

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import Ridge,Lasso
from sklearn.linear_model import ElasticNet
from sklearn import metrics
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.tree import plot_tree
```

```
In [2]: df=pd.read_csv(r"C:\Users\user\Downloads\csvs_per_year\csvs_per_year\madrid_2009\madrid_2009.csv")
df
```

Out[2]:

	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3	PM
0	2009-10-01 01:00:00	NaN	0.27	NaN	NaN	NaN	39.889999	48.150002	NaN	50.680000	18.260000
1	2009-10-01 01:00:00	NaN	0.22	NaN	NaN	NaN	21.230000	24.260000	NaN	55.880001	10.580000
2	2009-10-01 01:00:00	NaN	0.18	NaN	NaN	NaN	31.230000	34.880001	NaN	49.060001	25.190000
3	2009-10-01 01:00:00	0.95	0.33	1.43	2.68	0.25	55.180000	81.360001	1.57	36.669998	26.530000
4	2009-10-01 01:00:00	NaN	0.41	NaN	NaN	0.12	61.349998	76.260002	NaN	38.090000	23.760000
...	...	...	...	...	...	...	...	...	...	...	...
215683	2009-06-01 00:00:00	0.50	0.22	0.39	0.75	0.09	22.000000	24.510000	1.00	82.239998	10.830000
215684	2009-06-01 00:00:00	NaN	0.31	NaN	NaN	NaN	76.110001	101.099998	NaN	41.220001	9.920000
215685	2009-06-01 00:00:00	0.13	NaN	0.86	NaN	0.23	81.050003	99.849998	NaN	24.830000	12.460000
215686	2009-06-01 00:00:00	0.21	NaN	2.96	NaN	0.10	72.419998	82.959999	NaN	NaN	13.030000
215687	2009-06-01 00:00:00	0.37	0.32	0.99	1.36	0.14	54.290001	64.480003	1.06	56.919998	15.360000

215688 rows × 17 columns



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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 215688 entries, 0 to 215687
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        215688 non-null  object
1   BEN         60082 non-null   float64
2   CO          190801 non-null  float64
3   EBE         60081 non-null   float64
4   MXY         24846 non-null   float64
5   NMHC        74748 non-null   float64
6   NO_2        214562 non-null  float64
7   NOx         214565 non-null  float64
8   OXY         24854 non-null   float64
9   O_3         204482 non-null  float64
10  PM10        196331 non-null  float64
11  PM25        55822 non-null   float64
12  PXY         24854 non-null   float64
13  SO_2        212671 non-null  float64
14  TCH         75213 non-null   float64
15  TOL         59920 non-null   float64
16  station     215688 non-null  int64
dtypes: float64(15), int64(1), object(1)
memory usage: 28.0+ MB
```

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

Out[4]:

	date	BEN	CO	EBE	MXY	NMHC	NO_2	NOx	OXY	O_3	PM10
<b>3</b>	2009-10-01 01:00:00	0.95	0.33	1.43	2.68	0.25	55.180000	81.360001	1.57	36.669998	26.530000
<b>20</b>	2009-10-01 01:00:00	0.38	0.32	0.32	0.89	0.01	17.969999	19.240000	1.00	65.870003	10.520000
<b>24</b>	2009-10-01 01:00:00	0.55	0.24	0.65	1.79	0.18	36.619999	43.919998	1.28	48.070000	19.150000
<b>28</b>	2009-10-01 02:00:00	0.65	0.21	1.20	2.04	0.18	37.169998	48.869999	1.21	26.950001	32.200000
<b>45</b>	2009-10-01 02:00:00	0.38	0.30	0.50	1.15	0.00	17.889999	19.299999	1.00	60.009998	12.260000
...	...	...	...	...	...	...	...	...	...	...	...
<b>215659</b>	2009-05-31 23:00:00	0.54	0.27	1.00	0.69	0.09	28.280001	29.490000	0.86	78.750000	15.170000
<b>215663</b>	2009-05-31 23:00:00	0.74	0.35	1.13	1.65	0.15	56.410000	69.870003	1.26	56.799999	11.800000
<b>215667</b>	2009-06-01 00:00:00	0.78	0.29	0.99	1.96	0.04	64.870003	82.629997	1.13	58.000000	12.670000
<b>215683</b>	2009-06-01 00:00:00	0.50	0.22	0.39	0.75	0.09	22.000000	24.510000	1.00	82.239998	10.830000
<b>215687</b>	2009-06-01 00:00:00	0.37	0.32	0.99	1.36	0.14	54.290001	64.480003	1.06	56.919998	15.360000

