

```
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	MXV	NMHC	NO_2	NOx	OXY	O_3
0	2003-03-01 01:00:00	NaN	1.72	NaN	NaN	NaN	73.900002	316.299988	NaN	10.550000
1	2003-03-01 01:00:00	NaN	1.45	NaN	NaN	0.26	72.110001	250.000000	0.73	6.720000
2	2003-03-01 01:00:00	NaN	1.57	NaN	NaN	NaN	80.559998	224.199997	NaN	21.049999
3	2003-03-01 01:00:00	NaN	2.45	NaN	NaN	NaN	78.370003	450.399994	NaN	4.220000
4	2003-03-01 01:00:00	NaN	3.26	NaN	NaN	NaN	96.250000	479.100006	NaN	8.460000
...
243979	2003-10-01 00:00:00	0.20	0.16	2.01	3.17	0.02	31.799999	32.299999	1.68	34.049999
243980	2003-10-01 00:00:00	0.32	0.08	0.36	0.72	NaN	10.450000	14.760000	1.00	34.610001
243981	2003-10-01 00:00:00	NaN	NaN	NaN	NaN	0.07	34.639999	50.810001	NaN	32.160000
243982	2003-10-01 00:00:00	NaN	NaN	NaN	NaN	0.07	32.580002	41.020000	NaN	NaN
243983	2003-10-01 00:00:00	1.00	0.29	2.15	6.41	0.07	37.150002	56.849998	2.28	21.480000

243984 rows × 16 columns

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 243984 entries, 0 to 243983
Data columns (total 16 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   date        243984 non-null  object
 1   BEN         69745 non-null   float64
 2   CO          225340 non-null  float64
 3   EBE         61244 non-null   float64
 4   MXY         42045 non-null   float64
 5   NMHC        111951 non-null  float64
 6   NO_2        242625 non-null  float64
 7   NOx         242629 non-null  float64
 8   OXY         42072 non-null   float64
 9   O_3         234131 non-null  float64
10  PM10        240896 non-null  float64
11  PXY         42063 non-null   float64
12  SO_2        242729 non-null  float64
13  TCH         111991 non-null  float64
14  TOL         69439 non-null   float64
15  station     243984 non-null  int64
dtypes: float64(14), int64(1), object(1)
memory usage: 29.8+ MB
```

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

Out [4]:

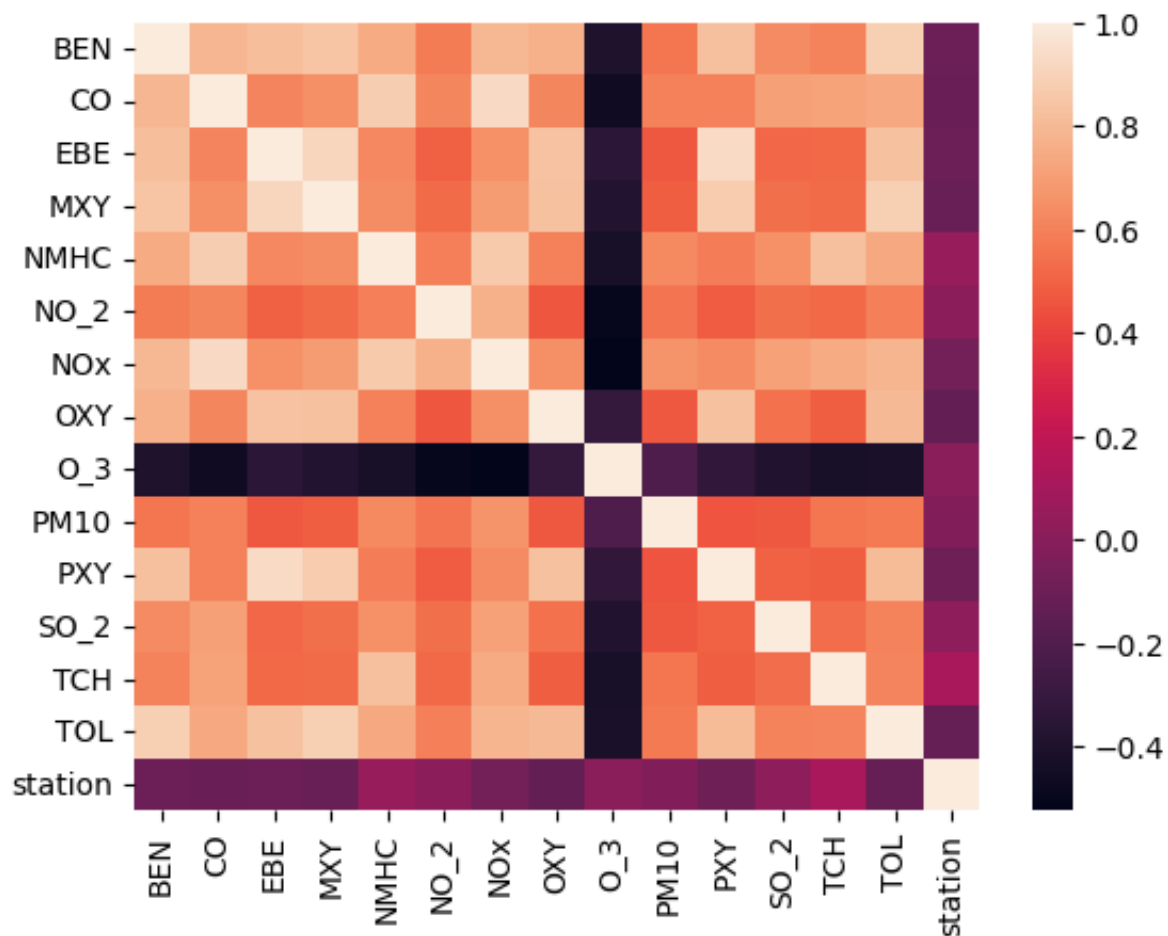
	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3
5	2003-03-01 01:00:00	8.41	1.94	9.83	21.49	0.45	90.300003	384.899994	9.48	9.950000
23	2003-03-01 01:00:00	3.46	1.27	3.43	7.08	0.18	54.250000	173.300003	3.37	6.540000
27	2003-03-01 01:00:00	6.39	1.79	5.75	10.88	0.33	75.459999	281.100006	3.68	6.690000
33	2003-03-01 02:00:00	7.42	1.47	10.63	24.73	0.35	83.309998	277.200012	11.00	9.900000
51	2003-03-01 02:00:00	3.62	1.29	3.20	7.08	0.19	42.209999	166.300003	3.41	6.380000
...
243955	2003-09-30 23:00:00	1.75	0.41	3.07	9.38	0.09	46.290001	77.709999	3.11	18.280001
243957	2003-10-01 00:00:00	2.35	0.60	3.88	10.86	0.11	61.240002	133.100006	0.89	10.900000
243961	2003-10-01 00:00:00	2.97	0.82	4.53	10.88	0.05	36.529999	131.300003	5.52	12.940000
243979	2003-10-01 00:00:00	0.20	0.16	2.01	3.17	0.02	31.799999	32.299999	1.68	34.049999
243983	2003-10-01 00:00:00	1.00	0.29	2.15	6.41	0.07	37.150002	56.849998	2.28	21.480000

33010 rows × 16 columns

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

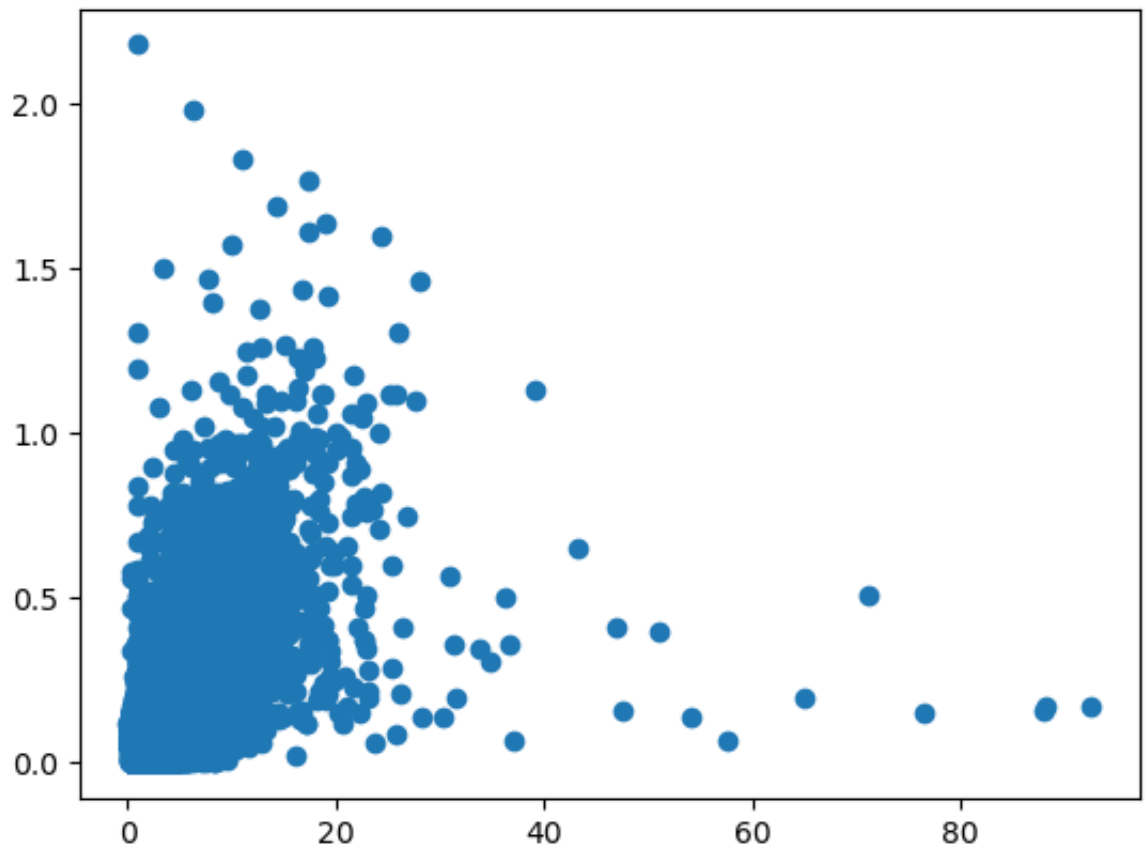
```
In [6]: sns.heatmap(df1.corr())
```

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



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

```
Out[7]: [<matplotlib.lines.Line2D at 0x7fd78926b370>]
```



```
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)
```

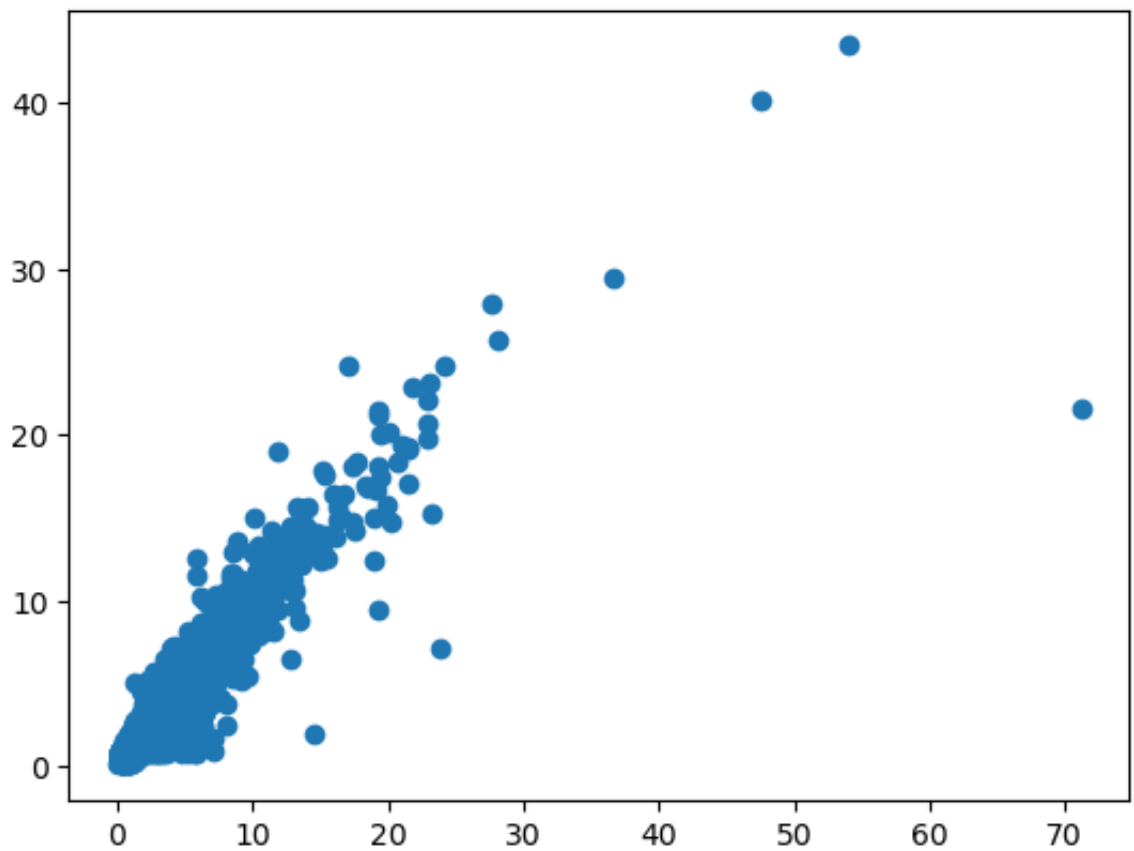
Linear

```
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 0x7fd789242350>
```



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

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

```
Out[13]: 1.30    1344
         1.31    1342
         1.32    1281
         1.27    1279
         1.29    1262
         ...
         3.58         1
         3.28         1
         3.43         1
         3.03         1
         3.59         1
         Name: TCH, Length: 243, 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]: 1.0    21614  
        2.0    11396  
        Name: TCH, dtype: int64
```

Lasso

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

```
Out[15]: 

▼ Lasso

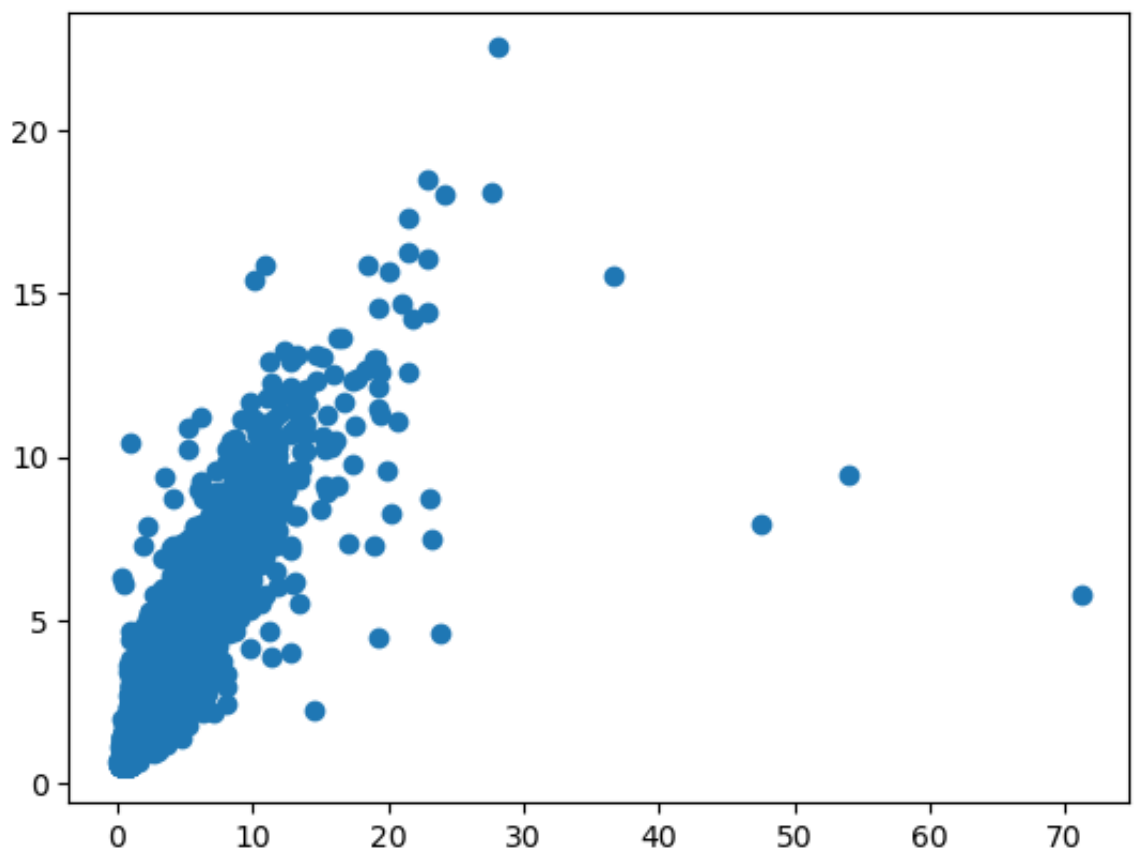


Lasso(alpha=5)


```

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

```
Out[16]: <matplotlib.collections.PathCollection at 0x7fd7894ed1e0>
```



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

Ridge

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

```
Out[18]:
```

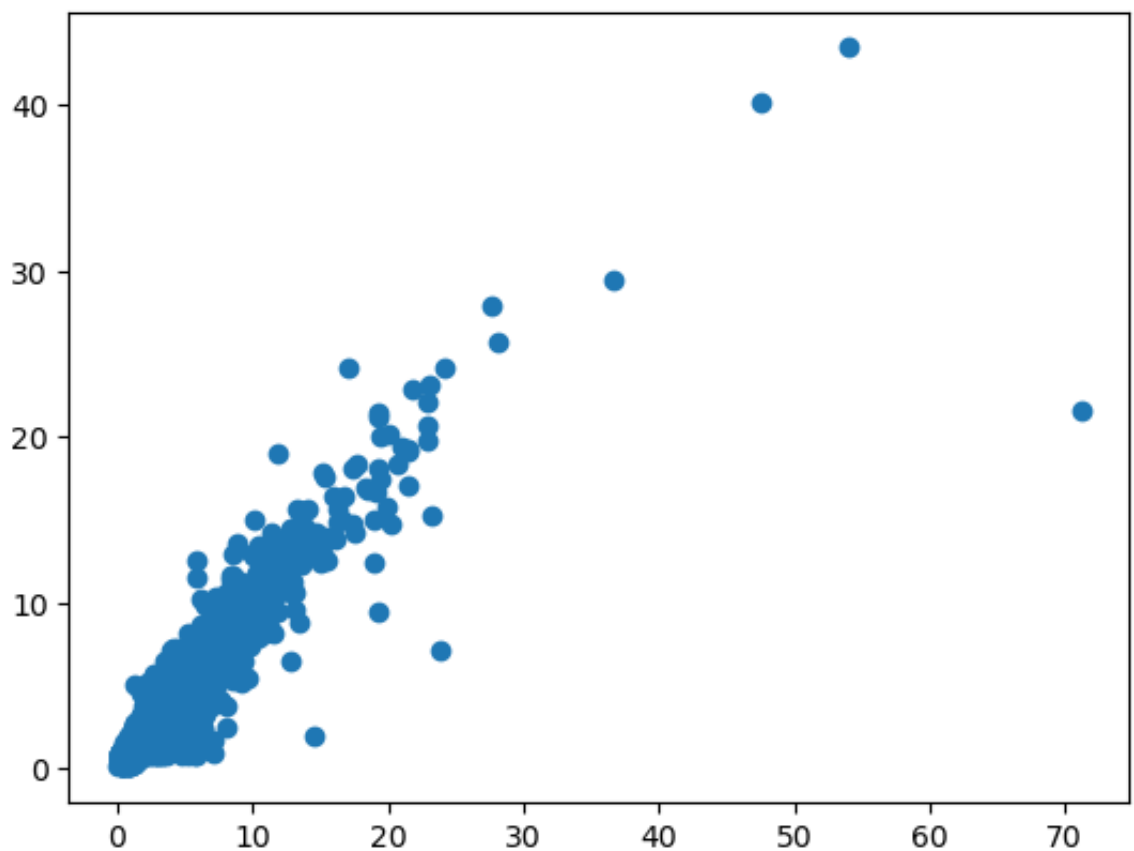
▼

Ridge

Ridge(alpha=1)

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

```
Out[19]: <matplotlib.collections.PathCollection at 0x7fd7892e4f10>
```



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

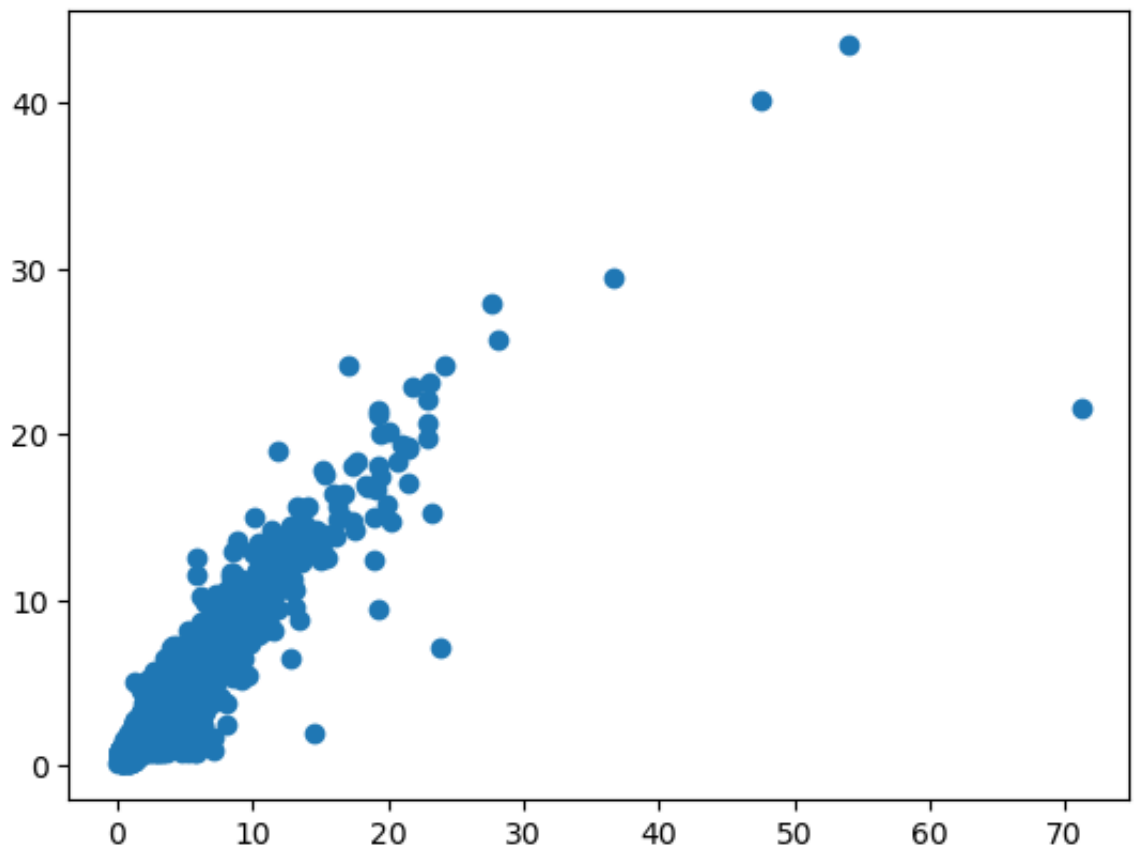
ElasticNet


```
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 0x7fd7896283a0>
```



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

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

```
0.9046329773630202
```

```
Out [24]: 0.9191186610151827
```

Logistic

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

```
Out[25]: Low      21614
        High      11396
        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 0x7fd7896fe560>
```



```
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_nam
```

```
Out[38]: [Text(0.5080818965517241, 0.9285714285714286, 'CO <= 0.865\ngini =
0.449\nsamples = 14577\nvalue = [15259, 7848]\nclass = Yes'),
Text(0.2510775862068966, 0.7857142857142857, 'TOL <= 7.455\ngini
= 0.305\nsamples = 10478\nvalue = [13484, 3119]\nclass = Yes'),
Text(0.12284482758620689, 0.6428571428571429, 'NOx <= 88.445\ngin
i = 0.174\nsamples = 6820\nvalue = [9744, 1036]\nclass = Yes'),
Text(0.06896551724137931, 0.5, 'O_3 <= 19.005\ngini = 0.116\nsamp
les = 5846\nvalue = [8699, 572]\nclass = Yes'),
Text(0.034482758620689655, 0.35714285714285715, 'NO_2 <= 21.0\ngi
ni = 0.35\nsamples = 555\nvalue = [710, 208]\nclass = Yes'),
Text(0.017241379310344827, 0.21428571428571427, 'NOx <= 11.645\ng
ini = 0.412\nsamples = 20\nvalue = [9, 22]\nclass = No'),
Text(0.008620689655172414, 0.07142857142857142, 'gini = 0.26\nsam
ples = 10\nvalue = [2, 11]\nclass = No'),
Text(0.02586206896551724, 0.07142857142857142, 'gini = 0.475\nsam
ples = 10\nvalue = [7, 11]\nclass = No'),
Text(0.05172413793103448, 0.21428571428571427, 'PM10 <= 32.735\ng
ini = 0.331\nsamples = 535\nvalue = [701, 186]\nclass = Yes'),
Text(0.04310344827586207, 0.07142857142857142, 'gini = 0.282\nsam
ples = 14577\nvalue = [15259, 7848]\nclass = Yes')]
```

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

```
Linear: 0.9046345936860489
Lasso: 0.7357327274811196
Ridge: 0.9046329773630202
ElasticNet: 0.8875959784048247
Logistic: 0.6531354135110573
Random Forest: 0.8828925435665907
```

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	MXV	NMHC	NO_2	NOx	OXY	O_3
0	2004-08-01 01:00:00	NaN	0.66	NaN	NaN	NaN	89.550003	118.900002	NaN	40.020000
1	2004-08-01 01:00:00	2.66	0.54	2.99	6.08	0.18	51.799999	53.860001	3.28	51.689999
2	2004-08-01 01:00:00	NaN	1.02	NaN	NaN	NaN	93.389999	138.600006	NaN	20.860001
3	2004-08-01 01:00:00	NaN	0.53	NaN	NaN	NaN	87.290001	105.000000	NaN	36.730000
4	2004-08-01 01:00:00	NaN	0.17	NaN	NaN	NaN	34.910000	35.349998	NaN	86.269997
...
245491	2004-06-01 00:00:00	0.75	0.21	0.85	1.55	0.07	59.580002	64.389999	0.66	33.029999
245492	2004-06-01 00:00:00	2.49	0.75	2.44	4.57	NaN	97.139999	146.899994	2.34	7.740000
245493	2004-06-01 00:00:00	NaN	NaN	NaN	NaN	0.13	102.699997	132.600006	NaN	17.809999
245494	2004-06-01 00:00:00	NaN	NaN	NaN	NaN	0.09	82.599998	102.599998	NaN	NaN
245495	2004-06-01 00:00:00	3.01	0.67	2.78	5.12	0.20	92.550003	141.000000	2.60	11.460000

245496 rows × 11 columns

```
In [41]: df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 245496 entries, 0 to 245495
Data columns (total 17 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   date        245496 non-null object
 1   BEN         65158 non-null  float64
 2   CO          226043 non-null float64
 3   EBE         56781 non-null  float64
 4   MXY         39867 non-null  float64
 5   NMHC        107630 non-null float64
 6   NO_2        243280 non-null float64
 7   NOx         243283 non-null float64
 8   OXY         39882 non-null  float64
 9   O_3         233811 non-null float64
10  PM10        234655 non-null float64
11  PM25        58145 non-null  float64
12  PXY         39891 non-null  float64
13  SO_2        243402 non-null float64
14  TCH         107650 non-null float64
15  TOL         64914 non-null  float64
16  station     245496 non-null int64
dtypes: float64(15), int64(1), object(1)
memory usage: 31.8+ MB
```

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

Out[42]:

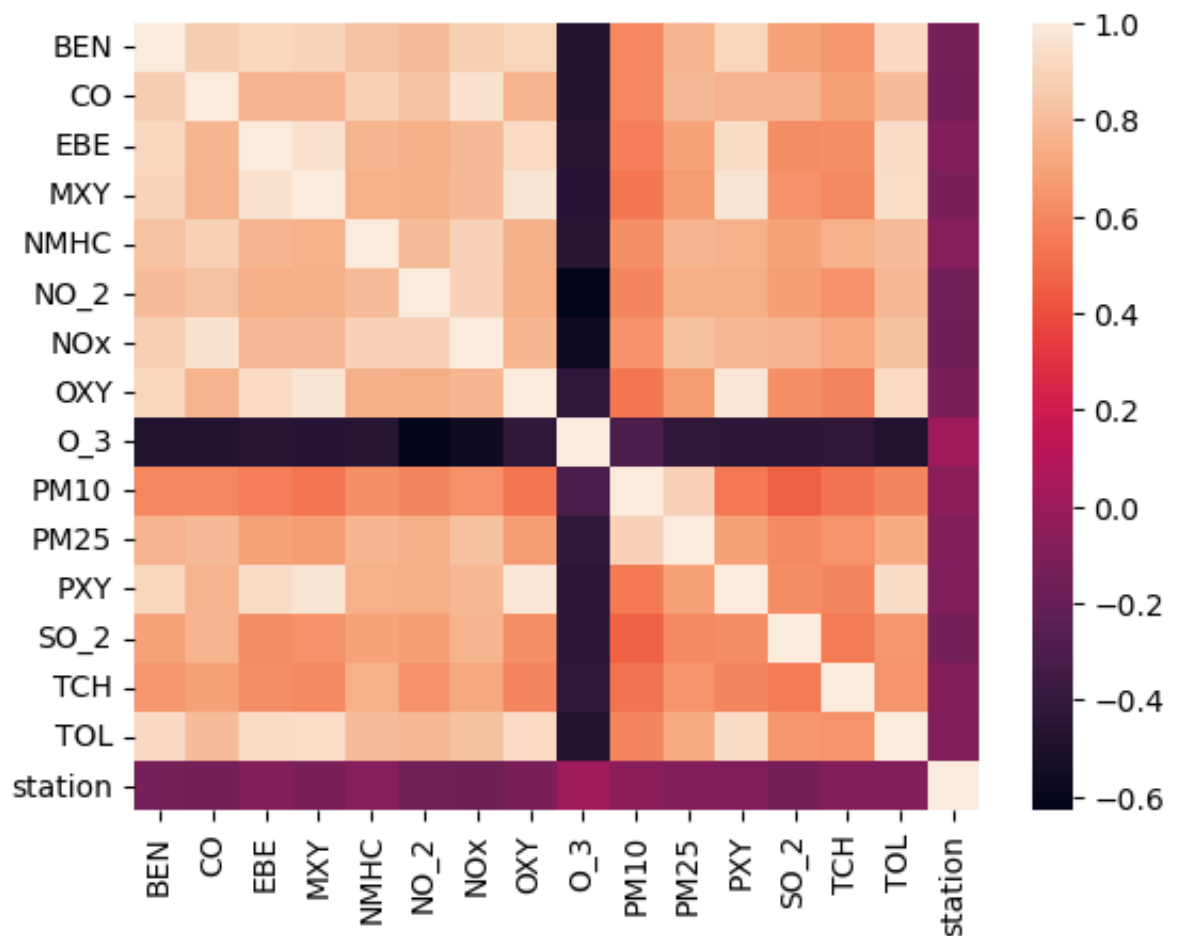
	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3
5	2004-08-01 01:00:00	3.24	0.63	5.55	9.72	0.06	103.800003	144.800003	5.04	32.480000
22	2004-08-01 01:00:00	0.55	0.36	0.54	0.86	0.07	31.980000	32.799999	0.50	79.040001
26	2004-08-01 01:00:00	1.80	0.46	2.28	4.62	0.21	62.259998	75.470001	2.47	54.419998
32	2004-08-01 02:00:00	1.94	0.67	3.14	4.91	0.06	113.500000	165.800003	2.56	26.980000
49	2004-08-01 02:00:00	0.29	0.30	0.47	0.76	0.07	33.919998	34.840000	0.46	75.570000
...
245463	2004-05-31 23:00:00	0.62	0.08	0.54	0.70	0.04	44.360001	45.450001	0.42	43.419998
245467	2004-05-31 23:00:00	2.39	0.67	2.49	3.92	0.20	89.809998	132.800003	2.09	14.740000
245473	2004-06-01 00:00:00	3.72	1.12	4.33	8.79	0.24	113.900002	253.600006	4.51	9.380000
245491	2004-06-01 00:00:00	0.75	0.21	0.85	1.55	0.07	59.580002	64.389999	0.66	33.029999
245495	2004-06-01 00:00:00	3.01	0.67	2.78	5.12	0.20	92.550003	141.000000	2.60	11.460000

19397 rows × 17 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

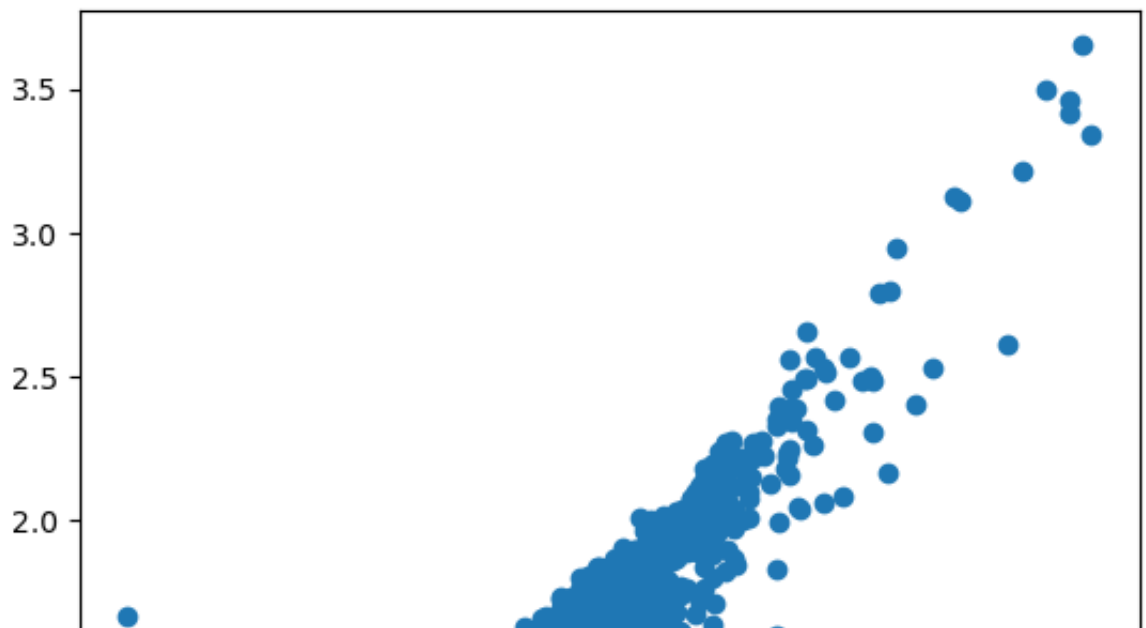
```
In [46]: li=LinearRegression()
li.fit(x_train,y_train)
```

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



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

```
Out[47]: <matplotlib.collections.PathCollection at 0x7fd7b95cd690>
```



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

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

```
Out[49]: 1.34    740
1.33    714
1.35    708
1.37    688
1.36    679
...
3.65     1
2.86     1
2.87     1
3.86     1
2.66     1
Name: TCH, Length: 191, 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    11861
2.0     7536
Name: TCH, dtype: int64
```

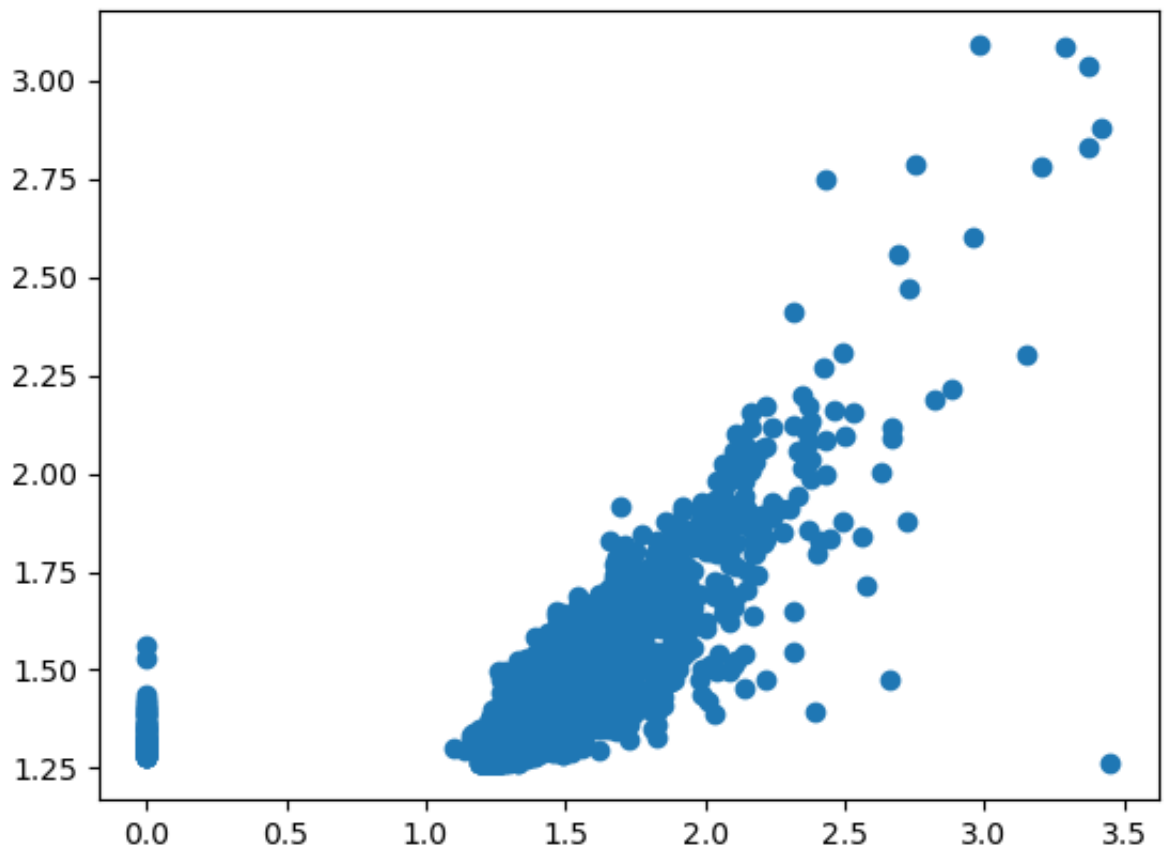
Lasso

```
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 0x7fd789f275b0>
```



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

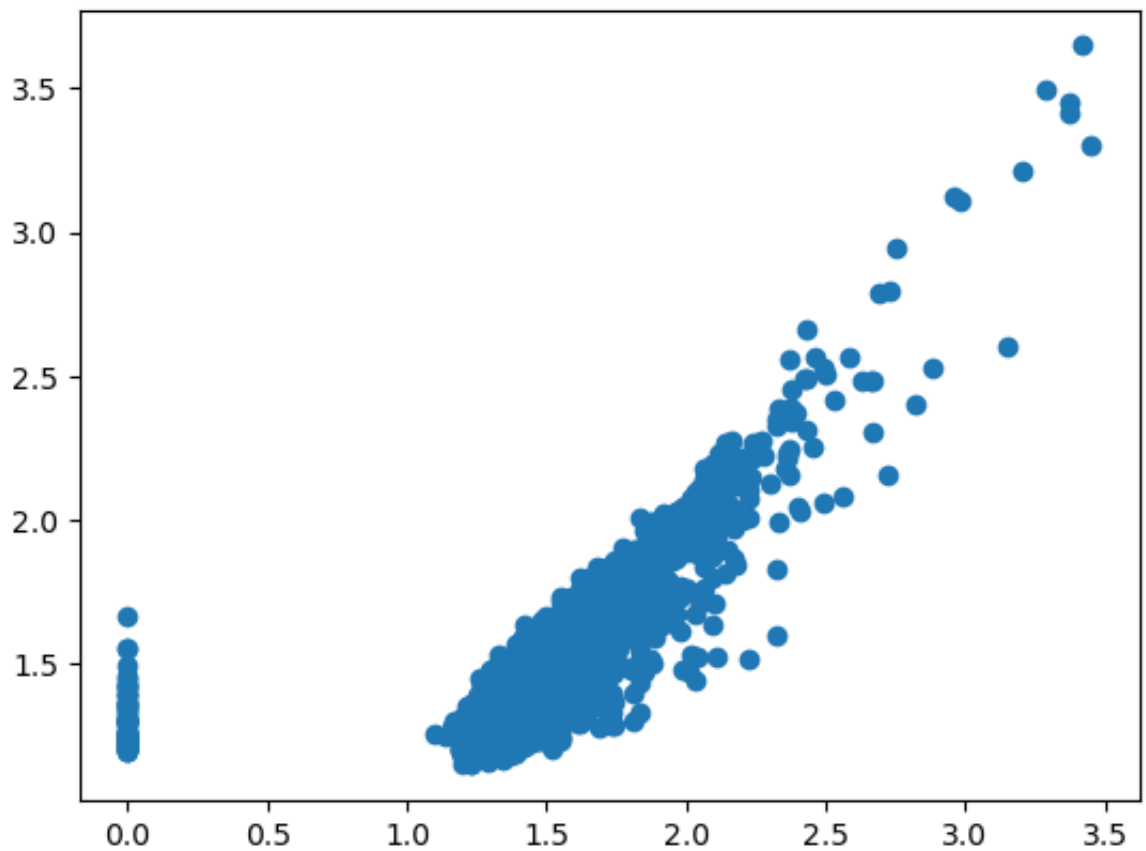
Ridge

```
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 0x7fd789f77fd0>
```



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

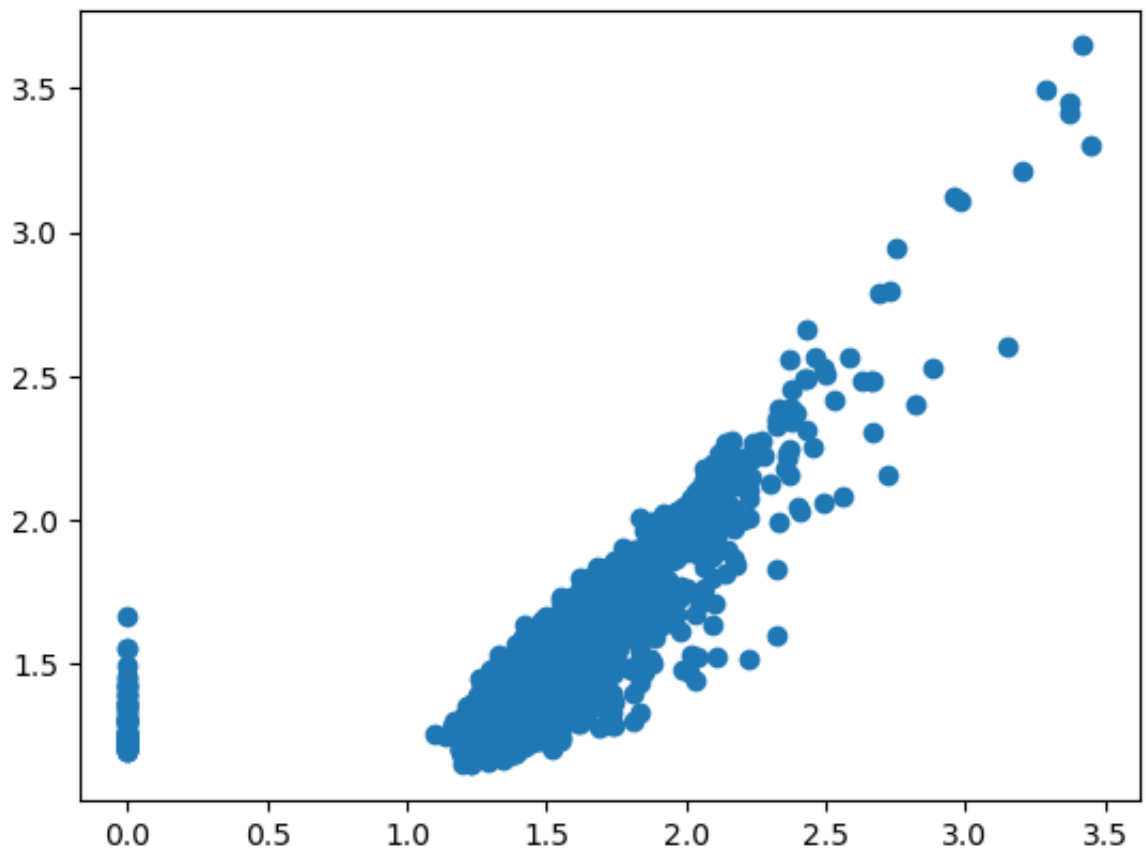
ElasticNet

```
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 0x7fd789dd48b0>
```



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

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

```
0.5675301031753455
```

```
Out[60]: 0.5998074340539438
```

Logistic

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

```
Out[61]: Low      11861
High       7536
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 0x7fd7b95cde40>
```



```
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(0.49848790322580644, 0.9285714285714286, 'NO_2 <= 66.86\ngini = 0.474\nsamples = 8629\nvalue = [8336, 5241]\nclass = Yes'),
Text(0.2550403225806452, 0.7857142857142857, 'TOL <= 9.875\ngini = 0.204\nsamples = 5025\nvalue = [6998, 910]\nclass = Yes'),
Text(0.12903225806451613, 0.6428571428571429, 'NO_2 <= 40.035\ngini = 0.138\nsamples = 4535\nvalue = [6608, 531]\nclass = Yes'),
Text(0.06451612903225806, 0.5, 'BEN <= 1.685\ngini = 0.037\nsamples = 2569\nvalue = [4007, 77]\nclass = Yes'),
Text(0.03225806451612903, 0.35714285714285715, 'OXY <= 1.155\ngini = 0.029\nsamples = 2516\nvalue = [3937, 59]\nclass = Yes'),
Text(0.016129032258064516, 0.21428571428571427, 'NO_2 <= 23.005\ngini = 0.012\nsamples = 1905\nvalue = [3020, 19]\nclass = Yes'),
Text(0.008064516129032258, 0.07142857142857142, 'gini = 0.0\nsamples = 1173\nvalue = [1877, 0]\nclass = Yes'),
Text(0.024193548387096774, 0.07142857142857142, 'gini = 0.032\nsamples = 732\nvalue = [1143, 19]\nclass = Yes'),
Text(0.04838709677419355, 0.21428571428571427, 'NMHC <= 0.165\ngini = 0.08\nsamples = 611\nvalue = [917, 40]\nclass = Yes'),
Text(0.04032258064516129, 0.07142857142857142, 'gini = 0.068\nsam
```

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

```
Linear: 0.5675063934890017
Lasso: 0.47670581670145473
Ridge: 0.5675301031753455
ElasticNet: 0.5009515759428786
Logistic: 0.6051546391752577
Random Forest: 0.8972527875997542
```

Best model is Random Forest

