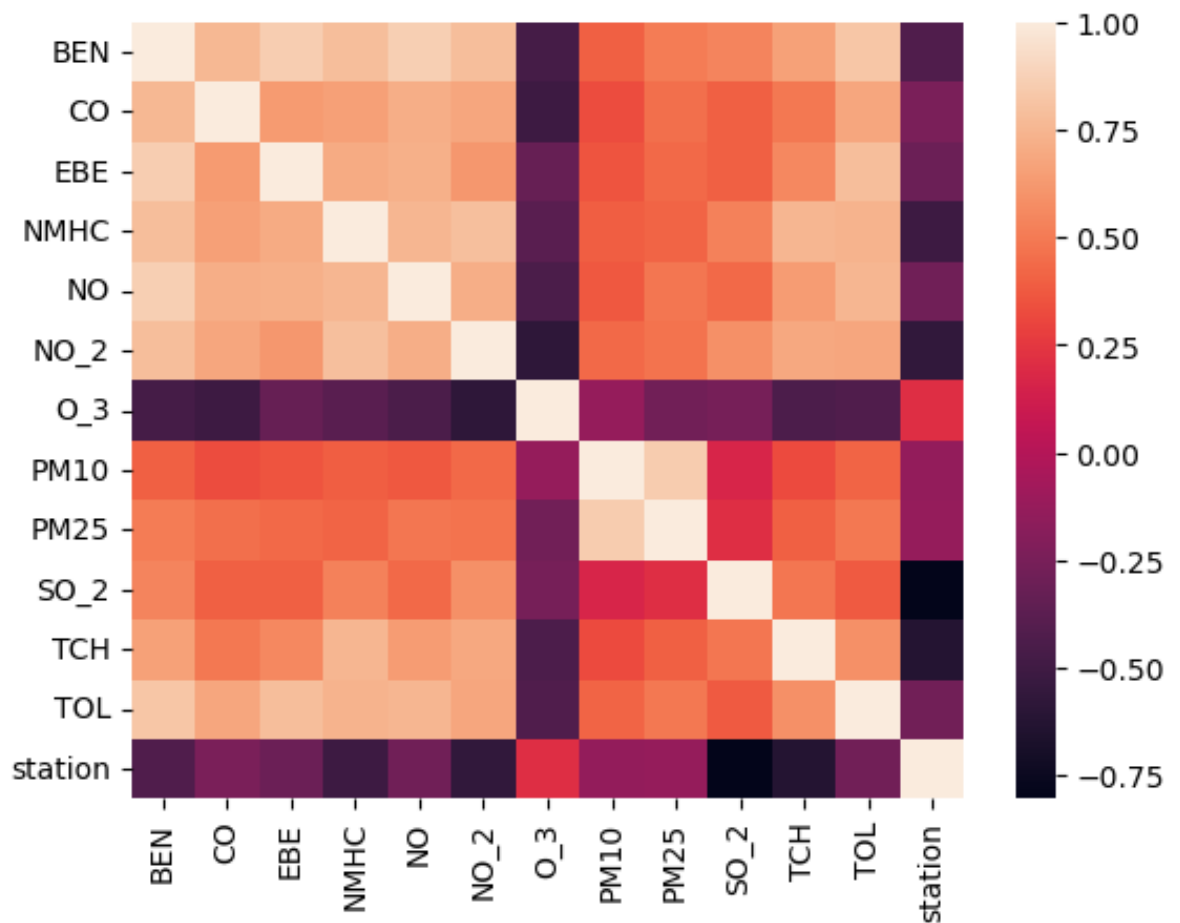
	date	BEN	CO	EBE	NMHC	NO	NO_2	O_3	PM10	PM25	SO_2	TCH	T
1	2016-11-01 01:00:00	3.1	1.1	2.0	0.53	260.0	144.0	4.0	46.0	24.0	18.0	2.44	1
6	2016-11-01 01:00:00	0.7	0.8	0.4	0.13	57.0	66.0	3.0	23.0	15.0	4.0	1.35	1
25	2016-11-01 02:00:00	2.7	1.0	2.1	0.40	139.0	114.0	4.0	37.0	21.0	14.0	2.30	1
30	2016-11-01 02:00:00	0.7	0.7	0.4	0.13	48.0	59.0	3.0	23.0	15.0	3.0	1.35	1
49	2016-11-01 03:00:00	1.7	0.8	1.4	0.25	53.0	90.0	4.0	31.0	19.0	10.0	1.95	1
...
209430	2016-06-30 22:00:00	0.1	0.2	0.1	0.02	1.0	5.0	97.0	19.0	12.0	2.0	1.15	1
209449	2016-06-30 23:00:00	0.6	0.4	0.3	0.15	14.0	63.0	54.0	29.0	13.0	16.0	1.48	1
209454	2016-06-30 23:00:00	0.1	0.2	0.1	0.02	1.0	7.0	91.0	16.0	9.0	2.0	1.15	1
209473	2016-07-01 00:00:00	0.6	0.4	0.3	0.16	11.0	68.0	45.0	24.0	14.0	16.0	1.50	1
209478	2016-07-01 00:00:00	0.1	0.2	0.1	0.02	1.0	6.0	89.0	16.0	9.0	2.0	1.15	1

16932 rows × 14 columns

```
In [43]: df3=df3.drop(["date"],axis=1)
```

```
In [44]: sns.heatmap(df3.corr())
```

```
Out [44]: <Axes: >
```



```
In [45]: 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

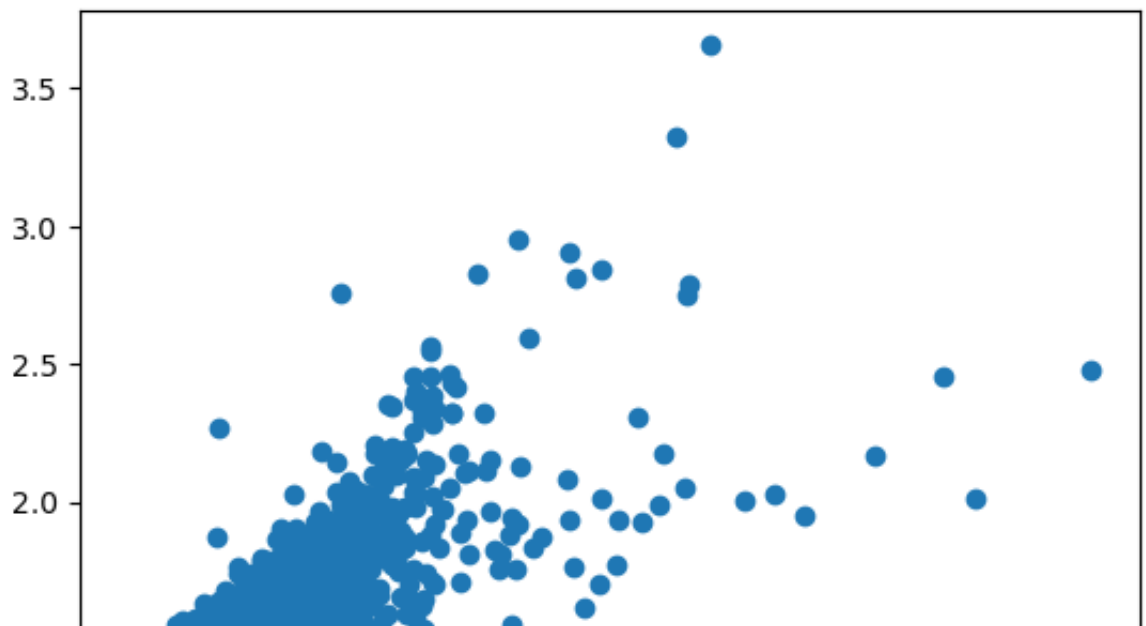
```
In [46]: li=LinearRegression()
li.fit(x_train,y_train)
```

```
Out [46]: ▼ LinearRegression
LinearRegression()
```



```
In [47]: prediction=li.predict(x_test)
plt.scatter(y_test,prediction)
```

```
Out[47]: <matplotlib.collections.PathCollection at 0x7fa501f85330>
```



```
In [48]: lis=li.score(x_test,y_test)
```

```
In [49]: df3["TCH"].value_counts()
```

```
Out[49]: 1.16    757
         1.18    701
         1.17    683
         1.19    618
         1.15    577
         ...
         3.28     1
         3.96     1
         4.82     1
         2.69     1
         4.55     1
         Name: TCH, Length: 217, dtype: int64
```

```
In [50]: df3.loc[df3["TCH"]<1.40,"TCH"]=1
df3.loc[df3["TCH"]>1.40,"TCH"]=2
df3["TCH"].value_counts()
```

```
Out[50]: 1.0    10002
         2.0     6930
         Name: TCH, dtype: int64
```

```
In [ ]:
```

Lasso

```
In [51]: la=Lasso(alpha=5)  
la.fit(x_train,y_train)
```

```
Out [51]: 

▼



Lasso

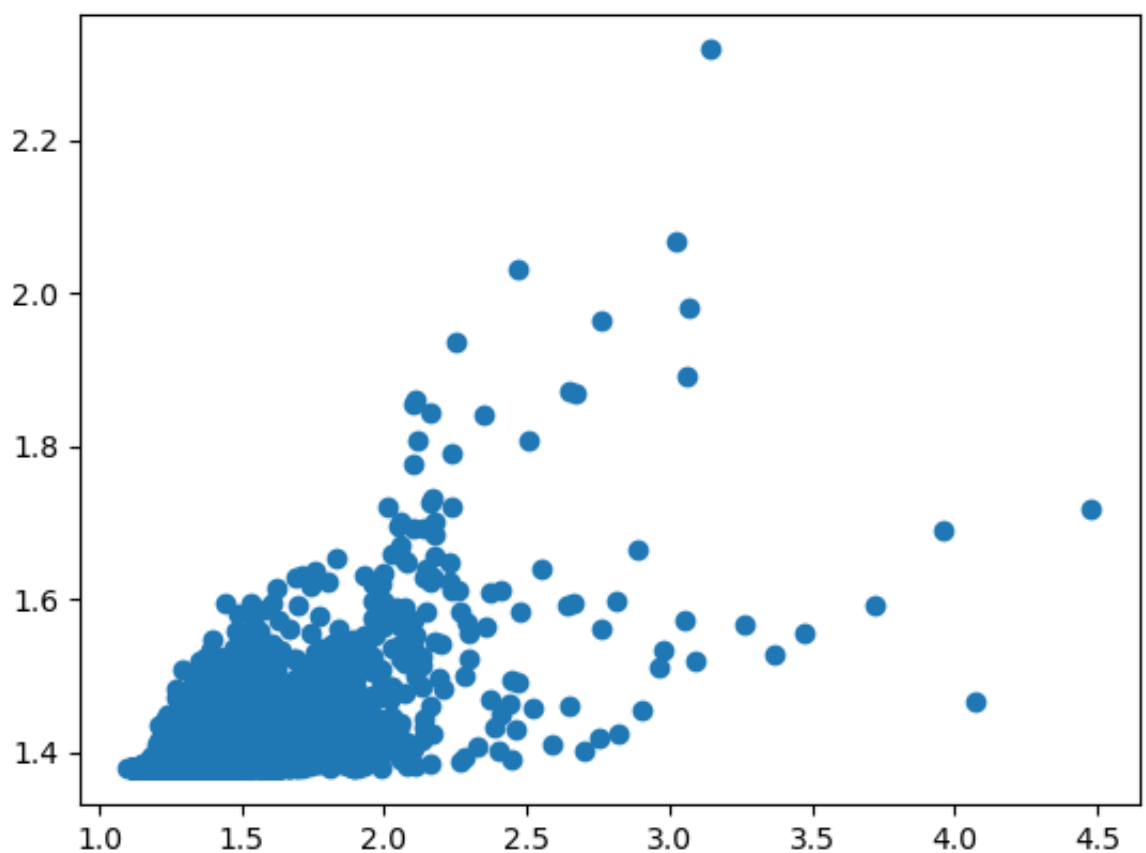


Lasso(alpha=5)


```

```
In [52]: prediction1=la.predict(x_test)  
plt.scatter(y_test,prediction1)
```

```
Out [52]: <matplotlib.collections.PathCollection at 0x7fa520f1c7c0>
```



```
In [53]: las=la.score(x_test,y_test)
```

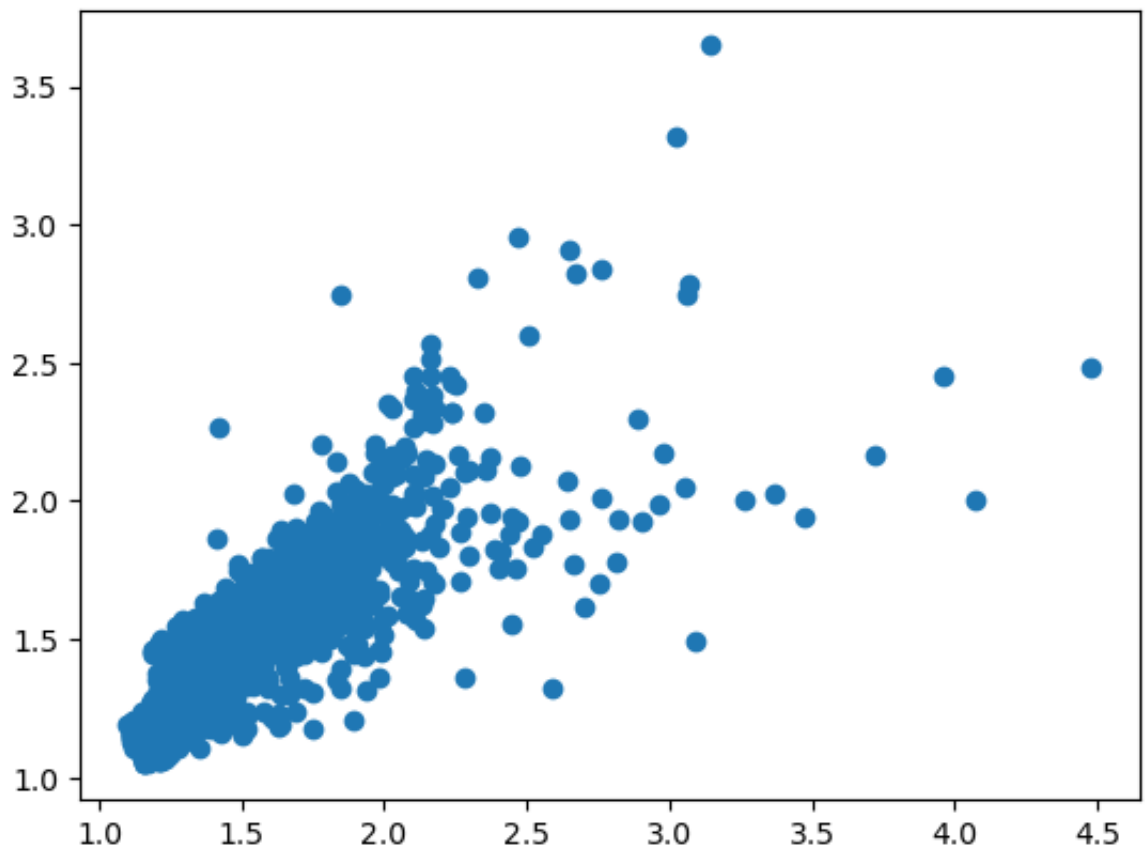
Ridge

```
In [54]: rr=Ridge(alpha=1)
         rr.fit(x_train,y_train)
```

```
Out [54]: ▼ Ridge
          Ridge(alpha=1)
```

```
In [55]: prediction2=rr.predict(x_test)
         plt.scatter(y_test,prediction2)
```

```
Out [55]: <matplotlib.collections.PathCollection at 0x7fa5143c2320>
```



```
In [56]: rrs=rr.score(x_test,y_test)
```

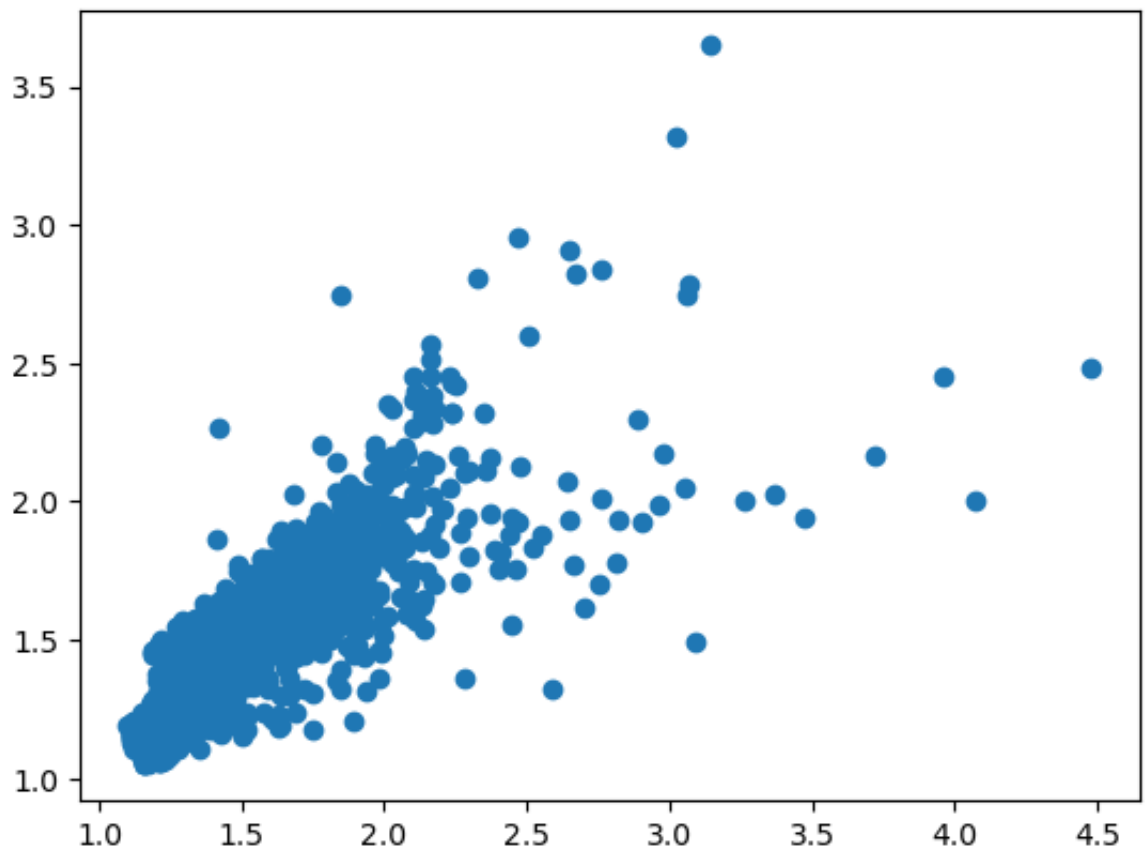
ElasticNet

```
In [57]: en=ElasticNet()
         en.fit(x_train,y_train)
```

```
Out [57]: ▼ ElasticNet
          ElasticNet()
```

```
In [58]: prediction2=rr.predict(x_test)
plt.scatter(y_test,prediction2)
```

Out[58]: <matplotlib.collections.PathCollection at 0x7fa5143c2470>



```
In [59]: ens=en.score(x_test,y_test)
```

```
In [60]: print(rr.score(x_test,y_test))
rr.score(x_train,y_train)
```

0.7612362215703397

Out[60]: 0.7491431239624469

Logistic

```
In [61]: g={"TCH":{1.0:"Low",2.0:"High"}}
df3=df3.replace(g)
df3["TCH"].value_counts()
```

Out[61]: Low 10002
High 6930
Name: TCH, dtype: int64

```
In [62]: 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 [63]: lo=LogisticRegression()
         lo.fit(x_train,y_train)
```

```
Out [63]: ▾ LogisticRegression
          LogisticRegression()
```

```
In [64]: prediction3=lo.predict(x_test)
         plt.scatter(y_test,prediction3)
```

```
Out [64]: <matplotlib.collections.PathCollection at 0x7fa52275df90>
```



```
In [65]: los=lo.score(x_test,y_test)
```

Random Forest

```
In [66]: from sklearn.ensemble import RandomForestClassifier
         from sklearn.model_selection import GridSearchCV
```

```
In [67]: g1={"TCH":{"Low":1.0,"High":2.0}}
df3=df3.replace(g1)
```

```
In [68]: 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 [69]: rfc=RandomForestClassifier()
rfc.fit(x_train,y_train)
```

```
Out [69]: ▾ RandomForestClassifier
RandomForestClassifier()
```

```
In [70]: parameter={
    'max_depth':[1,2,4,5,6],
    'min_samples_leaf':[5,10,15,20,25],
    'n_estimators':[10,20,30,40,50]
}
```

```
In [71]: grid_search=GridSearchCV(estimator=rfc,param_grid=parameter,cv=2,sc
grid_search.fit(x_train,y_train)
```

```
Out [71]: ▸ GridSearchCV
▸ estimator: RandomForestClassifier
    ▸ RandomForestClassifier
```

```
In [72]: rfcs=grid_search.best_score_
```

```
In [73]: rfc_best=grid_search.best_estimator_
```

```
In [74]: from sklearn.tree import plot_tree

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

```
Out [74]: [Text(0.49604430379746833, 0.9285714285714286, 'station <= 2807901
6.0\ngini = 0.482\nsamples = 7503\nvalue = [7057, 4795]\nclass = Y
es'),
Text(0.22784810126582278, 0.7857142857142857, 'PM10 <= 13.5\ngini
= 0.368\nsamples = 3763\nvalue = [1434, 4470]\nclass = No'),
Text(0.10126582278481013, 0.6428571428571429, 'PM10 <= 6.5\ngini
= 0.498\nsamples = 1344\nvalue = [977, 1110]\nclass = No'),
Text(0.05063291139240506, 0.5, 'SO_2 <= 12.5\ngini = 0.46\nsampl
s = 374\nvalue = [374, 209]\nclass = Yes'),
Text(0.02531645569620253, 0.35714285714285715, 'NO <= 2.5\ngini =
0.013\nsamples = 100\nvalue = [1, 152]\nclass = No'),
Text(0.012658227848101266, 0.21428571428571427, 'gini = 0.062\nsa
mples = 19\nvalue = [1, 30]\nclass = No'),
Text(0.0379746835443038, 0.21428571428571427, 'gini = 0.0\nsampl
s = 81\nvalue = [0, 122]\nclass = No'),
Text(0.0759493670886076, 0.35714285714285715, 'NO_2 <= 63.5\ngini
= 0.23\nsamples = 274\nvalue = [373, 57]\nclass = Yes'),
Text(0.06329113924050633, 0.21428571428571427, 'SO_2 <= 15.5\ngin
i = 0.19\nsamples = 259\nvalue = [361, 43]\nclass = Yes'),
Text(0.05063291139240506, 0.35714285714285715, 'gini = 0.137\nsam
```

```
In [75]: print("Linear:",lis)
print("Lasso:",las)
print("Ridge:",rrs)
print("ElasticNet:",ens)
print("Logistic:",los)
print("Random Forest:",rfcs)
```

```
Linear: 0.7614033018124887
Lasso: 0.2149866611637189
Ridge: 0.7612362215703397
ElasticNet: 0.5893135218045396
Logistic: 0.5923228346456693
Random Forest: 0.9214478569017888
```

Best model is Random Forest