24717 rows × 17 columns



In [5]: `df.isnull().sum()`

```
Out[5]: date      0
      BEN      0
      CO      0
      EBE      0
      MXY      0
      NMHC     0
      NO_2     0
      NOx      0
      OXY      0
      O_3      0
      PM10     0
      PM25     0
      PXY      0
      SO_2     0
      TCH      0
      TOL      0
      station  0
      dtype: int64
```

In [6]: `df.describe()`

```
Out[6]:
```

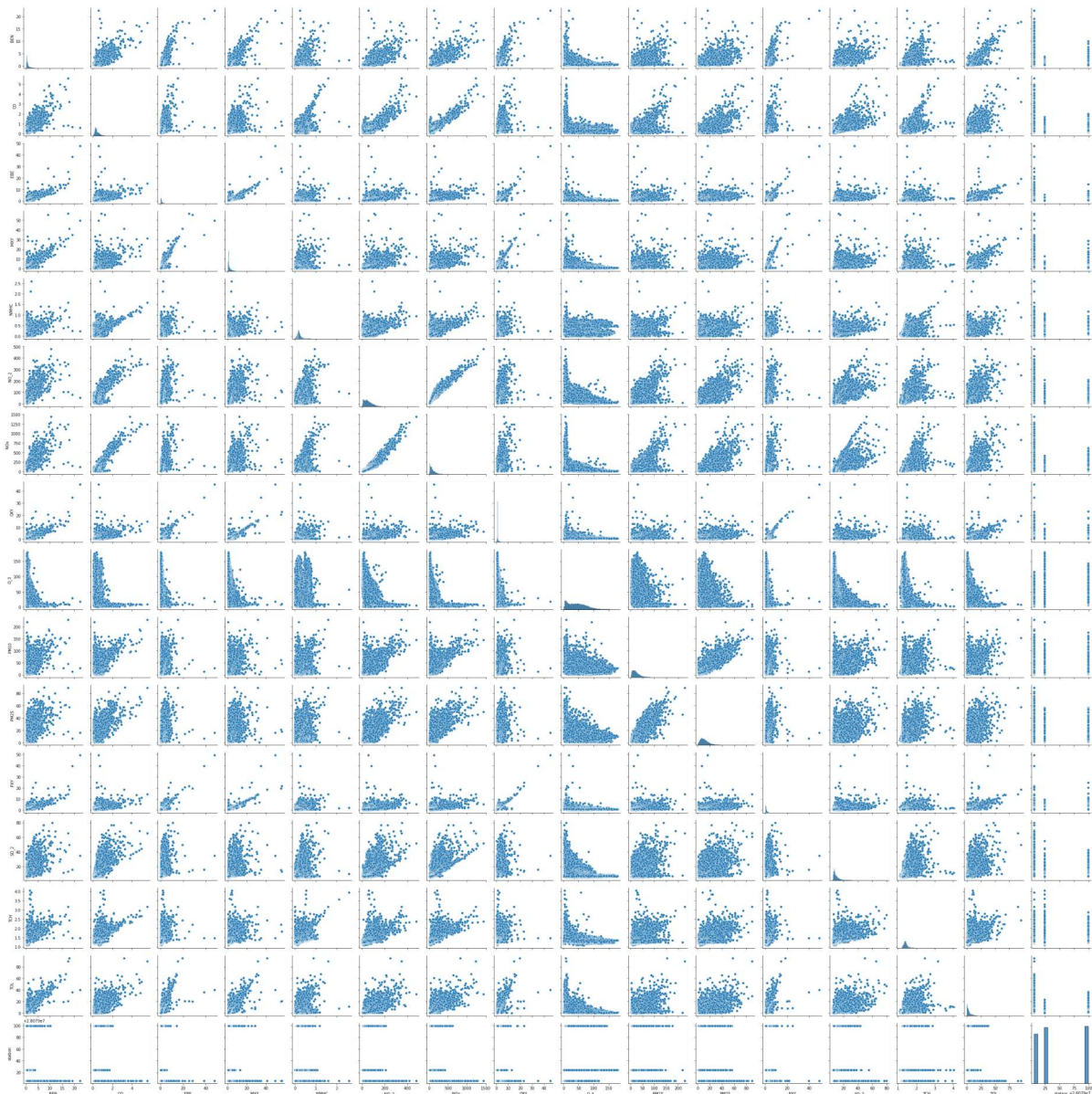
	BEN	CO	EBE	MXY	NMHC	NO_2	
<b>count</b>	24717.000000	24717.000000	24717.000000	24717.000000	24717.000000	24717.000000	247
<b>mean</b>	1.010583	0.448056	1.262430	2.244469	0.219582	55.563929	
<b>std</b>	1.007345	0.291706	1.074768	2.242214	0.141661	38.911677	
<b>min</b>	0.170000	0.060000	0.250000	0.240000	0.000000	0.600000	
<b>25%</b>	0.460000	0.270000	0.720000	0.990000	0.140000	26.510000	
<b>50%</b>	0.670000	0.370000	1.000000	1.490000	0.190000	47.930000	
<b>75%</b>	1.180000	0.570000	1.430000	2.820000	0.260000	76.269997	1
<b>max</b>	22.379999	5.570000	47.669998	56.500000	2.580000	477.399994	14

In [7]: `df.columns`

```
Out[7]: Index(['date', 'BEN', 'CO', 'EBE', 'MXY', 'NMHC', 'NO_2', 'NOx', 'OXY', 'O_3',
              'PM10', 'PM25', 'PXY', 'SO_2', 'TCH', 'TOL', 'station'],
              dtype='object')
```

```
In [8]: sns.pairplot(df)
```

```
Out[8]: <seaborn.axisgrid.PairGrid at 0x1a819cae7f0>
```

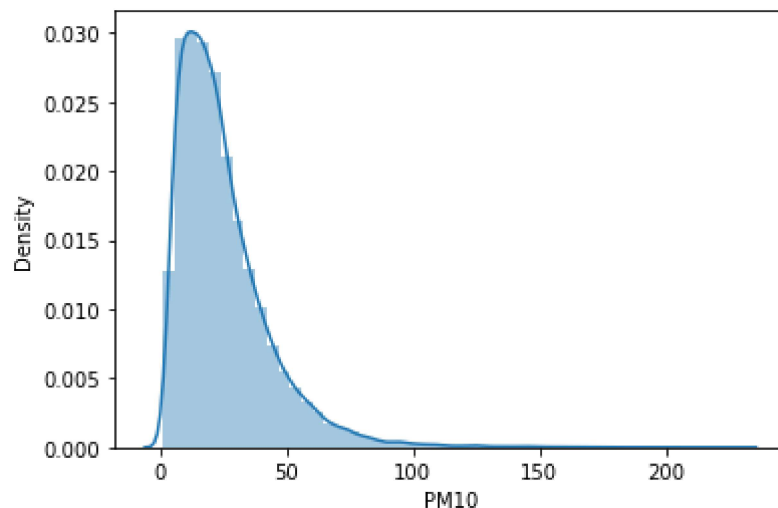


```
In [9]: sns.distplot(df['PM10'])
```

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).

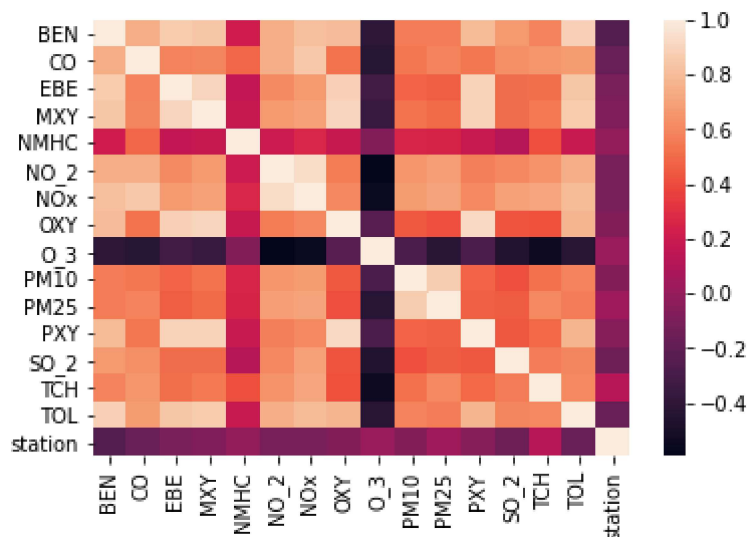
```
warnings.warn(msg, FutureWarning)
```

```
Out[9]: <AxesSubplot:xlabel='PM10', ylabel='Density'>
```



```
In [10]: sns.heatmap(df.corr())
```

```
Out[10]: <AxesSubplot:>
```



```
In [11]: df.loc[df['NMHC']<1, 'NMHC']=0
df.loc[df['NMHC']>1, 'NMHC']=1
df['NMHC']=df['NMHC'].astype(int)
```

C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\indexing.py:1720: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

```
self._setitem_single_column(loc, value, pi)
```

C:\ProgramData\Anaconda3\lib\site-packages\pandas\core\indexing.py:1720: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

```
self._setitem_single_column(loc, value, pi)
```

<ipython-input-11-c5145d14383f>:3: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

```
df['NMHC']=df['NMHC'].astype(int)
```

## LogisticRegression

```
In [12]: x=df[['BEN', 'CO', 'EBE', 'MXY', 'NO_2', 'NOx', 'OXY', 'O_3',
               'PM10', 'PXY', 'SO_2', 'TCH', 'TOL', 'station']]
y=df['NMHC']
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
lgr=LogisticRegression()
lgr.fit(x_train,y_train)
```

Out[12]: LogisticRegression()

```
In [13]: lgr.predict(x_test)
```

Out[13]: array([0, 0, 0, ..., 0, 0, 0])

```
In [14]: lgr.score(x_test,y_test)
```

Out[14]: 0.9989212513484358



```
In [15]: fs=StandardScaler().fit_transform(x)
logr=LogisticRegression()
logr.fit(fs,y)
```

Out[15]: LogisticRegression()

```
In [16]: o=[[1,2,3,4,5,6,7,8,9,10,11,12,13,14]]
prediction=logr.predict(o)
print(prediction)

[1]
```

```
In [17]: logr.classes_
```

Out[17]: array([0, 1])

```
In [18]: logr.predict_proba(o)[0][0]
```

Out[18]: 0.04088401509996609

```
In [19]: logr.predict_proba(o)[0][1]
```

Out[19]: 0.9591159849000339

## LinearRegression

```
In [20]: lr=LinearRegression()
lr.fit(x_train,y_train)
```

Out[20]: LinearRegression()

```
In [21]: print(lr.intercept_)
```

34.33733794159383

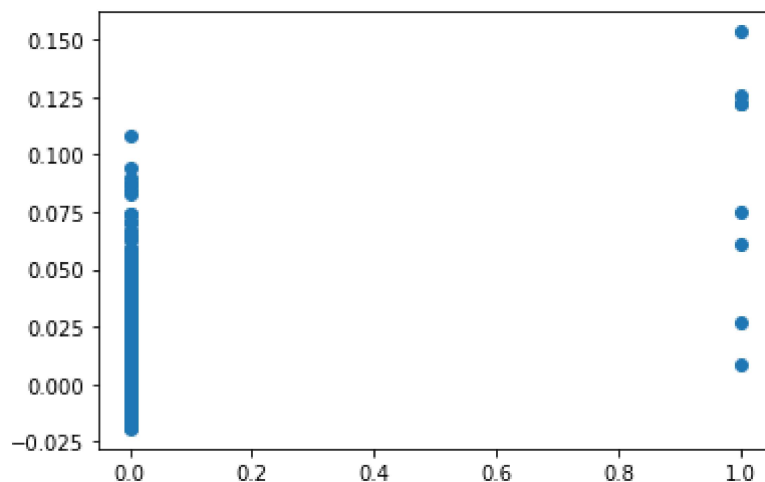
```
In [22]: coeff = pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

Out[22]:

	Co-efficient
<b>BEN</b>	0.002254
<b>CO</b>	0.017944
<b>EBE</b>	-0.002287
<b>MXY</b>	-0.000089
<b>NO_2</b>	-0.000266
<b>NOx</b>	0.000114
<b>OXY</b>	0.002818
<b>O_3</b>	0.000078
<b>PM10</b>	-0.000022
<b>PXY</b>	-0.001675
<b>SO_2</b>	-0.000443
<b>TCH</b>	0.021697
<b>TOL</b>	0.000114
<b>station</b>	-0.000001

```
In [23]: prediction=lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[23]: <matplotlib.collections.PathCollection at 0x1a832608f10>



```
In [24]: print(lr.score(x_test,y_test))
```

0.10116974932462697

## Ridge,Lasso

```
In [25]: rr=Ridge(alpha=10)
rr.fit(x_train,y_train)
```

```
Out[25]: Ridge(alpha=10)
```

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

```
Out[26]: 0.10063691567159239
```

```
In [27]: la=Lasso(alpha=10)
la.fit(x_train,y_train)
```

```
Out[27]: Lasso(alpha=10)
```

```
In [28]: la.score(x_test,y_test)
```

```
Out[28]: -4.160849085832474e-05
```

## ElasticNet

```
In [29]: en=ElasticNet()
en.fit(x_train,y_train)
```

```
Out[29]: ElasticNet()
```

```
In [30]: print(en.coef_)
```

```
[ 0.  0.  0.  0.  0.  0.  0. -0.  0.  0.  0.  0.  0. -0.]
```

```
In [31]: print(en.intercept_)
```

```
0.0008670019074041963
```

```
In [32]: print(en.predict(x_train))
```

```
[0.000867 0.000867 0.000867 ... 0.000867 0.000867 0.000867]
```

```
In [33]: print(en.score(x_train,y_train))
```

```
0.0
```

```
In [34]: print("Mean Absolytre Error:",metrics.mean_absolute_error(y_test,prediction))
```

```
Mean Absolytre Error: 0.005823880933805836
```

```
In [35]: print("Mean Square Error:",metrics.mean_squared_error(y_test,prediction))
```

```
Mean Square Error: 0.0009685659533489445
```

```
In [36]: print("Root Mean Square Error:",np.sqrt(metrics.mean_absolute_error(y_test,pred)))
```

Root Mean Square Error: 0.0763143560138316

## RandomForest

```
In [37]: rfc=RandomForestClassifier()  
rfc.fit(x_train,y_train)
```

Out[37]: RandomForestClassifier()

```
In [38]: parameters={'max_depth':[1,2,3,4,5],  
                    'min_samples_leaf':[5,10,15,20,25],  
                    'n_estimators':[10,20,30,40,50]}
```

```
In [39]: grid_search=GridSearchCV(estimator=rfc,param_grid=parameters,cv=2,scoring="accuracy")  
grid_search.fit(x_train,y_train)
```

Out[39]: GridSearchCV(cv=2, estimator=RandomForestClassifier(),  
 param\_grid={'max\_depth': [1, 2, 3, 4, 5],  
 'min\_samples\_leaf': [5, 10, 15, 20, 25],  
 'n\_estimators': [10, 20, 30, 40, 50]},  
 scoring='accuracy')

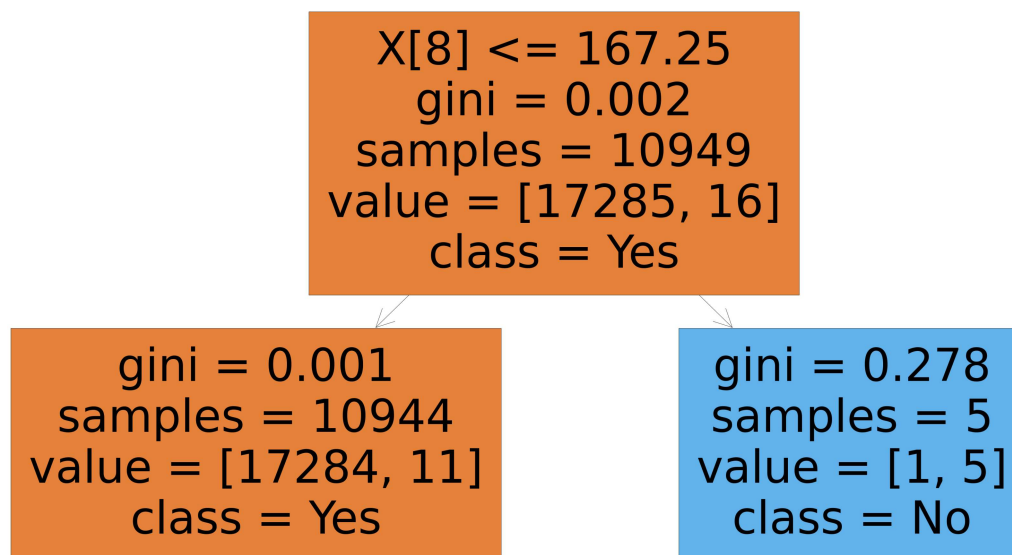
```
In [40]: grid_search.best_score_
```

Out[40]: 0.9991908048987621

```
In [41]: rfc_best=grid_search.best_estimator_
```

```
In [42]: plt.figure(figsize=(80,40))  
plot_tree(rfc_best.estimators_[5],class_names=['Yes','No'],'Yes','No'],filled=True)
```

```
Out[42]: [Text(2232.0, 1630.8000000000002, 'X[8] <= 167.25\nngini = 0.002\nnsamples = 10949\nnvalue = [17285, 16]\nnclass = Yes'),  
Text(1116.0, 543.5999999999999, 'gini = 0.001\nnsamples = 10944\nnvalue = [17284, 11]\nnclass = Yes'),  
Text(3348.0, 543.5999999999999, 'gini = 0.278\nnsamples = 5\nnvalue = [1, 5]\nnclass = No')]
```



Best model: RandomForest

In [ ]: