

```
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_2008\madrid_2008.csv")
df
```

Out[2]:

	date	BEN	CO	EBE	MXV	NMHC	NO_2	NOx	OXY	O_3	P
0	2008-06-01 01:00:00	NaN	0.47	NaN	NaN	NaN	83.089996	120.699997	NaN	16.990000	16.889
1	2008-06-01 01:00:00	NaN	0.59	NaN	NaN	NaN	94.820000	130.399994	NaN	17.469999	19.040
2	2008-06-01 01:00:00	NaN	0.55	NaN	NaN	NaN	75.919998	104.599998	NaN	13.470000	20.270
3	2008-06-01 01:00:00	NaN	0.36	NaN	NaN	NaN	61.029999	66.559998	NaN	23.110001	10.850
4	2008-06-01 01:00:00	1.68	0.80	1.70	3.01	0.30	105.199997	214.899994	1.61	12.120000	37.160
...
226387	2008-11-01 00:00:00	0.48	0.30	0.57	1.00	0.31	13.050000	14.160000	0.91	57.400002	5.450
226388	2008-11-01 00:00:00	NaN	0.30	NaN	NaN	NaN	41.880001	48.500000	NaN	35.830002	15.020
226389	2008-11-01 00:00:00	0.25	NaN	0.56	NaN	0.11	83.610001	102.199997	NaN	14.130000	17.540
226390	2008-11-01 00:00:00	0.54	NaN	2.70	NaN	0.18	70.639999	81.860001	NaN	NaN	11.910
226391	2008-11-01 00:00:00	0.75	0.36	1.20	2.75	0.16	58.240002	74.239998	1.64	31.910000	12.690

226392 rows × 17 columns

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

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 226392 entries, 0 to 226391
Data columns (total 17 columns):
#   Column      Non-Null Count  Dtype
---  -
0   date        226392 non-null object
1   BEN         67047 non-null float64
2   CO          208109 non-null float64
3   EBE         67044 non-null float64
4   MXY         25867 non-null float64
5   NMHC        85079 non-null float64
6   NO_2        225315 non-null float64
7   NOx         225311 non-null float64
8   OXY         25878 non-null float64
9   O_3         215716 non-null float64
10  PM10        220179 non-null float64
11  PM25        67833 non-null float64
12  PXY         25877 non-null float64
13  SO_2        225405 non-null float64
14  TCH         85107 non-null float64
15  TOL         66940 non-null float64
16  station     226392 non-null int64
dtypes: float64(15), int64(1), object(1)
memory usage: 29.4+ MB
```

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

Out[4]:

	date	BEN	CO	EBE	MXY	NMHC	NO_2	NOx	OXY	O_3	P
4	2008-06-01 01:00:00	1.68	0.80	1.70	3.01	0.30	105.199997	214.899994	1.61	12.120000	37.160
21	2008-06-01 01:00:00	0.32	0.37	1.00	0.39	0.33	21.580000	22.180000	1.00	35.770000	7.900
25	2008-06-01 01:00:00	0.73	0.39	1.04	1.70	0.18	64.839996	86.709999	1.31	23.379999	14.760
30	2008-06-01 02:00:00	1.95	0.51	1.98	3.77	0.24	79.750000	143.399994	2.03	18.090000	31.139
47	2008-06-01 02:00:00	0.36	0.39	0.39	0.50	0.34	26.790001	27.389999	1.00	33.029999	7.620
...
226362	2008-10-31 23:00:00	0.47	0.35	0.65	1.00	0.33	22.480000	25.020000	1.00	33.509998	10.200
226366	2008-10-31 23:00:00	0.92	0.46	1.21	2.75	0.19	78.440002	106.199997	1.70	18.320000	14.140
226371	2008-11-01 00:00:00	1.83	0.53	2.22	4.51	0.17	93.260002	158.399994	2.38	18.770000	20.750
226387	2008-11-01 00:00:00	0.48	0.30	0.57	1.00	0.31	13.050000	14.160000	0.91	57.400002	5.450
226391	2008-11-01 00:00:00	0.75	0.36	1.20	2.75	0.16	58.240002	74.239998	1.64	31.910000	12.690

25631 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]:
```

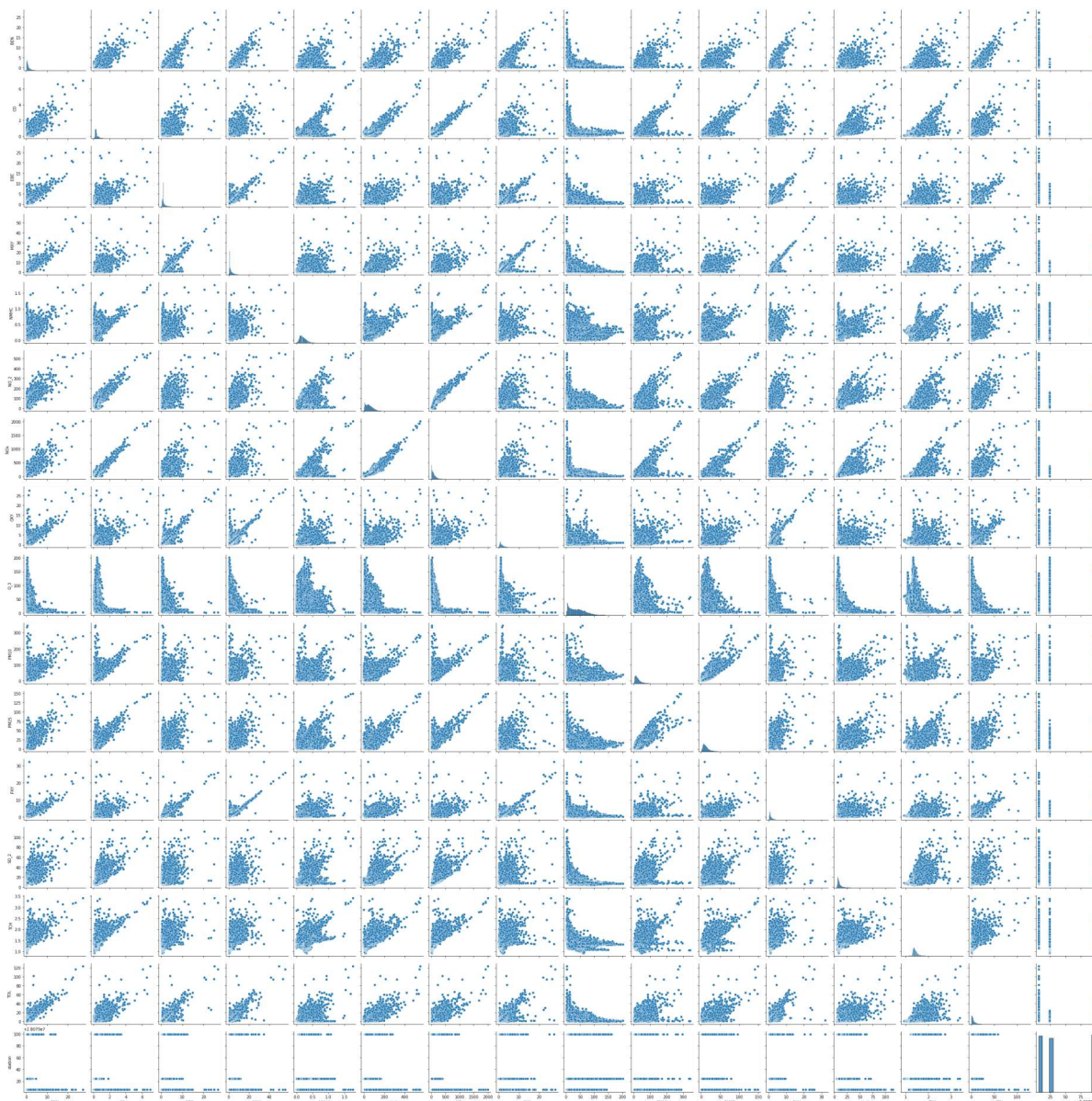
	BEN	CO	EBE	MXY	NMHC	NO_2	
count	25631.000000	25631.000000	25631.000000	25631.000000	25631.000000	25631.000000	256
mean	1.090541	0.440632	1.352355	2.446045	0.213323	54.225261	
std	1.146461	0.317853	1.118191	2.390023	0.123409	38.164647	1
min	0.100000	0.060000	0.170000	0.240000	0.000000	0.240000	
25%	0.430000	0.260000	0.740000	1.000000	0.130000	25.719999	
50%	0.750000	0.350000	1.000000	1.620000	0.190000	48.000000	
75%	1.320000	0.510000	1.580000	3.105000	0.270000	74.924999	1
max	27.230000	7.030000	26.740000	55.889999	1.760000	554.900024	20

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 0x2843983dfd0>
```

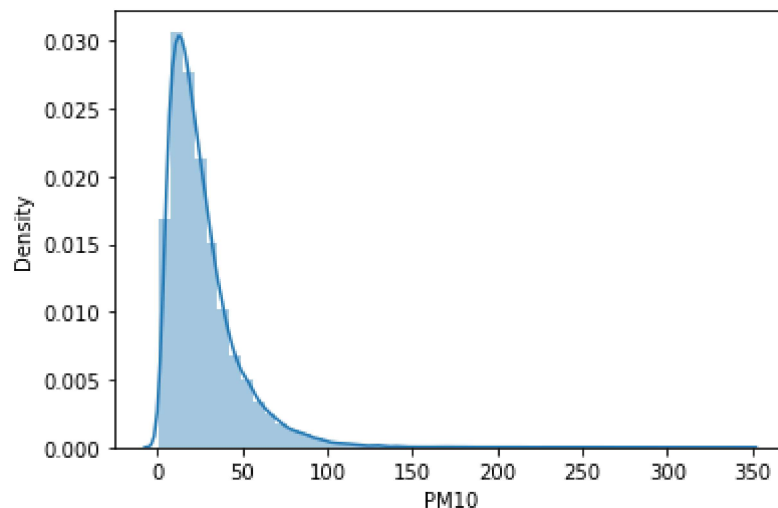


```
In [9]: sns.distplot(df['PM10'])
```

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).

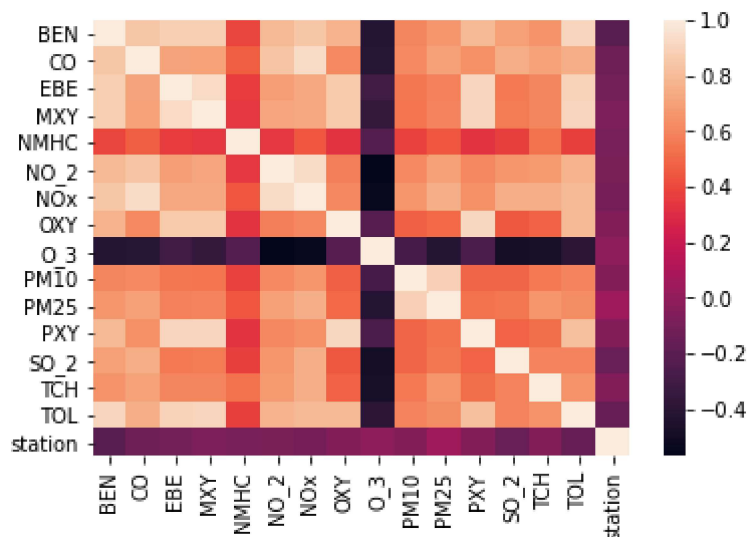
```
warnings.warn(msg, FutureWarning)
```

```
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)

```
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)

```
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)

```
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]: 0.999479843953186


```
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)

[0]
```

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

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

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

```
Out[18]: 0.9100270731923569
```

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

```
Out[19]: 0.08997292680764313
```

LinearRegression

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

```
Out[20]: LinearRegression()
```

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

```
-412.1100134283094
```

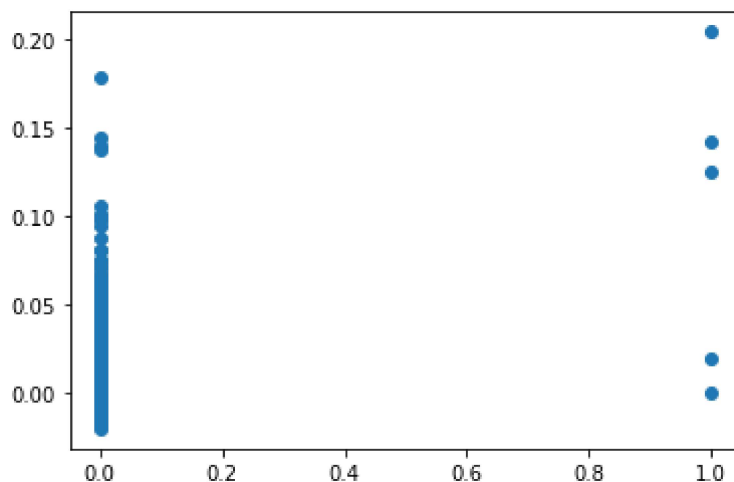
```
In [22]: coeff = pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
coeff
```

Out[22]:

	Co-efficient
BEN	0.004979
CO	-0.003074
EBE	-0.001586
MXY	-0.001169
NO_2	-0.000367
NOx	0.000203
OXY	0.000539
O_3	0.000084
PM10	0.000081
PXY	0.003170
SO_2	-0.000365
TCH	0.007760
TOL	-0.000487
station	0.000015

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

Out[23]: <matplotlib.collections.PathCollection at 0x2845159d8e0>



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

0.04642909562319342

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]: 0.04653308357451391
```

```
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]: -0.0006142951697187815
```

ElasticNet

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

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

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

```
[ 0.00000000e+00  0.00000000e+00  0.00000000e+00  0.00000000e+00
 0.00000000e+00  2.35965016e-05  0.00000000e+00  0.00000000e+00
 0.00000000e+00  0.00000000e+00  0.00000000e+00  0.00000000e+00
 0.00000000e+00 -0.00000000e+00]
```

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

```
-0.0010352759043499716
```

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

```
[0.00354716 0.00824287 0.00375009 ... 0.0004841 0.00118209 0.00570389]
```

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

```
0.022993895120210106
```

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

```
Mean Absolytre Error: 0.005522805540824685
```

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

```
Mean Square Error: 0.0006196039644257703
```

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

Root Mean Square Error: 0.07431558074068105

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')

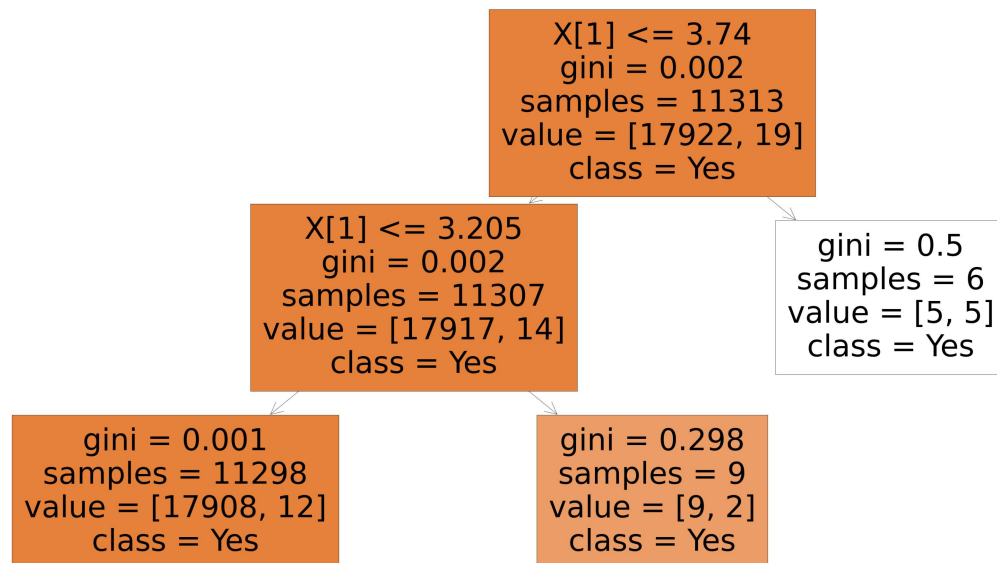
```
In [40]: grid_search.best_score_
```

Out[40]: 0.9989967176534522

```
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(2678.3999999999996, 1812.0, 'X[1] <= 3.74\ngini = 0.002\nsamples = 11313\nvalue = [17922, 19]\nnclass = Yes'),
Text(1785.6, 1087.2, 'X[1] <= 3.205\ngini = 0.002\nsamples = 11307\nvalue = [17917, 14]\nnclass = Yes'),
Text(892.8, 362.39999999999986, 'gini = 0.001\nsamples = 11298\nvalue = [17908, 12]\nnclass = Yes'),
Text(2678.3999999999996, 362.39999999999986, 'gini = 0.298\nsamples = 9\nvalue = [9, 2]\nnclass = Yes'),
Text(3571.2, 1087.2, 'gini = 0.5\nsamples = 6\nvalue = [5, 5]\nnclass = Yes')]
```



Best model: RandomForest

In []: