

Assignment 4

Charvi Upreti

charvi.upreti2021@vitstudent.ac.in (<mailto:charvi.upreti2021@vitstudent.ac.in>)

21BCE1440

ASSIGNMENT 4

Project Title:

Grapes to Greatness: Machine Learning in Wine Quality Prediction

Description:

Predicting wine quality using machine learning is a common and valuable application in the field of data science and analytics. Wine quality prediction involves building a model that can assess and predict the quality of a wine based on various input features, such as chemical composition, sensory characteristics, and environmental factors.

The two datasets are related to red and white variants of the Portuguese "Vinho Verde" wine. For more details, consult the reference [Cortez et al., 2009]. Due to privacy and logistic issues, only physicochemical (inputs) and sensory (the output) variables are available (e.g. there is no data about grape types, wine brand, wine selling price, etc.).

These datasets can be viewed as classification or regression tasks. The classes are ordered and not balanced (e.g. there are much more normal wines than excellent or poor ones).

Dataset: [link](#)

Task:

- Load the Dataset
- Data preprocessing including visualization
- Machine Learning Model building
- Evaluate the model
- Test with random observation

Importing required libraries

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 from sklearn import metrics
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import MinMaxScaler
6 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
7 from sklearn.linear_model import LinearRegression, Ridge, Lasso, LogisticRegression
8 from sklearn.tree import DecisionTreeClassifier, export_graphviz
9 from IPython.display import Image
10 import pydotplus
11 from io import StringIO
12 from sklearn.ensemble import RandomForestClassifier
13 import seaborn as sns
14 import matplotlib.pyplot as plt
15 from matplotlib import rcParams
df.head().T
```

Out[4]:

Loading the data

	0	1	2	3	4
fixed acidity	7.4000	7.8000	7.800	11.200	7.4000
volatile acidity	0.7000	0.8800	0.760	0.280	0.7000
citric acid	0.0000	0.0000	0.040	0.560	0.0000
chlorides	0.0760	0.0980	0.092	0.075	0.0760
free sulfur dioxide	11.0000	25.0000	15.000	17.000	11.0000
total sulfur dioxide	34.0000	67.0000	54.000	60.000	34.0000

```
In [2]: 1 df=pd.read_csv("C:/Users/Charvi Upreti/Desktop/Assignments/Assignment 4/winequality-red.csv")
```

```
Out[3]: (1599, 12)
```

pH	3.5100	3.2000	3.260	3.160	3.5100
sulphates	0.5600	0.6800	0.650	0.580	0.5600
alcohol	9.4000	9.8000	9.800	9.800	9.4000
quality	5.0000	5.0000	5.000	6.000	5.0000


```
citric acid    1599.0   0.270976   0.194801   0.00000   0.0900   0.26000   0.420000  1.00000
```

```
Checking for null values
```

```
chlorides     1599.0   0.087467   0.047065   0.01200   0.0700   0.07900   0.090000  0.61100
```

```
In [7]: 1 free sulfur dioxide sum()
2 #no dioxides values
```

```
Out[7]: fixed acidity 1599.0 46.467792 32.895324 6.00000 22.0000 38.00000 62.000000 289.00000
volatile acidity      0
citric acid           0.990747 0.001887 0.99007 0.9956 0.99675 0.997835 1.00369
residual sugar        0.310113 0.154386 2.74000 3.2100 3.31000 3.400000 4.01000
chlorides              0
free sulfur dioxide  0.658049 0.169507 0.33000 0.5500 0.62000 0.730000 2.00000
total sulfur dioxide 0.10422983 1.065668 8.40000 9.5000 10.20000 11.100000 14.90000
density                0
alcohol                 0
pH                      0
sulphates               0
alcohol                  0
quality                  0
dtype: int64
```

```
In [8]: 1 df.quality.value_counts()
```

```
Out[8]: 5    681
6    638
7    199
4    53
8    18
3    10
Name: quality, dtype: int64
```

No categorical column to perform encoding

```
In [9]: 1 df.dtypes
```

```
Out[9]: fixed acidity          float64
volatile acidity         float64
citric acid              float64
residual sugar            float64
chlorides                 float64
free sulfur dioxide       float64
total sulfur dioxide      float64
density                   float64
pH                        float64
sulphates                 float64
alcohol                   float64
quality                   int64
dtype: object
```

```
In [10]: 1 from sklearn.preprocessing import LabelEncoder
2 le=LabelEncoder()
3 #no need here
```

Duplicates

```
In [11]: 1 df.duplicated()
```

```
Out[11]: 0      False
1      False
2      False
3      False
Index
```

```
In [14]: 4 1 df.inplace=True
```

```
...  
Out[14]: RangeIndex(start=0, stop=1599, step=1)  
1595      False  
1596      True  
1597      False  
1598      False  
Length: 1599, dtype: bool
```

```
In [12]: 1 df1=df.drop_duplicates()
```

```
In [13]: 1 df1.shape
2 ## We will not drop the duplicates as its greatly reducing the data set first
```

```
Out[13]: (1359, 12)
```

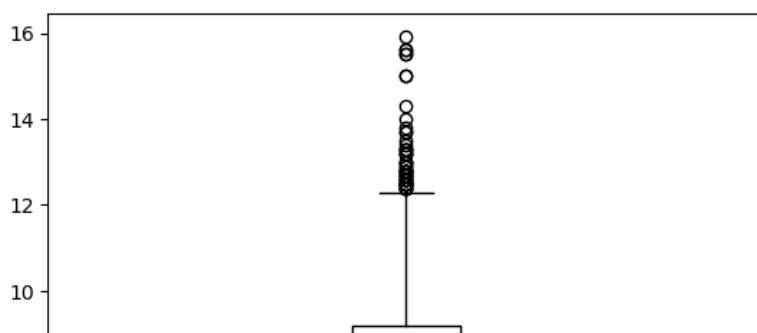
```
Out[11]: 0      False
         1      False
         2      False
         3      False
In [14]: 4 1 df.index
          ...
Out[14]: RangeIndex(start=0, stop=1599, step=1)
        1595    False
        1596     True
        1597    False
        1598    False
Length: 1599, dtype: bool

In [12]: 1 df1 = df.drop_duplicates()

In [13]: 1 df1.shape
2 ## We will not drop the duplicates as its greatly reducing the data set for analysis
Out[13]: (1359, 12)
```

Outlier Detection and replacement

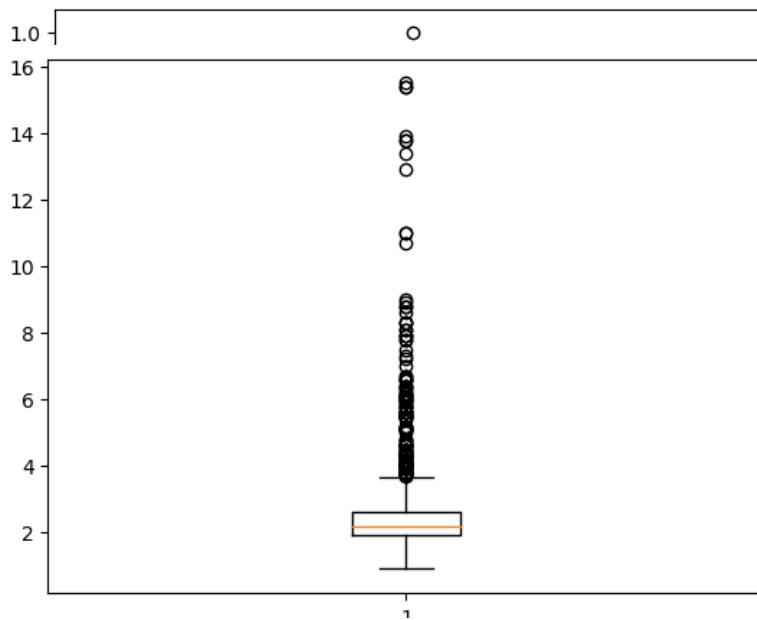
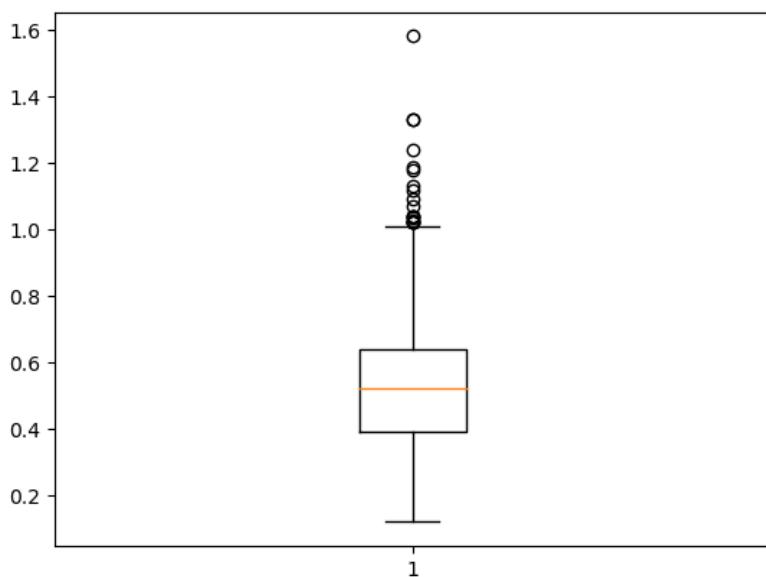
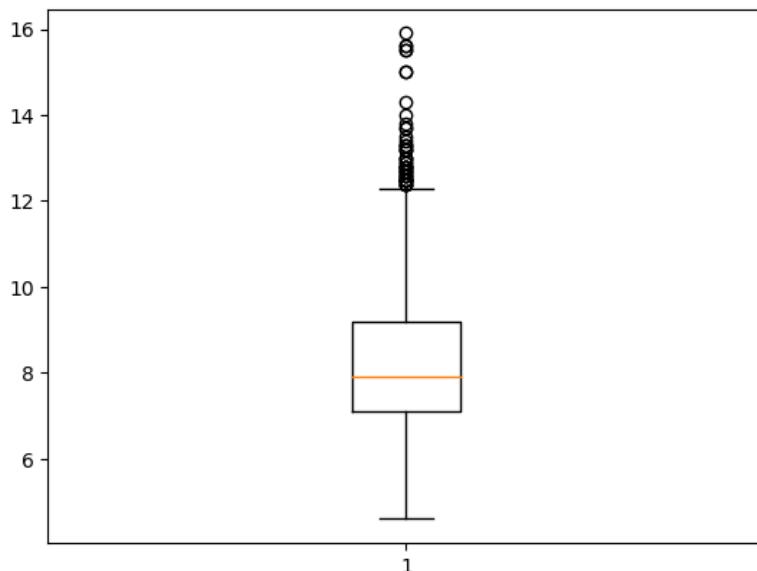
```
In [15]: 1 columns = df.columns
2 for col in columns:
3     if col != 'quality':
4         plt.figure()
5         plt.boxplot(df[col])
```

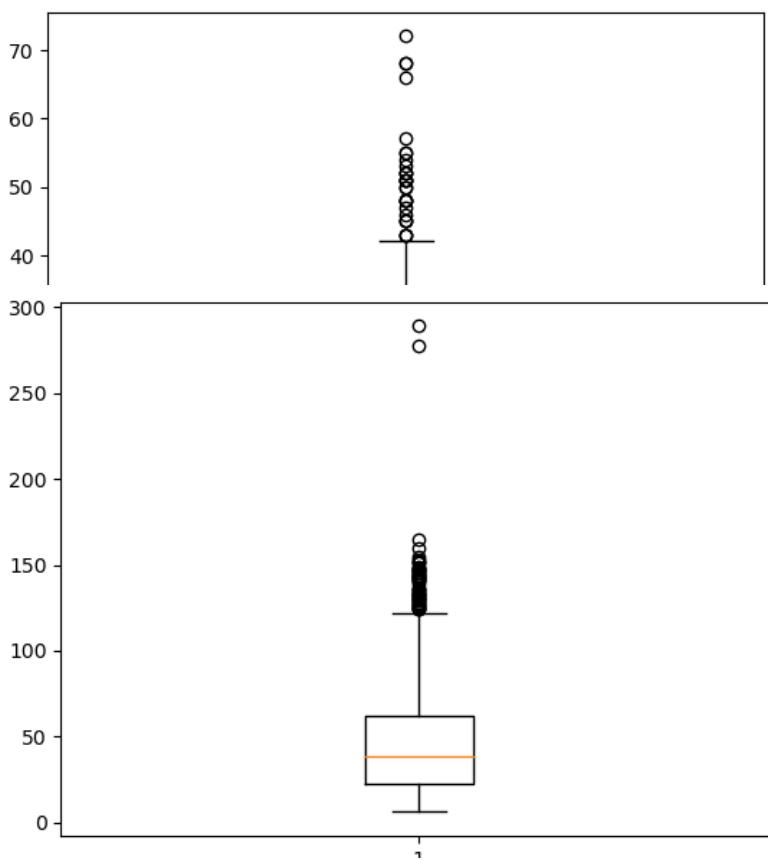
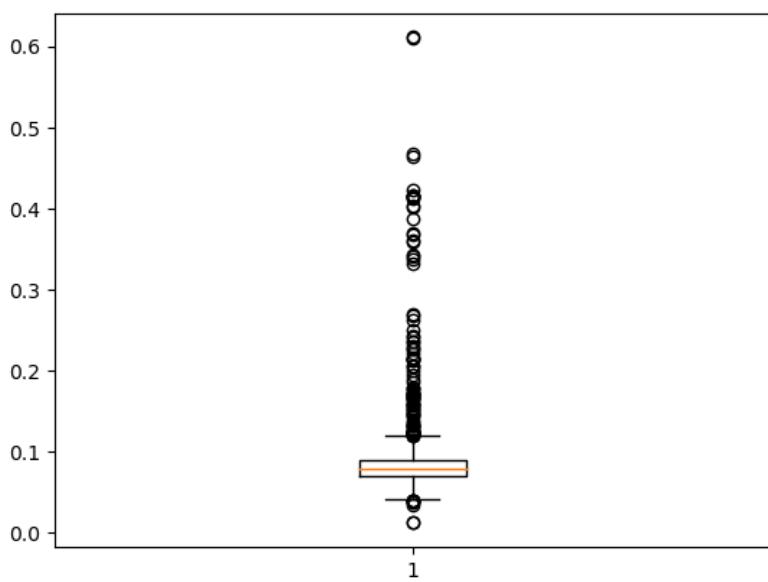
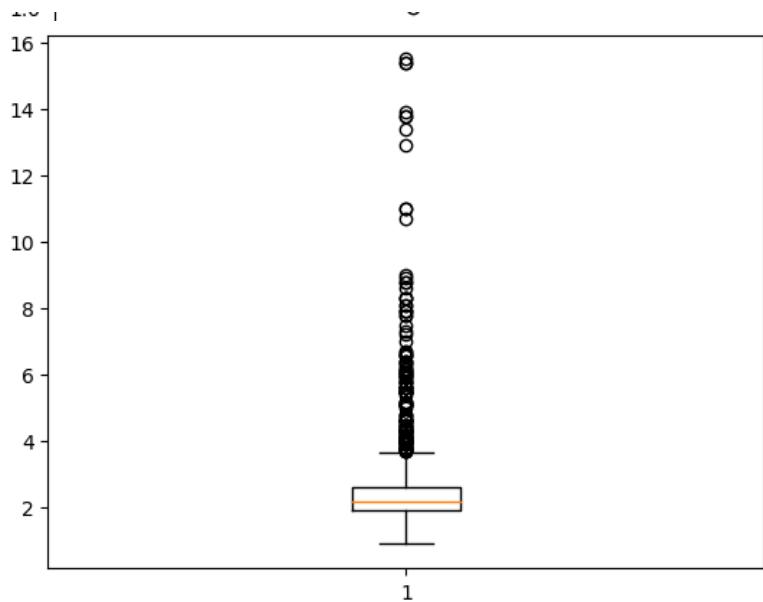


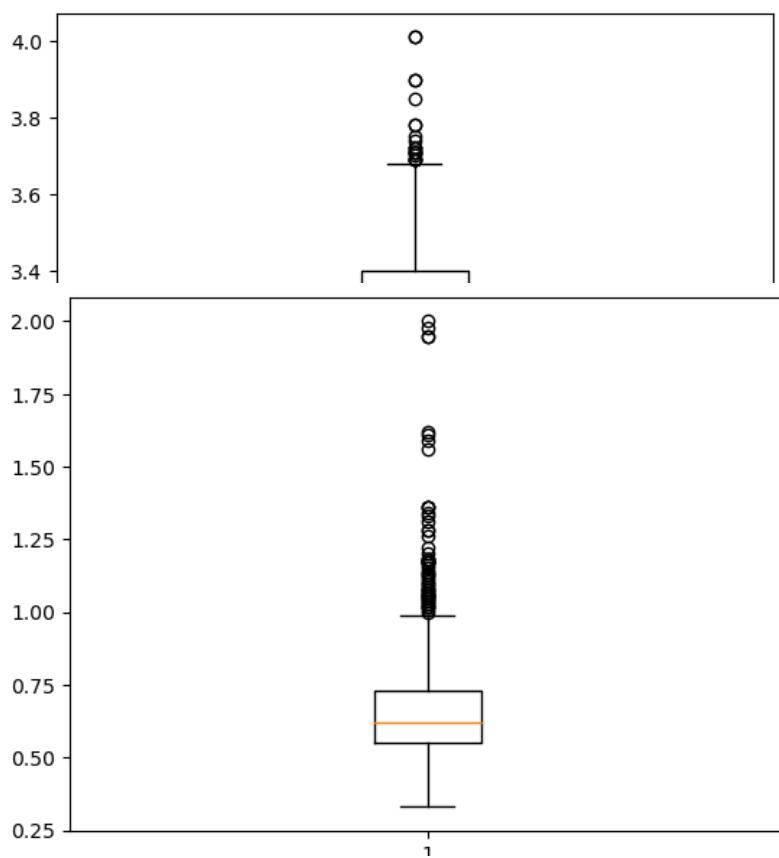
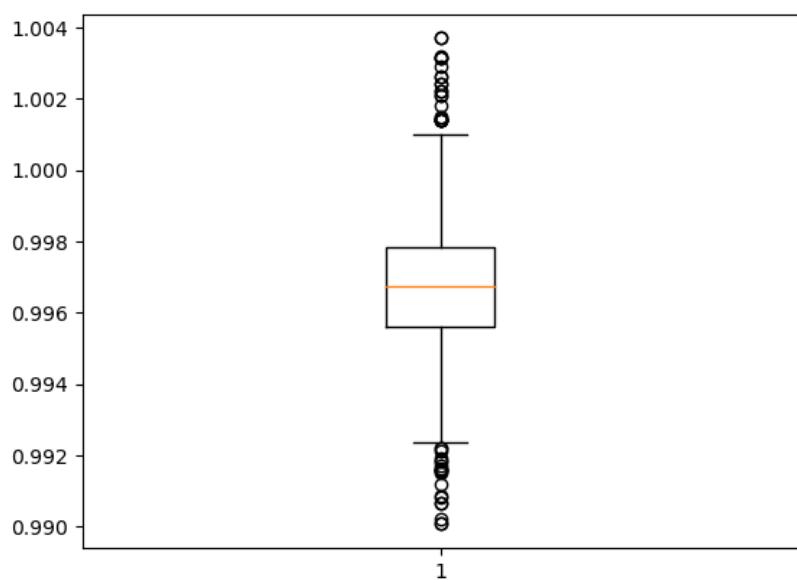
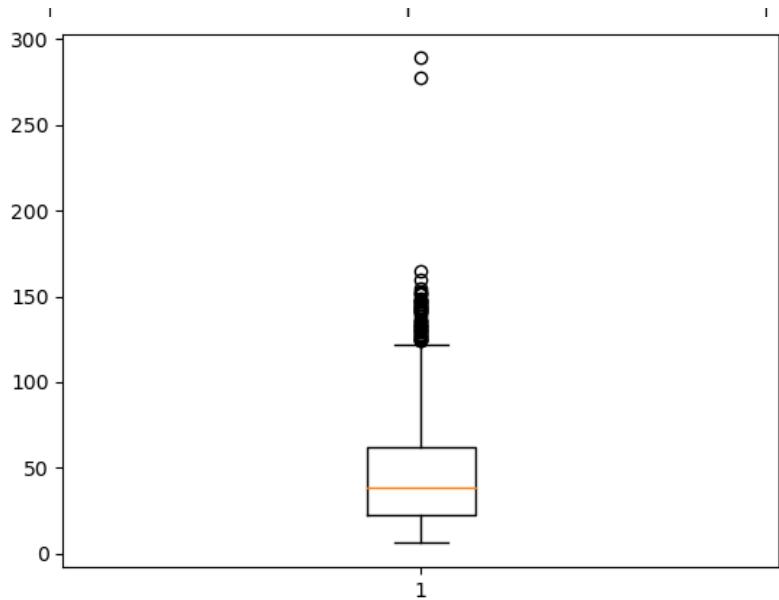
Outlier Detection and replacement

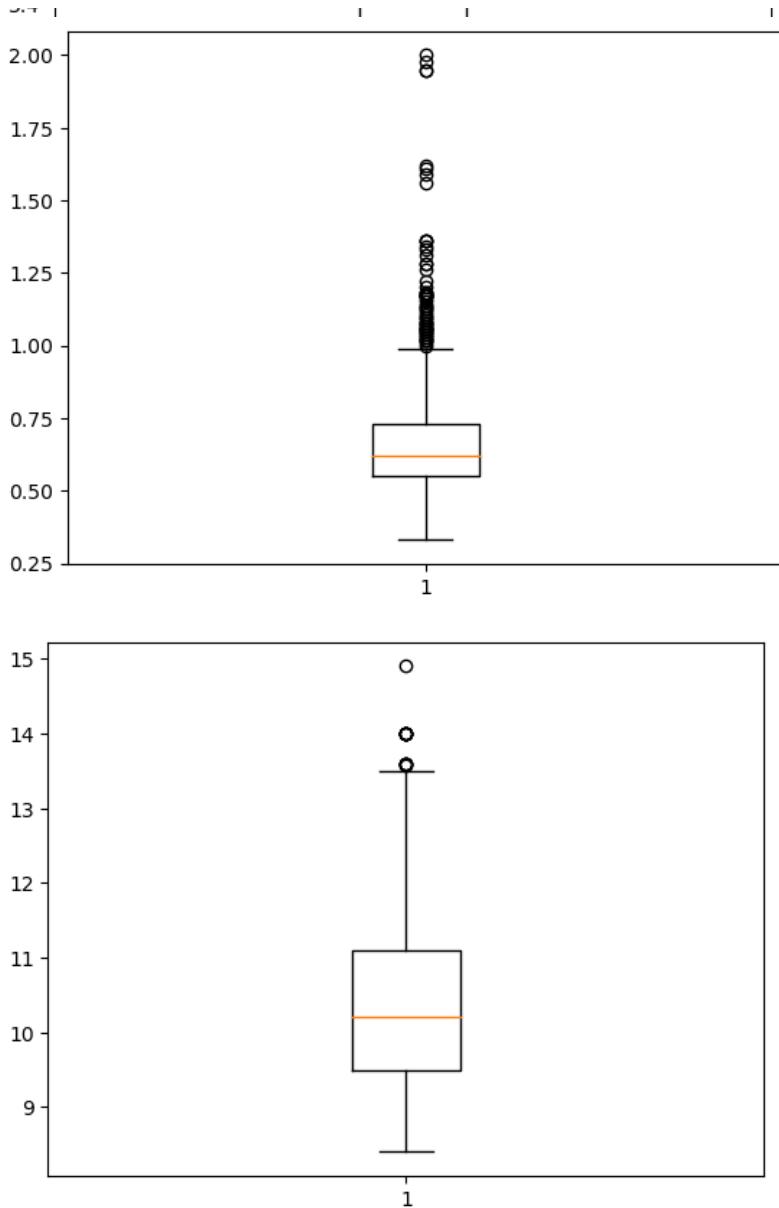
In [15]:

```
1 columns = df.columns
2 for col in columns:
3     if col != 'quality':
4         plt.figure()
5         plt.boxplot(df[col])
```









```
In [16]: 1 #Removing Outliers
2 for col in columns:
3     if col != 'quality':
4         q1 = df[col].quantile(0.25)
5         q3 = df[col].quantile(0.75)
6         IQR = q3 - q1
7         upper_limit = q3 + 1.5 * IQR
8         lower_limit = q1 - 1.5 * IQR # Corrected this line
9         median=df[col].median()
10        df[col] = np.where(df[col] > upper_limit,median , df[col])
11        df[col] = np.where(df[col] < lower_limit, median, df[col])
12
```

Correlation

```
In [17]: 1 df.corr().T
```

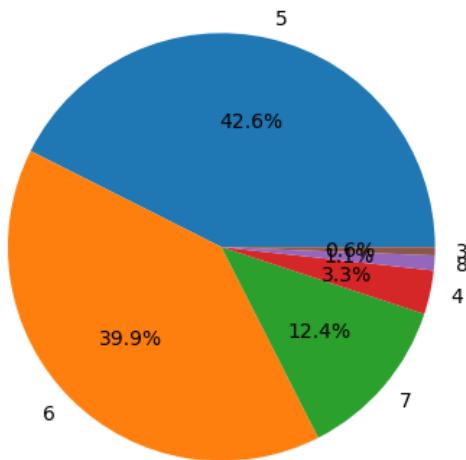
Out[17]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density
fixed acidity	1.000000	-0.258476	0.608682	0.225436	0.227814	-0.145813	-0.113789	0.557213
volatile acidity	-0.258476	1.000000	-0.564257	0.018550	0.138565	0.004853	0.090253	0.000445
citric acid	0.608682	-0.564257	1.000000	0.156708	0.091785	-0.060658	-0.010996	0.338521
residual sugar	0.225436	0.018550	0.156708	1.000000	0.229396	0.040756	0.126300	0.372047
chlorides	0.227814	0.138565	0.091785	0.229396	1.000000	-0.006299	0.099765	0.375573
free sulfur dioxide	-0.145813	0.004853	-0.060658	0.040756	-0.006299	1.000000	0.596640	-0.018912

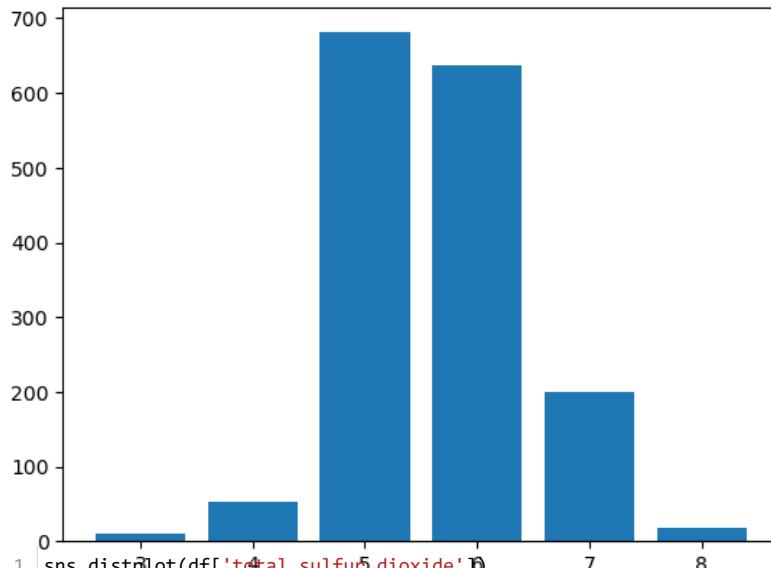
Visualizations

Univariate

```
In [19]: 1 counts = df['quality'].value_counts()
2 plt.figure()
3 plt.pie(counts, labels=counts.index, autopct='%1.1f%%')
4 plt.show()
```



```
In [20]: 1 plt.bar(counts.index, counts)
2 plt.show()
```



```
In [21]: 1 sns.distplot(df['total sulfur dioxide'])
```

C:\Users\Charvi Upreti\AppData\Local\Temp\ipykernel_10484\2700524578.py:1: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
sns.distplot(df['total sulfur dioxide'])
```

```
Out[21]: <Axes: xlabel='total sulfur dioxide', ylabel='Density'>
```



```
In [21]: 1 sns.distplot(df['total sulfur dioxide'])
```

C:\Users\Charvi Upreti\AppData\Local\Temp\ipykernel_10484\2700524578.py:1: UserWarning:

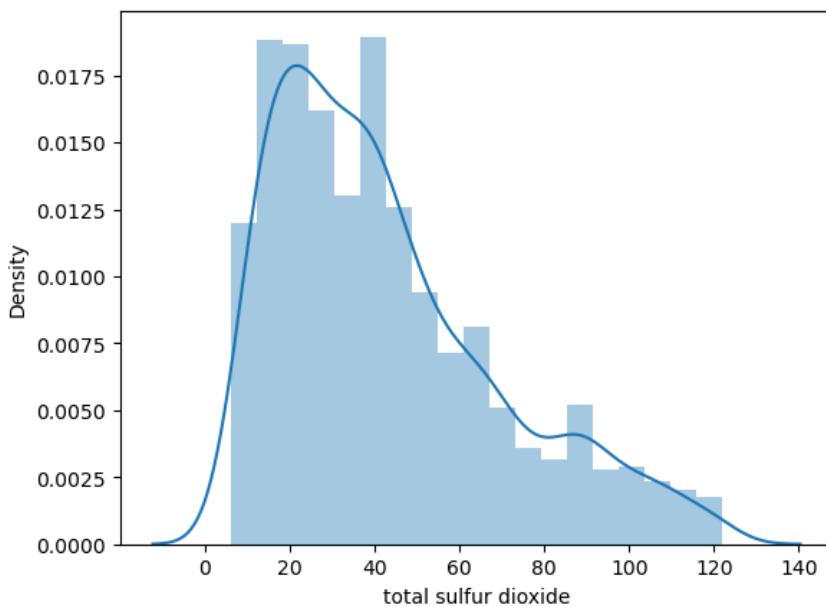
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see
<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
sns.distplot(df['total sulfur dioxide'])
```

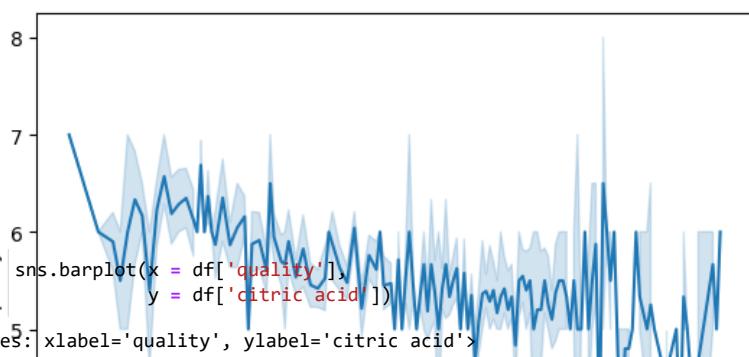
```
Out[21]: <Axes: xlabel='total sulfur dioxide', ylabel='Density'>
```



Bivariate

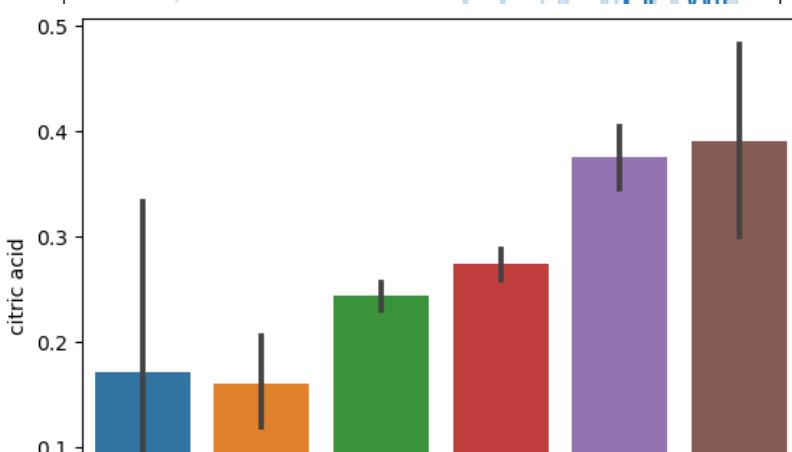
```
In [22]: 1 sns.lineplot(x = df['volatile acidity'],y=df['quality'])
```

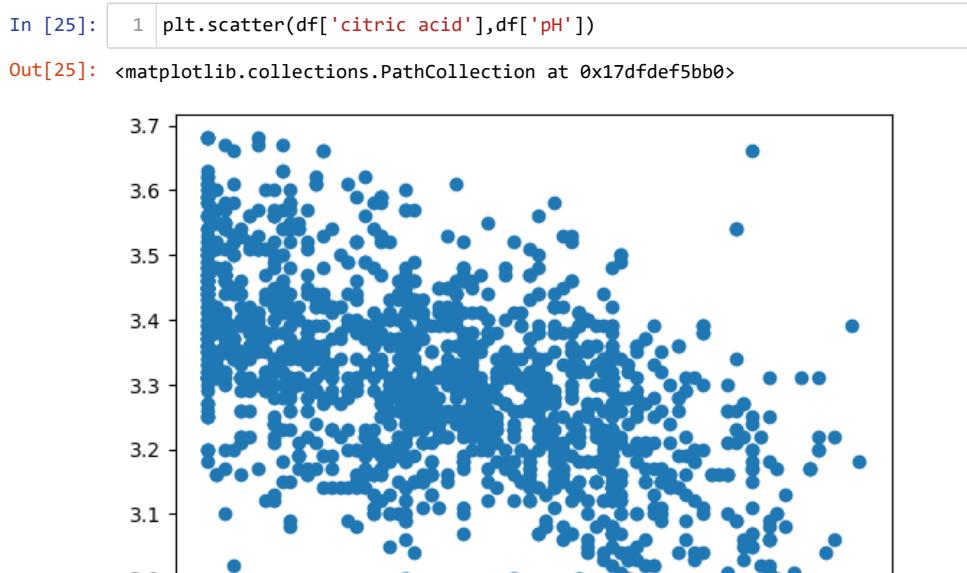
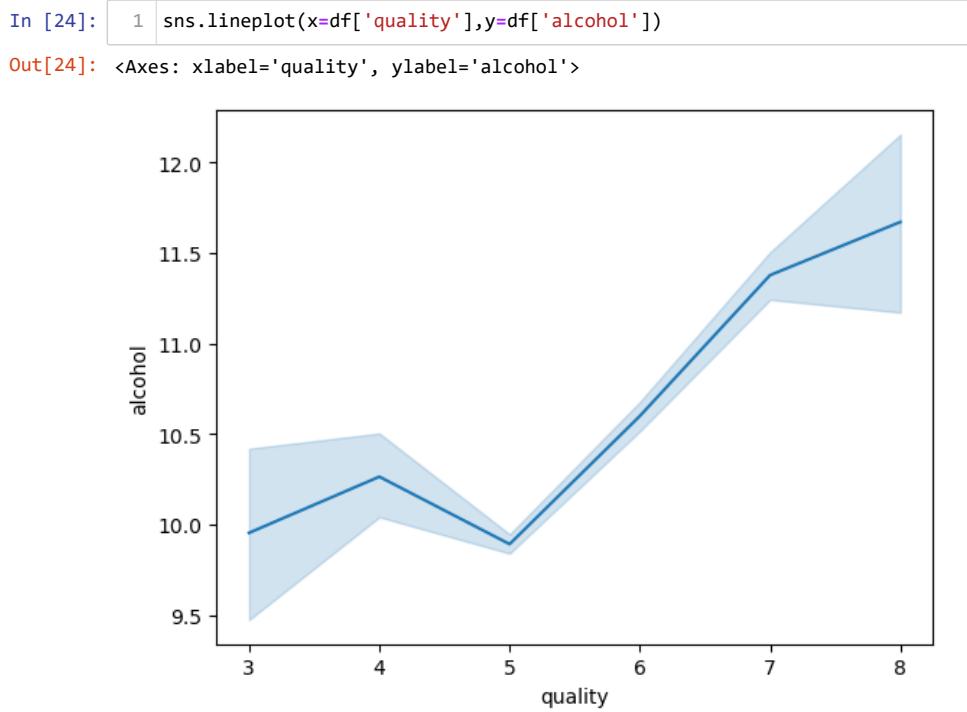
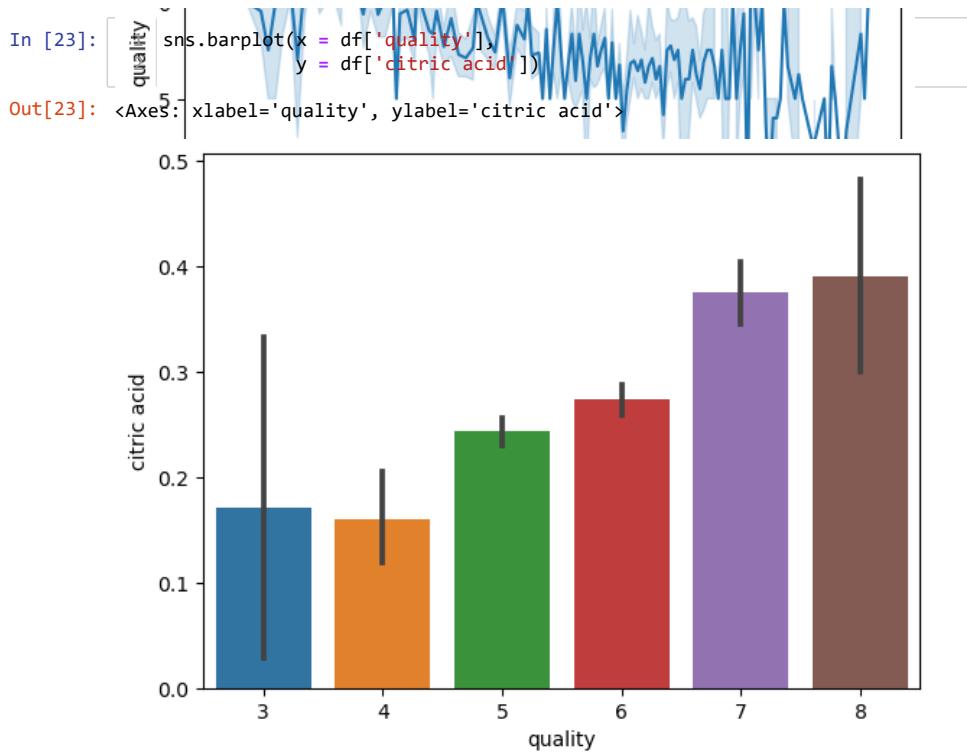
```
Out[22]: <Axes: xlabel='volatile acidity', ylabel='quality'>
```



```
In [23]: 1 sns.barplot(x = df['quality'],  
y = df['citric acid'])
```

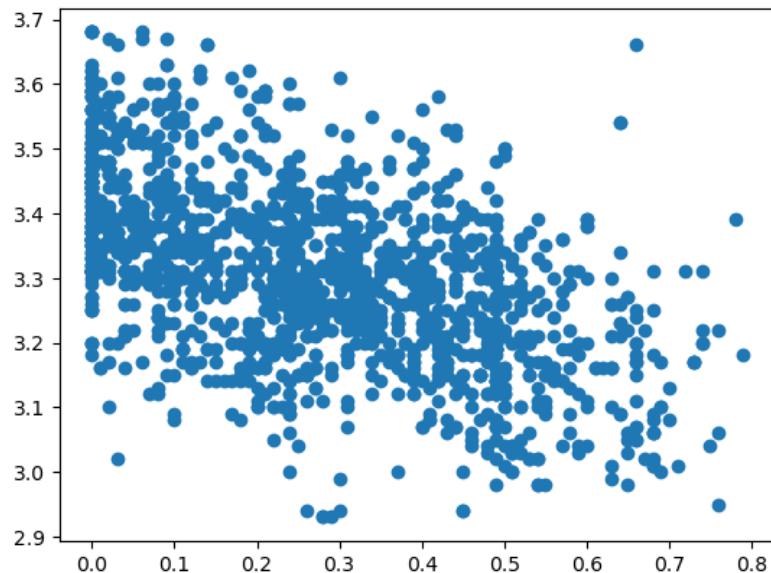
```
Out[23]: <Axes: xlabel='quality', ylabel='citric acid'>
```





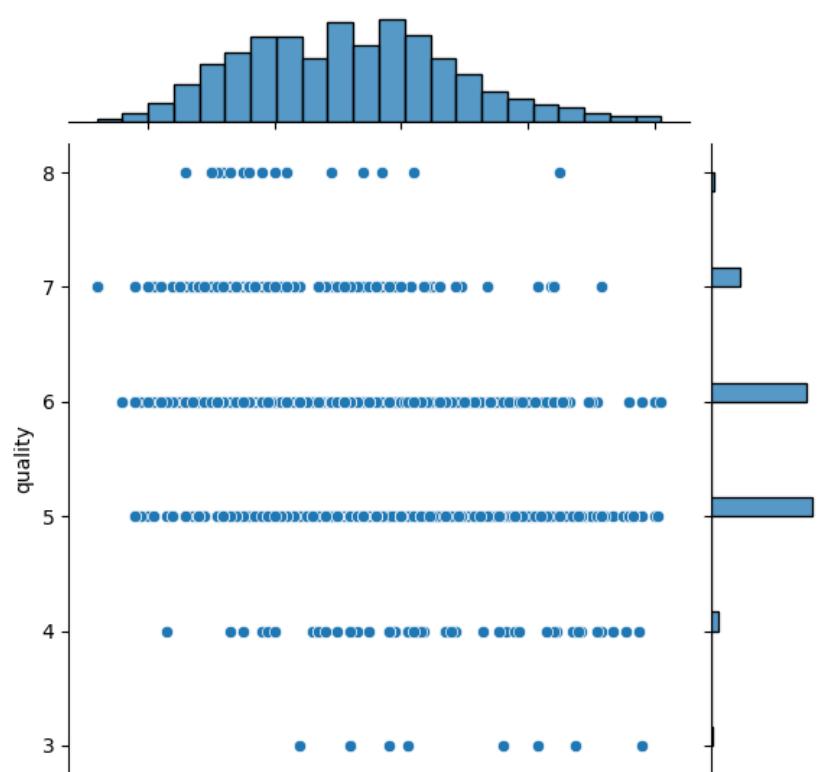
```
In [25]: 1 plt.scatter(df['citric acid'],df['pH'])
```

```
Out[25]: <matplotlib.collections.PathCollection at 0x17dfdef5bb0>
```

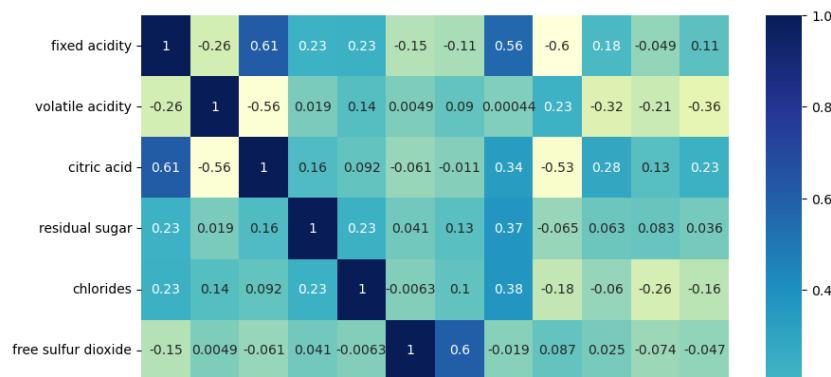


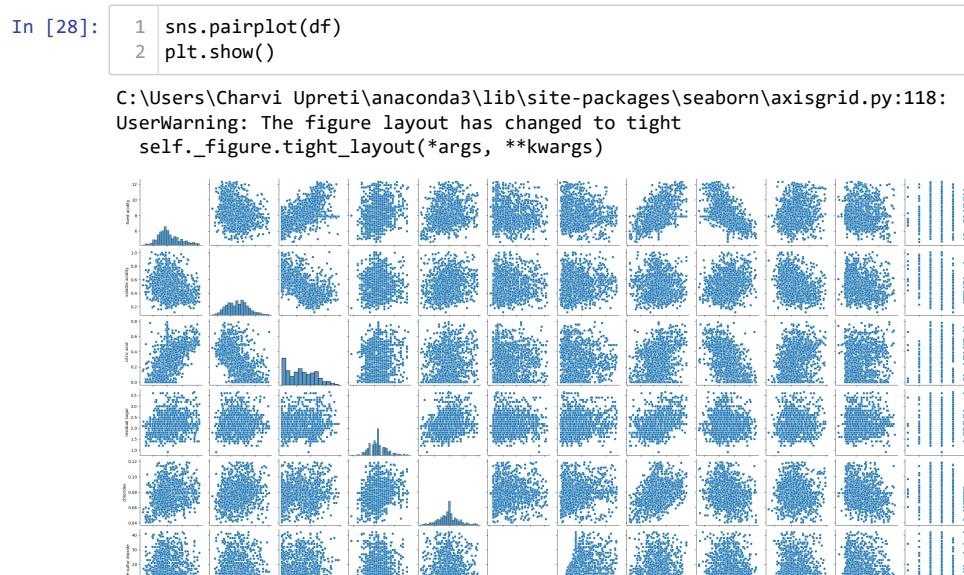
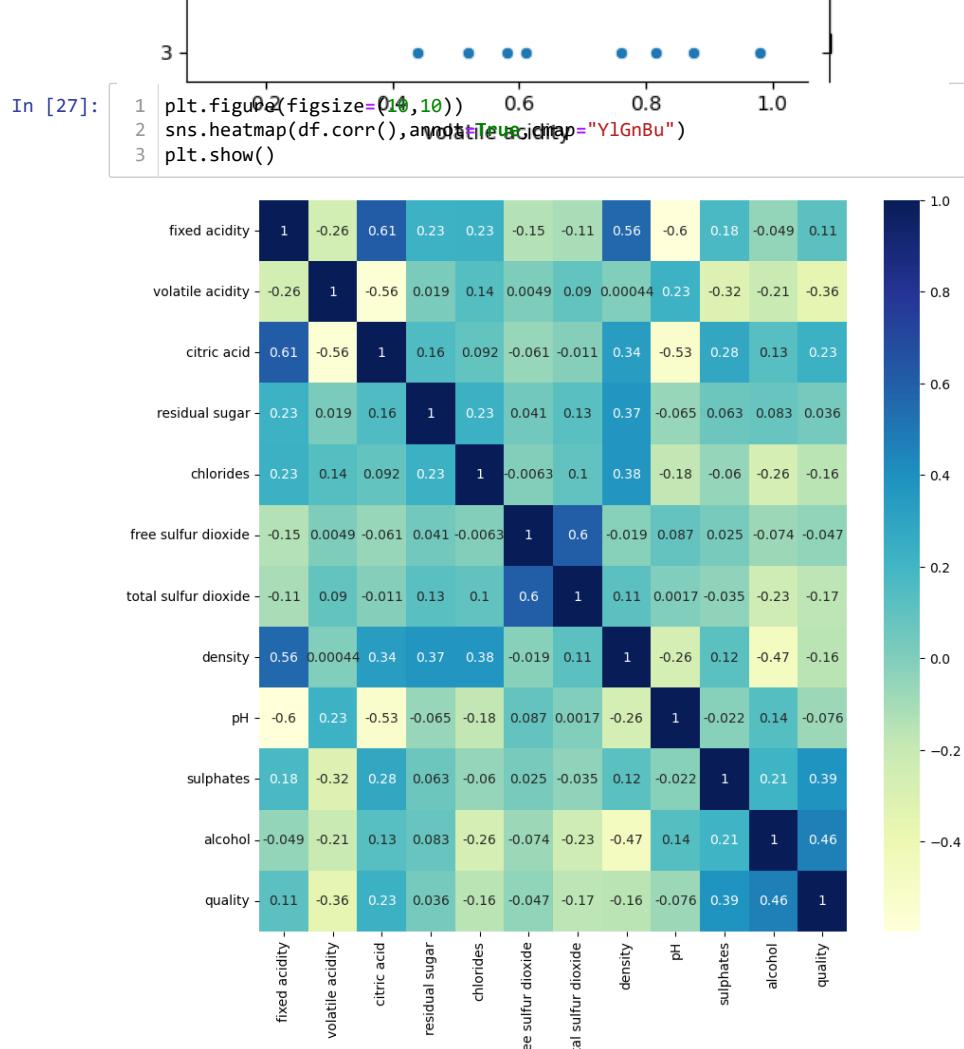
```
In [26]: 1 sns.jointplot(x = df['volatile acidity'],y=df['quality'])
```

```
Out[26]: <seaborn.axisgrid.JointGrid at 0x17d81b56d30>
```



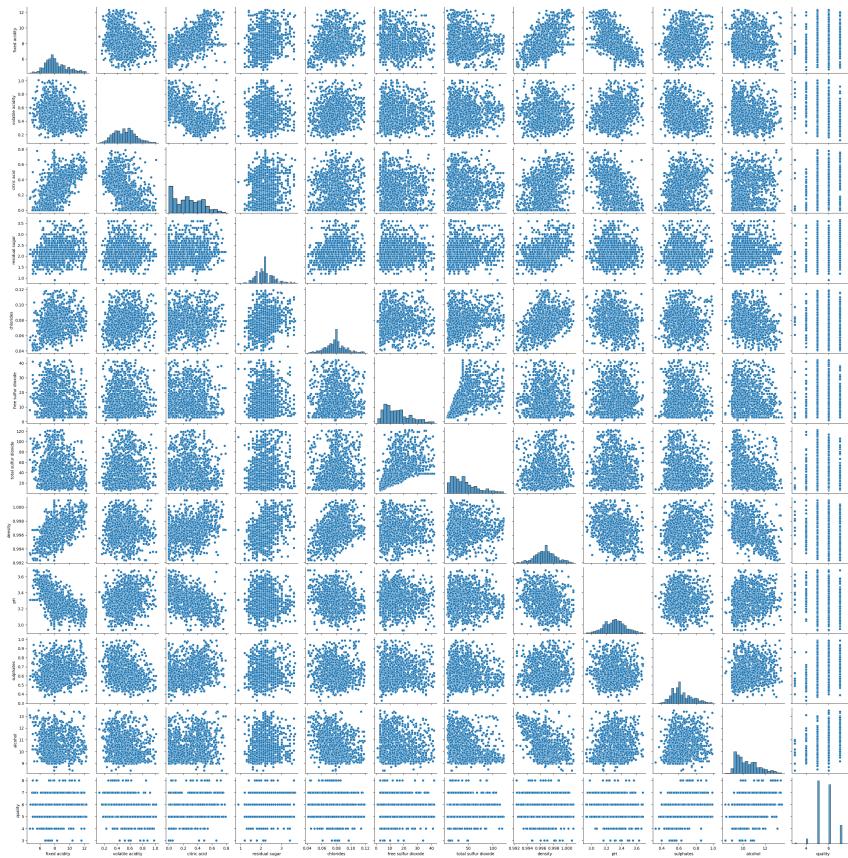
```
In [27]: 1 plt.figure(figsize=(10,10))
2 sns.heatmap(df.corr(),annot=True,cmap="YlGnBu")
3 plt.show()
```





```
In [28]: 1 sns.pairplot(df)
2 plt.show()
```

C:\Users\Charvi Upreti\anaconda3\lib\site-packages\seaborn\axisgrid.py:118:
UserWarning: The figure layout has changed to tight
self._figure.tight_layout(*args, **kwargs)



Splitting

```
In [29]: 1 #x and y split
2 x=df.drop(columns=['quality'],axis=1)
3 print(x.head().T)
4 y=df['quality']
5 y.head()
```

	0	1	2	3	4
fixed acidity	7.4000	7.8000	7.800	11.200	7.4000
volatile acidity	0.7000	0.8800	0.760	0.280	0.7000
citric acid	0.0000	0.0000	0.040	0.560	0.0000
residual sugar	1.9000	2.6000	2.300	1.900	1.9000
chlorides	0.0760	0.0980	0.092	0.075	0.0760
free sulfur dioxide	11.0000	25.0000	15.000	17.000	11.0000
total sulfur dioxide	34.0000	67.0000	54.000	60.000	34.0000
density	0.9978	0.9968	0.997	0.998	0.9978
pH	3.5100	3.2000	3.260	3.160	3.5100
sulphates	0.5600	0.6800	0.650	0.580	0.5600
alcohol	9.4000	9.8000	9.800	9.800	9.4000

Scaling

```
Out[30]: 0 1 scale=MinMaxScaler()
1 2 x_scaled=pd.DataFrame(scale.fit_transform(x),columns=x.columns)
2 3 x_scaled.head().T
3 4
4 5
```

	1	2	3	4	
fixed acidity	0.363636	0.415584	0.415584	0.857143	0.363636
volatile acidity	0.651685	0.853933	0.719101	0.179775	0.651685
citric acid	0.000000	0.000000	0.050633	0.708861	0.000000
residual sugar	0.363636	0.618182	0.509091	0.363636	0.363636
chlorides	0.448718	0.730769	0.653846	0.435897	0.448718
free sulfur dioxide	0.243902	0.585366	0.341463	0.390244	0.243902
total sulfur dioxide	0.241379	0.525862	0.413793	0.465517	0.241379
density	0.630058	0.514451	0.537572	0.653179	0.630058
pH	0.773333	0.360000	0.440000	0.306667	0.773333

```

-----,
pH                      3.5100  3.2000  3.260   3.160  3.5100
sulphates                0.5600  0.6800  0.650   0.580  0.5600
alcohol                  9.4000  9.8000  9.800   9.800  9.4000

```

```

Out[30]: 0 1
1 2
2 3
3 4
4 5

```

```

Name: quality, dtype: float64   1      2      3      4
fixed acidity    0.363636  0.415584  0.415584  0.857143  0.363636
volatile acidity 0.651685  0.853933  0.719101  0.179775  0.651685
citric acid      0.000000  0.000000  0.050633  0.708861  0.000000
residual sugar   0.363636  0.618182  0.509091  0.363636  0.363636
chlorides        0.448718  0.730769  0.653846  0.435897  0.448718
free sulfur dioxide 0.243902  0.585366  0.341463  0.390244  0.243902
total sulfur dioxide 0.241379  0.525862  0.413793  0.465517  0.241379
density          0.630058  0.514451  0.537572  0.653179  0.630058
pH               0.773333  0.360000  0.440000  0.306667  0.773333
sulphates        0.348485  0.530303  0.484848  0.378788  0.348485
alcohol          0.196078  0.274510  0.274510  0.274510  0.196078

```

```

In [31]: 1 x_train,x_test,y_train,y_test=train_test_split(x_scaled,y,test_size=0.2,random_state=42)
2 print(x_train.shape)
3 print(y_train.shape)
4 print(x_test.shape)
5 print(y_test.shape)

(1279, 11)
(1279,)
(320, 11)
(320,)

```

Machine Learning Model Building, Evaluation of the model and Test with random observation

Linear Regression

```

In [32]: 1 lr=LinearRegression()
2 lr.fit(x_train,y_train)

```

```

Out[32]: ▾ LinearRegression
LinearRegression()

```

```

In [33]: 1 #testing
2 print("Test result")
3 y_predict=lr.predict(x_test)
4 #R2 Score
5 print("R2 score",metrics.r2_score(y_test,y_predict))
6 #MSE (Mean Square Error)
7 print("MSE",metrics.mean_squared_error(y_test,y_predict))
8 #RMSE (Root Mean Square Error)
9 print("RMSE",np.sqrt(metrics.mean_squared_error(y_test,y_predict)))
10
11 print('\n')
12
13 #to check if overfitting
14 y_predict1=lr.predict(x_train)
15 #printing to check if overfitting or underfitting
16 #R2 Score
17 print("R2 Score of training",metrics.r2_score(y_train,y_predict1))
18 print("R2 score of training and testing is not much different.")
19
20 print('\n')
21
22 #Checking a random observation:-
23 predicted_values = lr.predict(x_test.iloc[[0]])

```

```

In [33]: 1 #testing
2 print("Test result")
3 y_predict=lr.predict(x_test)
4 #R2 Score
5 print("R2 score",metrics.r2_score(y_test,y_predict))
6 #MSE (Mean Square Error)
7 print("MSE",metrics.mean_squared_error(y_test,y_predict))
8 #RMSE (Root Mean Square Error)
9 print("RMSE",np.sqrt(metrics.mean_squared_error(y_test,y_predict)))
10
11 print('\n')
12
13 #to check if overfitting
14 y_predict1=lr.predict(x_train)
15 #printing to check if overfitting or underfitting
16 #R2 Score
17 print("R2 Score of training",metrics.r2_score(y_train,y_predict1))
18 print("R2 score of training and testing is not much different.")
19
20 print('\n')
21
22 #Checking a random observation:-
23 predicted_values = lr.predict(x_test.iloc[[0]])
24 actual_values = y_test.iloc[[0]]
25 print("Predicted " + str(predicted_values))
26 print("Actual " + str(actual_values))
27
28 column_names = x_test.columns
29
30 # Input data for a single observation
31 input_data = [0.6, 0.5, 0.43, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
32 input_df = pd.DataFrame([input_data], columns=column_names)
33
34 # Make a prediction for the single observation
35 prediction_single_observation = lr.predict(input_df)
36
37 print("Prediction for the single observation:", prediction_single_observation)
38
39 input_df

```

Test result
R2 score 0.3318798167106596
MSE 0.3824727064884877
RMSE 0.6184437779527641

R2 Score of training 0.3515107836557766
R2 score of training and testing is not much different.

Predicted [6.05272976]
Actual 1109 6
Name: quality, dtype: int64
Prediction for the single observation: [5.63253129]

Out[33]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol
0	0.6	0.5	0.43	0.5	0.5	0.4	0.3	0.3	0.8	0.6	0.4

In [37]:

```

1 r_pred=r.predict(x_test)
2 r_pred_1=r.predict(x_train)
Lasso and Ridge
3 l_pred=l.predict(x_test)
4 l_pred_1=l.predict(x_train)

```

In [34]:

```

1 r=Ridge()
2 Qubaatly(pd.DataFrame({'Actual_Profit':y_test, 'Ridge_pred':r_pred, 'Lasso_1000':l_pred_1}), 10)
3 Quality.head(10)

```

Out[35]:

	Actual_Profit	Ridge_pred	Lasso_pred
1109	6	6.043328	5.646599
1032	5	4.950493	5.646599

In [36]:

```

1 l.fit(x_train,y_train)
487 l.fit(x_train,y_train)
5.56834 5.646599

```

Out[36]:

	Actual_Profit	Ridge_pred	Lasso_pred
979	5	5.827819	5.646599
1054	6	5.187421	5.646599
542	5	5.426785	5.646599
853	6	5.955455	5.646599

```

In [37]: 1 r_pred=r.predict(x_test)
          2 r_pred_1=r.predict(x_train)
          Lasso and Ridge
          3 l_pred=l.predict(x_test)
          4 l_pred_1=l.predict(x_train)

In [34]: 1 r=Ridge()
In [38]: 2 Q=Quality(pd.DataFrame({'Actual_Profit':y_test, 'Ridge_pred':r_pred, 'Lasso':l_pred_1}))
          3 Quality.head(10)

Out[38]: 1 r.fit(x_train,y_train)
          Actual_Profit  Ridge_pred  Lasso_pred
          +-----+
          1109           6   6.043328  5.646599
          Ridge()
          1032           5   4.950493  5.646599
          1002           7   6.681428  5.646599
          487            6   5.156334  5.646599

In [36]: 1 l.fit(x_train,y_train)
          Lasso()
          979            5   5.827819  5.646599
          1054           6   5.187421  5.646599
          542            5   5.426785  5.646599
          853            6   5.955455  5.646599
          1189           4   4.752772  5.646599
          412            5   5.121926  5.646599

```

```

In [39]: 1 print("For ridge\n")
          2
          3 #R2 Score
          4 print("R2 score",metrics.r2_score(y_test,r_pred))
          5 #printing to check if overfitting or underfitting
          6 #R2 Score
          7 print("R2 score for training",metrics.r2_score(y_train,r_pred_1))
          8 #MSE (Mean Square Error)
          9 print("MSE",metrics.mean_squared_error(y_test,r_pred))
          10 #RMSE (Root Mean Square Error)
          11 print("RMSE",np.sqrt(metrics.mean_squared_error(y_test,r_pred)))
          12
          13 print('\n')
          14 #Checking a random obseration:-
          15 predicted_values = r.predict(x_test.iloc[[0]])
          16 actual_values = y_test.iloc[[0]]
          17 print("Predicted " + str(predicted_values))
          18 print("Actual " + str(actual_values))
          19
          20 column_names = x_test.columns
          21
          22 # Input data for a single observation
          23 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]

```

```

In [39]: 1 print("For ridge\n")
2
3 #R2 Score
4 print("R2 score",metrics.r2_score(y_test,r_pred))
5 #printing to check if overfitting or underfitting
6 #R2 Score
7 print("R2 score for training",metrics.r2_score(y_train,r_pred_1))
8 #MSE (Mean Square Error)
9 print("MSE",metrics.mean_squared_error(y_test,r_pred))
10 #RMSE (Root Mean Square Error)
11 print("RMSE",np.sqrt(metrics.mean_squared_error(y_test,r_pred)))
12
13 print('\n')
14 #Checking a random observation:- 
15 predicted_values = r.predict(x_test.iloc[[0]])
16 actual_values = y_test.iloc[[0]]
17 print("Predicted " + str(predicted_values))
18 print("Actual " + str(actual_values))
19
20 column_names = x_test.columns
21
22 # Input data for a single observation
23 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
24 input_df = pd.DataFrame([input_data], columns=column_names)
25
26 # Make a prediction for the single observation
27 prediction_single_observation = r.predict(input_df)
28
29 print("Prediction for the single observation:", prediction_single_observation)
30
31
32
33 print('\n')
34 print("For Lasso")
35
36 #R2 Score
37 print("R2 score",metrics.r2_score(y_test,l_pred))
38 #printing to check if overfitting or underfitting
39 #R2 Score
40 print("R2 score with training",metrics.r2_score(y_train,l_pred_1))
41 #MSE (Mean Square Error)
42 print("MSE",metrics.mean_squared_error(y_test,l_pred))
43 #RMSE (Root Mean Square Error)
44 print("RMSE",np.sqrt(metrics.mean_squared_error(y_test,l_pred)))
45
46 print('\n')
47
48 #Checking a random observation:- 
49 predicted_values = l.predict(x_test.iloc[[0]])
50 actual_values = y_test.iloc[[0]]
51 print("Predicted " + str(predicted_values))
52 print("Actual " + str(actual_values))
53
54 column_names = x_test.columns
55
56 # Input data for a single observation
57 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
58 input_df = pd.DataFrame([input_data], columns=column_names)
59
60 # Make a prediction for the single observation
61 prediction_single_observation = l.predict(input_df)
62
63 print("Prediction for the single observation:", prediction_single_observation)
64 R2 score 0.33460003973813257
65 R2 score 0.35135292178989475
66 MSE 0.3809154850639714
67 RMSE 0.6171835100389279

```

```

Predicted [6.04332824]
Actual 1109      6
Name: quality, dtype: int64
Prediction for the single observation: [5.65709029]

```

```

For Lasso
R2 score -0.004878947397931155
R2 score with training 0.0
MSE 0.5752539443014328
RMSE 0.7584549718351333

```

```

Predicted [5.64659891]
Actual 1109      6

```

```

62
63 print("Prediction for the single observation:", prediction_single_observation)
64 score 0.33460003973813257
65 input_df
66 score for training 0.35135292178989475
MSE 0.3809154850639714
RMSE 0.6171835100389279

Predicted [6.04332824]
Actual 1109      6
Name: quality, dtype: int64
Prediction for the single observation: [5.65709029]

For Lasso
R2 score -0.004878947397931155
R2 score with training 0.0
MSE 0.5752539443014328
RMSE 0.7584549718351333

Predicted [5.64659891]
Actual 1109      6
Name: quality, dtype: int64
Prediction for the single observation: [5.64659891]

```

Out[39]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol
0	0.6	0.5	0.43	0.5	0.5	0.4	0.3	0.3	0.8	0.6	0.4

In [40]:

```

1 #Linear Regression and Ridge Regression both have similar R2 scores
2 #and perform better than Lasso Regression in terms of R2.

```

Logistic Regression

In [41]:

```

1 model=LogisticRegression()
2 model.fit(x_train,y_train)

```

```
C:\Users\Charvi Upadhyay\AppData\Roaming\Python\Python39\site-packages\sklearn
\linear_model\_logistic.py:460: ConvergenceWarning: lbfgs failed to converge
(status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.
```

```
Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html (https://scikit-learn.org/stable/modules/preprocessing.html)
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear\_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear\_model.html#logistic-regression)
    n_iter_i = _check_optimize_result()
```

Out[41]:

```

▼ LogisticRegression
LogisticRegression()

```

In [42]:

```

1 pred_log=model.predict(x_test)
2 pred_log_1=model.predict(x_train)
3
4 #Checking a random observation:
5 predicted_values = model.predict(x_test.iloc[[0]])
6 actual_values = y_test.iloc[[0]]
7
8 print("Predicted " + str(predicted_values))
9 print("Actual " + str(actual_values))
10
11 column_names = x_test.columns
12
13 # Input data for a single observation
14 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
15 input_df = pd.DataFrame([input_data], columns=column_names)
16
17 # Make a prediction for the single observation
18 prediction_single_observation = model.predict(input_df)
19
20 print("Prediction for the single observation:", prediction_single_observation)
21
22 input_df

```

```

In [42]: 1 pred_log=model.predict(x_test)
2 pred_log_1=model.predict(x_train)
3
4 #Checking a random observation:-
5 predicted_values = model.predict(x_test.iloc[[0]])
6 actual_values = y_test.iloc[[0]]
7
8 print("Predicted " + str(predicted_values))
9 print("Actual " + str(actual_values))
10
11 column_names = x_test.columns
12
13 # Input data for a single observation
14 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
15 input_df = pd.DataFrame([input_data], columns=column_names)
16
17 # Make a prediction for the single observation
18 prediction_single_observation = model.predict(input_df)
19
20 print("Prediction for the single observation:", prediction_single_observation)
21
22 input_df

```

```

Predicted [6]
Actual 1109      6
Name: quality, dtype: int64
Prediction for the single observation: [6]

```

Out[42]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol
0	0.6	0.5	0.43	0.5	0.5	0.4	0.3	0.3	0.8	0.6	0.4

```

In [43]: 1 print("Accuracy score: ",accuracy_score(y_test,pred_log))
2 print("Accuracy score with training: ",accuracy_score(y_train,pred_log_1)
3 print("Confusion matrix")
4 print(confusion_matrix(y_test,pred_log))
5 print("Crosstab")
6 pd.crosstab(y_test,pred_log)

```

```

Accuracy score: 0.621875
Accuracy score with training: 0.5957779515246286
Confusion matrix
[[ 0   0   2   0   0   0]
 [ 0   0   7   4   0   0]
 [ 0   0 105  29   1   0]
 [ 0   0  46  85  11   0]
 [ 0   0   2  16   9   0]
 [ 0   0   0   1   2   0]]
Crosstab

```

Out[43]:

col_0	5	6	7
quality			
3	2	0	0
4	7	4	0
5	105	29	1

```

In [44]: 1 print(classification_report(y_test,pred_log))
2
3           precision    recall  f1-score   support
4
5             0       31      2      0.00      0.00      0.00      2
6               4       0.00      0.00      0.00      0.00     11
7               5       0.65      0.78      0.71     135
8               6       0.63      0.60      0.61     142
9               7       0.39      0.33      0.36     27
10              8       0.00      0.00      0.00      0.00      3
11
12    accuracy                           0.62      320
13   macro avg       0.28      0.28      0.28      320
14 weighted avg       0.59      0.62      0.60      320

```

```

C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))

```

```

In [44]: 1 print(classification_report(y_test,pred_log))
          2 precision    recall   f1-score   support
          3      0.00     0.00     0.00      2
          4      0.00     0.00     0.00     11
          5      0.65     0.78     0.71    135
          6      0.63     0.60     0.61    142
          7      0.39     0.33     0.36     27
          8      0.00     0.00     0.00      3
          accuracy         0.62      320
          macro avg       0.28     0.28     0.28     320
          weighted avg    0.59     0.62     0.60     320

```

```

C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))

```

```

In [45]: 1 print("Here the columns 3,4,8 were lost, this is because there samples we
          Here the columns 3,4,8 were lost, this is because there samples were less.

```

Decision Tree

```

In [46]: 1 model1=DecisionTreeClassifier()
          2 model1.fit(x_train,y_train)

```

```

Out[46]: ▾ DecisionTreeClassifier
          DecisionTreeClassifier()

```

```

In [47]: 1 d_y_predict = model1.predict(x_test)
          2 d_y_predict_train = model1.predict(x_train)
          3
          4 #Checking a random observation:- 
          5 predicted_values = model1.predict(x_test.iloc[[0]])
          6 actual_values = y_test.iloc[[0]]
          7
          8 print("Predicted " + str(predicted_values))
          9 print("Actual " + str(actual_values))
          10
          11 column_names = x_test.columns
          12
          13 # Input data for a single observation

```

```

In [48]: 14 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
          15 print(pd.DataFrame([input_data],columns=x_test.columns))
          16 print("Accuracy score with training: ",accuracy_score(y_train,d_y_predict))
          17 print("Confusion matrix")
          18 print(confusion_matrix(y_test,d_y_predict))
          19 print("Crossstab")
          20 print(classification_report(y_test,d_y_predict))
          21 pd.crosstab(y_test,d_y_predict)
          22 input_df
          Accuracy score: 0.621875
          Accuracy score with training: 1.0
          Confusion matrix
          Name: Quality, dtype: int64
          Prediction for the single observation: [7]
          [ 0  6 89 33  6  1]

```

```

Out[47]: [ 0  3 25 91 23  0]
          [ 0  0  1 7 17 21]
          [fixed acidity, volatile acidity2, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates, alcohol]
          Crosstab
          0      0.6      0.5      0.5      0.4      0.3      0.3      0.3      0.2      0.1      0.1
          1      0.00     0.00     0.00     0.00      2
          2      0.18     0.18     0.18     0.18     11
          3      0.74     0.66     0.70     0.70    135

```

```
In [48]:  
14 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]  
15 print(pd.DataFrame([input_data], columns=column_names))  
16 print("Accuracy score with training: ", accuracy_score(y_train, d_y_predict))  
17 print("Confusion matrix")  
18 print(crosstab(y_test, d_y_predict))  
19 print("Crossstab")  
20 print(classification_report(y_test, d_y_predict)) prediction_single_observation  
21 pd.crosstab(y_test, d_y_predict)  
  
22 input_df  
Accuracy score: 0.621875
```

```
AUCROC score with training: 1.0  
Confusion matrix  
[[Name: Quality, Dtype: int64  
Prediction for the single observation: [7]  
[ 0 6 89 33 6 1]  
[ 0 3 25 91 23 0]  
[ 0 fixed acidity 17.0 residual sugar 1.0 chlorides 0.0 free sulfur dioxide 0.0 total sulfur dioxide 0.6 density 0.8 pH 0.6 sulphates 0.0 alcohol 0.4]  
Crosstab  
0 0.6 0.5 precision 0.5 recall 0.5 f1-score support 0.3 0.8 0.6 0.4  
0 3 0.00 0.00 0.00 2  
4 0.18 0.18 0.18 11  
5 0.74 0.66 0.70 135  
6 0.66 0.64 0.65 142  
7 0.35 0.63 0.45 27  
8 0.00 0.00 0.00 3  
  
accuracy 0.62 320  
macro avg 0.32 0.35 0.33 320  
weighted avg 0.64 0.62 0.63 320
```

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\networks\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\networks\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\networks\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))
```

```
Out[48]:  
col_0 4 5 6 7 8  
quality  
-----  
3 0 1 1 0 0  
4 2 4 5 0 0  
5 6 89 33 6 1  
6 3 25 91 23 0  
7 0 1 7 17 2  
8 0 0 1 2 0
```

```
In [49]: 1 print(classification_report(y_test, d_y_predict))
```

```
precision recall f1-score support  
  
3 0.00 0.00 0.00 2  
4 0.18 0.18 0.18 11  
5 0.74 0.66 0.70 135  
6 0.66 0.64 0.65 142  
7 0.35 0.63 0.45 27  
8 0.00 0.00 0.00 3  
  
accuracy 0.62 320  
macro avg 0.32 0.35 0.33 320  
weighted avg 0.64 0.62 0.63 320
```

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\networks\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\networks\_classification.py:1469: UndefinedMetricWarning: Precision and F-score
```

```
In [49]: 1 print(classification_report(y_test,d_y_predict))
```

	precision	recall	f1-score	support
3	0.00	0.00	0.00	2
4	0.18	0.18	0.18	11
5	0.74	0.66	0.70	135
6	0.66	0.64	0.65	142
7	0.35	0.63	0.45	27
8	0.00	0.00	0.00	3
accuracy			0.62	320
macro avg	0.32	0.35	0.33	320
weighted avg	0.64	0.62	0.63	320

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\network\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\network\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\network\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))
```

```
In [53]: 1 print("The model seems to be overfit.")  
2 print("Unlike logistic regression, No quality was lost in predictions.")
```

The model seems to be overfit.
Unlike logistic regression, No quality was lost in predictions.

```
In [51]: 1 dot_data = StringIO()  
2 export_graphviz(model1, out_file=dot_data, feature_names=x.columns, filled=True)  
3 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())  
4 graph.set_dpi(5)  
5 Image(graph.create_png())  
6
```



Random Forest Classifier

```
In [52]: 1 model2=RandomForestClassifier()  
2 model2.fit(x_train,y_train)
```

Out[52]:

```
RandomForestClassifier  
| RandomForestClassifier()  
| |
```

```
In [54]: 1 r_y_predict=model2.predict(x_test)  
2 r_y_predict_train=model2.predict(x_train)  
3  
4 #Checking a random observation:-  
5 predicted_values = model2.predict(x_test.iloc[[0]])  
6 actual_values = y_test.iloc[[0]]  
7  
8 print("Predicted " + str(predicted_values))  
9 print("Actual " + str(actual_values))  
10  
11 column_names = x_test.columns  
12  
13 # Input data for a single observation  
14 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]  
15 input_df = pd.DataFrame([input_data], columns=column_names)  
16  
17 # Make a prediction for the single observation  
18 prediction_single_observation = model2.predict(input_df)  
19  
20 print("Prediction for the single observation:", prediction_single_observation)  
21  
22 input_df
```

```

In [54]: 1 r_y_predict=model2.predict(x_test)
2 r_y_predict_train=model2.predict(x_train)
3
4 #Checking a random observation:- 
5 predicted_values = model2.predict(x_test.iloc[[0]])
6 actual_values = y_test.iloc[[0]]
7
8 print("Predicted " + str(predicted_values))
9 print("Actual " + str(actual_values))
10
11 column_names = x_test.columns
12
13 # Input data for a single observation
14 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
15 input_df = pd.DataFrame([input_data], columns=column_names)
16
17 # Make a prediction for the single observation
18 prediction_single_observation = model2.predict(input_df)
19
20 print("Prediction for the single observation:", prediction_single_observation)
21
22 input_df

```

Predicted [6]
 Actual 1109 6
 Name: quality, dtype: int64
 Prediction for the single observation: [6]

Out[54]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol
0	0.6	0.5	0.43	0.5	0.5	0.4	0.3	0.3	0.8	0.6	0.4

```

In [55]: 1 print("Accuracy score: ",accuracy_score(y_test,r_y_predict ))
2 print("Accuracy score with training: ",accuracy_score(y_train,r_y_predict))
3 print("Confusion matrix")
4 print(confusion_matrix(y_test,r_y_predict))
5 print("Crosstab")
6 print(classification_report(y_test,r_y_predict))
7 pd.crosstab(y_test,r_y_predict)

```

Accuracy score: 0.728125
 Accuracy score with training: 1.0
 Confusion matrix
 [[0 0 0 2 0 0]
 [0 0 6 5 0 0]
 [0 0 112 21 2 0]
 [0 0 25 106 11 0]
 [0 0 1 11 15 0]
 [0 0 0 1 2 0]]
 Crosstab

	precision	recall	f1-score	support
3	0.00	0.00	0.00	2
4	0.00	0.00	0.00	11
5	0.78	0.83	0.80	135

```
In [55]: 1 print("Accuracy score: ",accuracy_score(y_test,r_y_predict ))
2 print("Accuracy score with training: ",accuracy_score(y_train,r_y_predict))
3 print("Confusion matrix")
4 print(confusion_matrix(y_test,r_y_predict))
5 print("Crosstab")
6 print(classification_report(y_test,r_y_predict))
7 pd.crosstab(y_test,r_y_predict)
```

```
Accuracy score: 0.728125
Accuracy score with training: 1.0
Confusion matrix
[[ 0  0  0  2  0  0]
 [ 0  0  6  5  0  0]
 [ 0  0 112 21  2  0]
 [ 0  0 25 106 11  0]
 [ 0  0  1 11 15  0]
 [ 0  0  0  1  2  0]]
Crosstab
      precision    recall   f1-score   support
          3       0.00     0.00     0.00       2
          4       0.00     0.00     0.00      11
          5       0.78     0.83     0.80     135
          6       0.73     0.75     0.74     142
          7       0.50     0.56     0.53      27
          8       0.00     0.00     0.00       3
   accuracy           0.73     320
macro avg       0.33     0.36     0.34     320
weighted avg    0.69     0.73     0.71     320
```

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
```

```
Out[55]:
  col_0    5    6    7
quality
-----
  3    0    2    0
  4    6    5    0
  5  112   21   2
  6   25  106  11
  7    1   11  15
  8    0    1    2
```

```
In [56]: 1 print(classification_report(y_test,r_y_predict))
```

```
      precision    recall   f1-score   support
          3       0.00     0.00     0.00       2
          4       0.00     0.00     0.00      11
          5       0.78     0.83     0.80     135
          6       0.73     0.75     0.74     142
          7       0.50     0.56     0.53      27
          8       0.00     0.00     0.00       3
   accuracy           0.73     320
macro avg       0.33     0.36     0.34     320
weighted avg    0.69     0.73     0.71     320
```

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
ore are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn
\metrics\_classification.py:1469: UndefinedMetricWarning: Precision and F-sc
```

```
In [56]: 1 print(classification_report(y_test,r_y_predict))
```

	precision	recall	f1-score	support
3	0.00	0.00	0.00	2
4	0.00	0.00	0.00	11
5	0.78	0.83	0.80	135
6	0.73	0.75	0.74	142
7	0.50	0.56	0.53	27
8	0.00	0.00	0.00	3
accuracy			0.73	320
macro avg	0.33	0.36	0.34	320
weighted avg	0.69	0.73	0.71	320

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\naive_bayes\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\naive_bayes\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\naive_bayes\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted sample
s. Use `zero_division` parameter to control this behavior.
    _warn_prf(average, modifier, msg_start, len(result))
```

```
In [57]: 1 #Random forest based on entropy
2 model3 = RandomForestClassifier(criterion='entropy')
3 model3.fit(x_train, y_train)
```

```
Out[57]: RandomForestClassifier
RandomForestClassifier(criterion='entropy')
```

```
In [58]: 1 r_y_predict_1 = model3.predict(x_test)
2 r_y_predict_train_1 = model3.predict(x_train)
3
4 #Checking a random observation:- 
5 predicted_values = model3.predict(x_test.iloc[[0]])
6 actual_values = y_test.iloc[[0]]
7
8 print("Predicted " + str(predicted_values))
9 print("Actual " + str(actual_values))
10
11 column_names = x_test.columns
12
13 # Input data for a single observation
14 input_data = [0.6, 0.5, 0.43, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
15 input_df = pd.DataFrame([input_data], columns=column_names)
16
17 # Make a prediction for the single observation
18 prediction_single_observation = model3.predict(input_df)
19
20 print("Prediction for the single observation:", prediction_single_observation)
In [59]: 21 print("Accuracy score: ",accuracy_score(y_test,r_y_predict_1))
22 print("Accuracy score with training: ",accuracy_score(y_train,r_y_predict_1))
23 print("Confusion matrix")
24 print(confusion_matrix(y_test,r_y_predict_1))
Actual 1199
25 print("Crosstab")
Name: quality, dtype: int64
26 print(pd.crosstab(y_test,r_y_predict_1))
Prediction for the single observation: [6]
27 pd.crosstab(y_test,r_y_predict_1)
```

```
Out[58]: Accuracy score: 0.734375
fixed acidity volatile acidity residual sugar
fixed acidity volatile acidity residual sugar
Confusion matrix
[[ 0  0.6  0  1.05  1  0.48  0] 0.5  0.5  0.4  0.3  0.3  0.8  0.6  0.4
 [ 0  1  7  3  0  0] 0
 [ 0  0  112  21  2  0] 0
 [ 0  0  27  105  10  0] 0
 [ 0  0  1  9  17  0] 0
 [ 0  0  0  1  2  0]] 0
Crosstab
precision recall f1-score support
3 0.00 0.00 0.00 2
4 1.00 0.09 0.17 11
5 0.76 0.83 0.79 135
```

```
In [59]:  
1 print("Accuracy score: ",accuracy_score(y_test,r_y_predict_1))  
2 input()  
3 print("Accuracy score with training: ",accuracy_score(y_train,r_y_predict_1))  
4 Predicted[6]  
5 print(confusion_matrix(y_test,r_y_predict_1))  
Actual[119]  
6 print("Crosstab")  
Name: quality, dtype: int64  
7 print(classification_report(y_test,r_y_predict_1))  
Prediction for the single observation: [6]  
pd.crosstab(y_test,r_y_predict_1)
```

```
Out[58]:  
Accuracy score: 0.734375  
fixed volatile citric acid residual chlrides free sulfur total  
Acidity acidity acid sugar dioxide sulphur dioxide density pH sulphates alcohol  
Confusion matrix  
[[ 0  0.6  1  0.5  1  0.43  0 ] 0.5] 0.5  0.4  0.3  0.3  0.8  0.6  0.4  
[[ 0  1  7  3  0  0 ]]  
[[ 0  0  112  21  2  0 ]]  
[[ 0  0  27  105  10  0 ]]  
[[ 0  0  1  9  17  0 ]]  
[[ 0  0  0  1  2  0 ]]  
Crosstab  
precision recall f1-score support  
  
3 0.00 0.00 0.00 2  
4 1.00 0.09 0.17 11  
5 0.76 0.83 0.79 135  
6 0.75 0.74 0.74 142  
7 0.55 0.63 0.59 27  
8 0.00 0.00 0.00 3  
  
accuracy 0.73 320  
macro avg 0.51 0.38 0.38 320  
weighted avg 0.73 0.73 0.72 320
```

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\naive_bayes\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\naive_bayes\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\naive_bayes\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))
```

```
Out[59]:  
col_0 4 5 6 7  
quality  
-----  
3 0 1 1 0  
4 1 7 3 0  
5 0 112 21 2  
6 0 27 105 10  
7 0 1 9 17  
8 0 0 1 2
```

```
In [60]:  
1 print(classification_report(y_test,r_y_predict_1))  
  
precision recall f1-score support  
  
3 0.00 0.00 0.00 2  
4 1.00 0.09 0.17 11  
5 0.76 0.83 0.79 135  
6 0.75 0.74 0.74 142  
7 0.55 0.63 0.59 27  
8 0.00 0.00 0.00 3  
  
accuracy 0.73 320  
macro avg 0.51 0.38 0.38 320  
weighted avg 0.73 0.73 0.72 320
```

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\naive_bayes\classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\naive_bayes\classification.py:1469: UndefinedMetricWarning: Precision and F-score
```

```
In [60]: 1 print(classification_report(y_test,r_y_predict_1))
```

	precision	recall	f1-score	support
3	0.00	0.00	0.00	2
4	1.00	0.09	0.17	11
5	0.76	0.83	0.79	135
6	0.75	0.74	0.74	142
7	0.55	0.63	0.59	27
8	0.00	0.00	0.00	3
accuracy			0.73	320
macro avg	0.51	0.38	0.38	320
weighted avg	0.73	0.73	0.72	320

```
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\numerical\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\numerical\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))  
C:\Users\Charvi Upreti\AppData\Roaming\Python\Python39\site-packages\sklearn\numerical\_classification.py:1469: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.  
_warn_prf(average, modifier, msg_start, len(result))
```

```
In [61]: 1 print("Maximum Accuracy reached was with Random forest classifier and is")
```

```
Maximum Accuracy reached was with Random forest classifier and is 0.734375
```

Trying binary classification to get a better accuracy (Good: 1 vs Bad: 0 Wine)

```
In [62]: 1 df['quality'].value_counts().sort_index(ascending=True)
```

```
Out[62]: 3    10  
4    53  
5   681  
6   638  
7   199  
8    18  
Name: quality, dtype: int64
```

```
In [63]: 1 y_train = y_train.replace([3, 4, 5, 6], 0)  
2 y_train = y_train.replace([7, 8], 1)  
3 y_train
```

```
Out[63]: 642    0  
679    0  
473    0  
390    1  
1096   0  
..
```

```
In [65]: 763 y_test = y_test.replace([3, 4, 5, 6], 0)  
835 y_test = y_test.replace([7, 8], 1)  
1216 y_test  
559  0
```

```
Out[65]: 649    0  
Name: quality, Length: 1279, dtype: int64
```

```
1002    1  
In [64]: 487 y_train.value_counts()  
979    0
```

```
Out[64]: 0    1092  
1    189  
Name: quality, dtype: int64  
1322    0  
704    0  
1023   0  
Name: quality, Length: 320, dtype: int64
```

```
In [66]: 1 y_test.value_counts()
```

```
Out[66]: 0    290  
1    30  
Name: quality, dtype: int64
```

```
In [65]: ..  
763 y_test[0] = y_test.replace([3, 4, 5, 6], 0)  
835 y_test[0] = y_test.replace([7, 8], 1)  
1216 y_test[0]  
559 0
```

```
Out[65]: 0  
Name: quality, Length: 1279, dtype: int64  
1002 1
```

```
In [64]: 487 y_train.value_counts()  
979 0
```

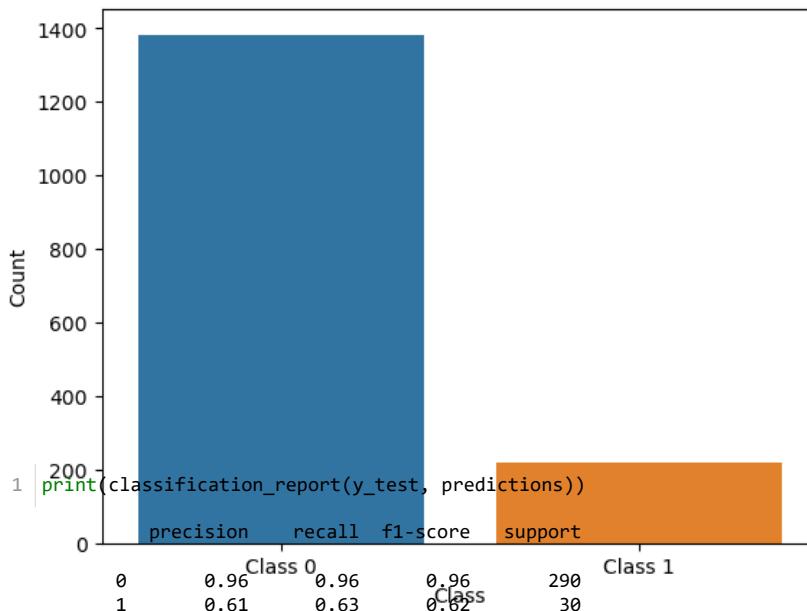
```
Out[64]: 0    1092  
1    180  
Name: quality, dtype: int64  
1322 0  
704 0  
1023 0  
Name: quality, Length: 320, dtype: int64
```

```
In [66]: 1 y_test.value_counts()
```

```
Out[66]: 0    290  
1    30  
Name: quality, dtype: int64
```

```
In [67]: 1 import seaborn as sns  
2 import matplotlib.pyplot as plt  
3  
4 combined_y = np.concatenate((y_test, y_train))
```

```
In [68]: 1 # Create a countplot using Seaborn  
2 sns.countplot(x=combined_y)  
3  
4 # Set labels for x and y axes  
5 plt.xlabel("Class")  
6 plt.ylabel("Count")  
7  
8 # Set class labels for the x-axis ticks (0 and 1)  
9 plt.xticks([0, 1], ["Class 0", "Class 1"])  
10  
11 # Show the plot  
12 plt.show()
```



```
In [71]: 1 print(classification_report(y_test, predictions))
```

```
In [69]: 1 accuracy_random = RandomForestClassifier(0.93)  
2 macro_avg = 0.79  
3 weighted_avg = 0.79  
4 model_random.fit(x_train, y_train)  
5 predictions = model_random.predict(x_test)  
6  
# Accuracy Score.
```

```
In [72]: 6 print(accuracy_score(y_test, predictions))
```

```
0.928125  
Input data for a single observation  
4 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]  
5 input_df = pd.DataFrame([input_data], columns=column_names)  
6  
Having only 2 classes increased the accuracy to: 0.928125  
7  
# Make a prediction for the single observation  
8 prediction_single_observation = model_random.predict(input_df)  
9  
print("Prediction for the single observation:", prediction_single_observation)
```

```
In [71]: 1 | print(classification_report(y_test, predictions))
0 |          precision    recall   f1-score   support
|           0         Class 0      0.96      0.96      0.96      290
|           1         Class 1      0.61      0.63      0.62      30
```

```
In [69]: 1 accuracy = 0.93
2 macro avg = 0.79
3 # Train the model on the training data
4 model_random = RandomForestClassifier()
5 model_random.fit(x_train, y_train)
6 predictions = model_random.predict(x_test)
7 # Accuracy Score.
8 print(column_accuracy(x_test,y_test,predictions))
9
10 0.928125 # Input data for a single observation
11 input_data = [0.6, 0.5, 0.43, 0.5, 0.5, 0.4, 0.3, 0.3, 0.8, 0.6, 0.4]
12 input_df = pd.DataFrame([input_data], columns=column_names)
13 print("Having only 2 classes increased the accuracy to ",accuracy_score(y_
14 Having only 2 classes increased the accuracy to 0.928125
15 prediction_single_observation = model_random.predict(input_df)
16
17 print("Prediction for the single observation:", prediction_single_observat
18
19 input_df
```

Prediction for the single observation: [0]

Out[72]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol
0	0.6	0.5	0.43	0.5	0.5	0.4	0.3	0.3	0.8	0.6	0.4

In []:

1