

```
[4]: # IMPORT ALL THE NECESSARY LIBRARIES

import warnings
warnings.filterwarnings(action="ignore")

import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns
from line import line

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

In [5]: # DATA PREPARATION

pd.set_option("display.max_rows", 12)
pd.set_option("display.max_columns", None)

# loading the Red wines dataset
df = pd.read_csv("winequality-red.csv", sep=";")

# display first few records
print(df.head())

# shape of the dataset
print(df.shape)

fixed acidity volatile acidity citric acid residual sugar chlorides \
0 7.4 0.70 0.88 34.0 0.9278 3.51 0.35
1 7.8 0.88 0.80 34.0 0.9978 3.26 0.45
2 7.8 0.76 0.84 34.0 0.9988 3.36 0.45
3 11.2 0.28 0.56 1.9 0.875
4 7.4 0.70 0.88 34.0 0.9978 3.51 0.35

free sulfur dioxide total sulfur dioxide density pH sulphates \
0 11.0 15.0 0.9968 3.20 0.40
1 25.0 67.0 0.9988 3.20 0.40
2 15.0 54.0 0.9978 3.26 0.45
3 17.0 60.0 0.9988 3.36 0.58
4 11.0 34.0 0.9978 3.51 0.35

alcohol quality
0 9.4 5
1 9.8 5
2 9.8 5
3 9.8 6
4 9.4 5
(1599, 12)

In [6]: # checking for the null values in the dataset


print(df.isnull().sum())
print("\n\n")

fixed acidity 0
volatile acidity 0
citric acid 0
residual sugar 0
chlorides 0
free sulfur dioxide 0
total sulfur dioxide 0
density 0
pH 0
sulphates 0
alcohol 0
quality 0
dtype: int64

In [7]: # checking for the outlier

fig, axes = plt.subplots(figsize=(20, 16))
df.plot(kind='box', subplots=True,
        layout=(3, 4),
        sharey=False,
        sharex=False, axes=axes)

plt.tight_layout()
plt.show()
```



```
In [8]: #EXPLORATORY DATA ANALYSIS(EDA)

#DATA ANALYSIS AND VASULIZATION

print(df.describe())
print("\n\n")

count    1599.000000    1599.000000    1599.000000    1599.000000    1599.000000
mean      8.319373      0.527821      0.276975      34.000000      2.338896
std       0.111056      0.173928      0.154001      0.099558      0.099558
min       4.600000      0.120000      0.000000      0.000000      0.000000
25%       7.200000      0.200000      0.000000      0.000000      0.000000
50%       7.900000      0.520000      0.260000      2.200000      0.000000
75%       9.200000      0.640000      0.420000      2.800000      0.000000
max      15.000000      1.580000      1.080000      15.500000      0.000000

chlorides free sulfur dioxide total sulfur dioxide density \
count    1599.000000    1599.000000    1599.000000    1599.000000
mean      8.981607      15.074822      40.407782      0.997177
std       0.847065      10.468157      32.895324      0.081887
min       0.010000      0.000000      0.000000      0.000000
25%       0.070000      0.000000      0.000000      0.000000
50%       0.070000      0.000000      0.000000      0.000000
75%       0.090000      0.000000      0.000000      0.000000
max       0.810000      0.000000      0.000000      0.000000

pH sulphates alcohol quality
count    1599.000000    1599.000000    1599.000000    1599.000000
mean      9.311133      0.658149      10.422983      5.636023
std       0.130305      0.189507      1.005068      0.081769
min       7.200000      0.120000      0.000000      0.000000
25%       8.700000      0.330000      0.400000      2.000000
50%       9.120000      0.500000      0.500000      3.000000
75%       9.800000      0.620000      0.600000      4.000000
max      10.400000      0.700000      1.100000      8.000000

In [9]: #checking for duplicate records

duplicate=df.duplicated()
print(duplicate.sum())
df(duplicate)
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
4	7.4	0.700	0.00	1.90	0.076	11.0	34.0	0.99780	3.51	0.56	9.4	5
31	7.5	0.500	0.36	6.10	0.071	17.0	102.0	0.99780	3.35	0.80	10.5	5
27	7.8	0.420	0.21	1.40	0.188	10.0	37.0	0.99680	3.17	0.01	9.5	5
46	7.3	0.460	0.36	5.90	0.074	12.0	97.0	0.99780	3.33	0.63	10.5	5
65	7.2	0.725	0.05	4.65	0.086	4.0	11.0	0.99546	3.41	0.39	10.5	5
...	...	...	...	...	...	...	...	...	...	...	...	...
1563	7.2	0.695	0.13	2.00	0.076	12.0	20.0	0.99546	3.29	0.54	10.1	5
1564	7.2	0.695	0.13	2.00	0.076	12.0	20.0	0.99546	3.29	0.54	10.1	5
1567	7.2	0.695	0.13	2.00	0.076	12.0	20.0	0.99546	3.29	0.54	10.1	5
1591	6.2	0.560	0.09	1.70	0.063	24.0	30.0	0.99402	3.54	0.60	11.3	5
1598	6.3	0.520	0.13	2.30	0.078	29.0	40.0	0.99574	3.42	0.75	11.0	6

240 rows x 12 columns


```
In [10]: #call drop_duplicate

df.drop_duplicates(inplace=True)
print(df.shape)

(1598, 12)


In [11]: # distribution of wine quality scores

plt.figure(figsize=(8, 5))
sns.countplot(x='quality', data=df)
plt.xlabel('Wine Quality')
plt.ylabel('Count')
plt.title('Distribution of Wine Quality Scores')
plt.show()
```




```
In [12]: #fixed acidity vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='fixed acidity', data=df)
print("\n\n")
```



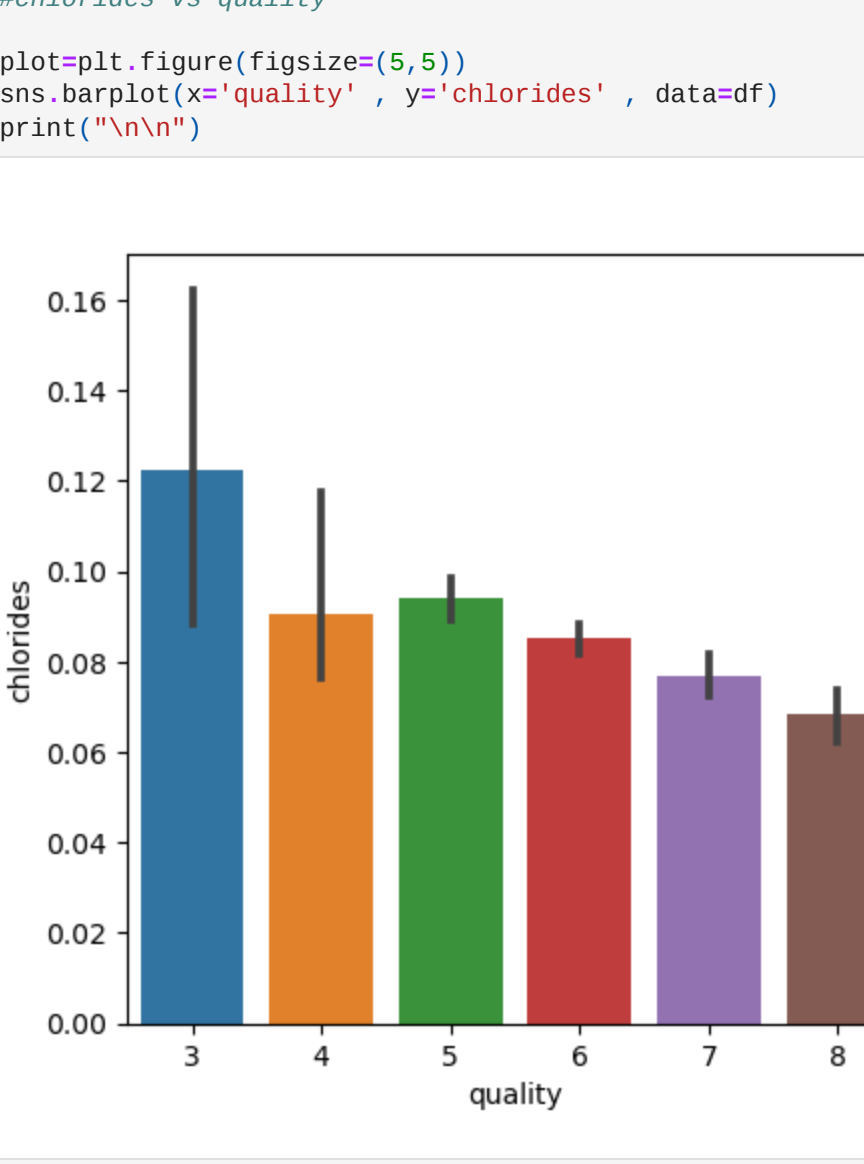
```
In [13]: #volatile acidity vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='volatile acidity', data=df)
print("\n\n")
```



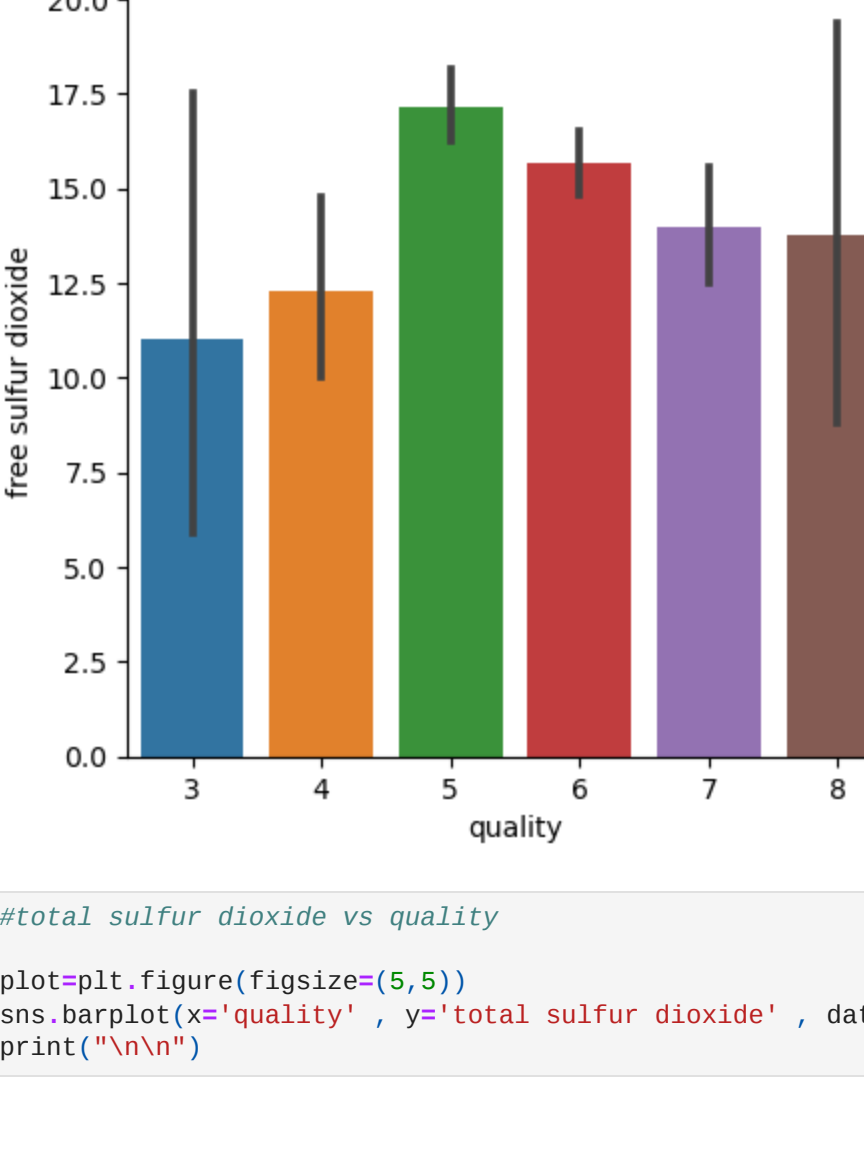
```
In [14]: #citric acid vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='citric acid', data=df)
print("\n\n")
```




```
In [15]: #residual sugar vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='residual sugar', data=df)
print("\n\n")
```




```
In [16]: #chlorides vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='chlorides', data=df)
print("\n\n")
```



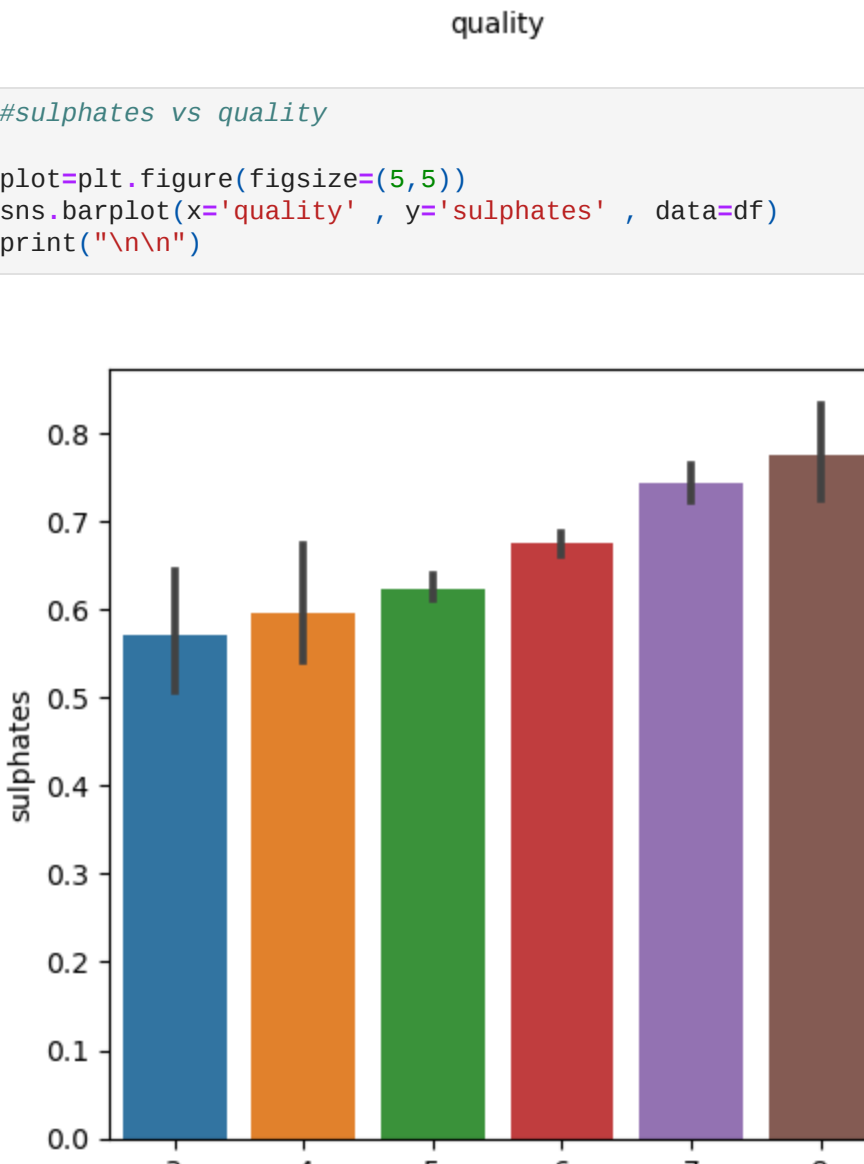
```
In [17]: #free sulfur dioxide vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='free sulfur dioxide', data=df)
print("\n\n")
```



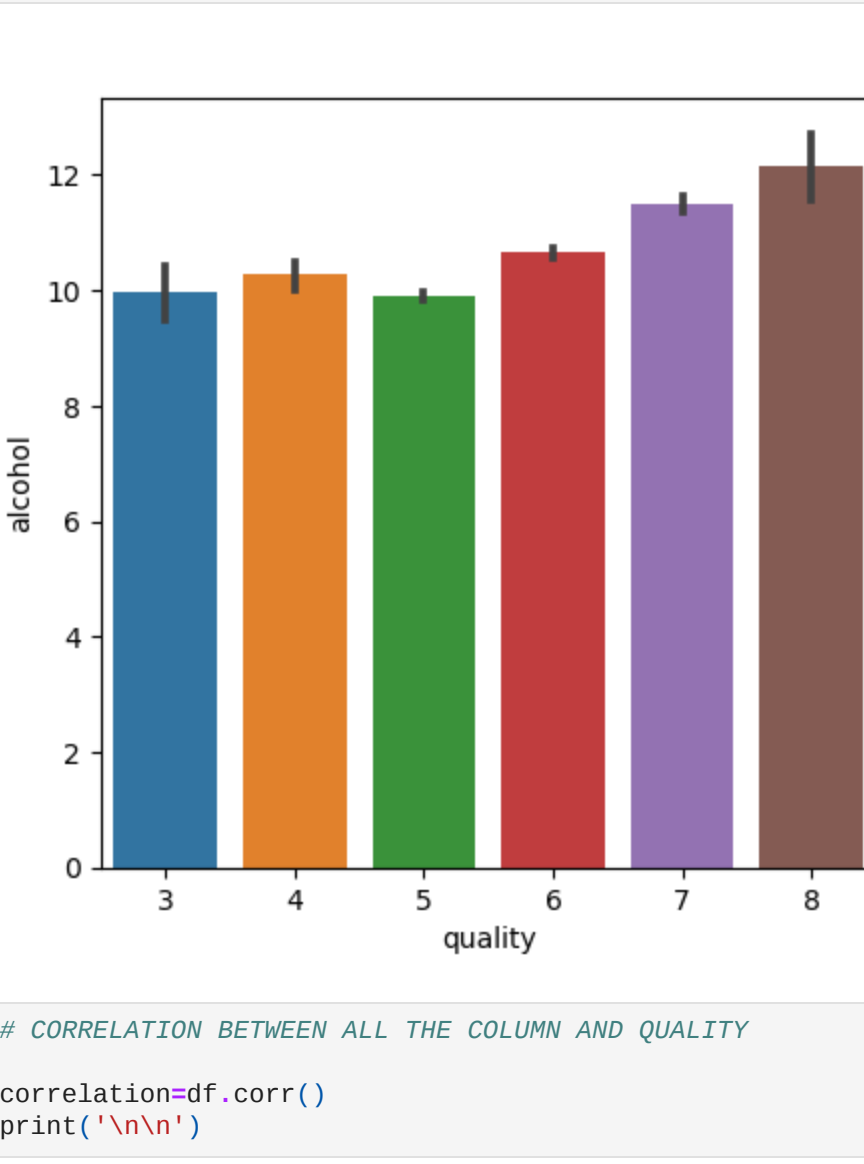
```
In [18]: #total sulfur dioxide vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='total sulfur dioxide', data=df)
print("\n\n")
```



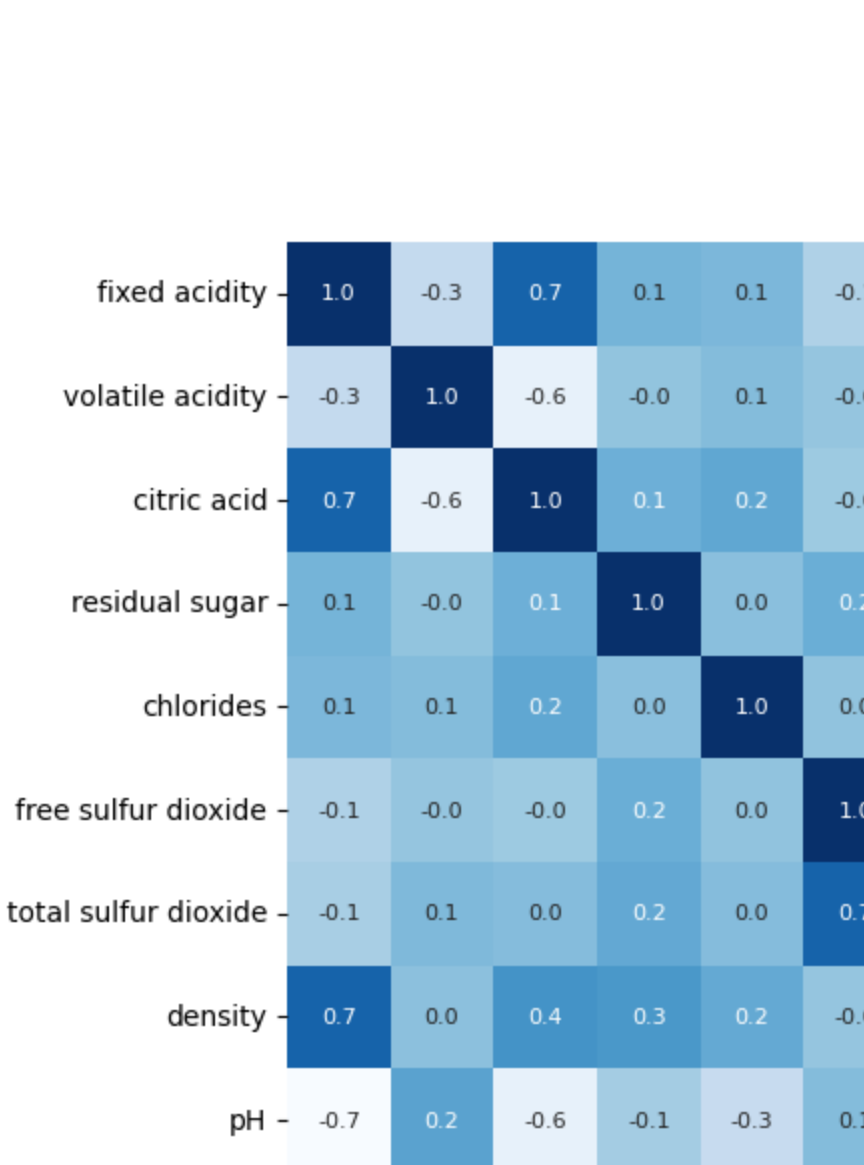
```
In [19]: #density vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='density', data=df)
print("\n\n")
```



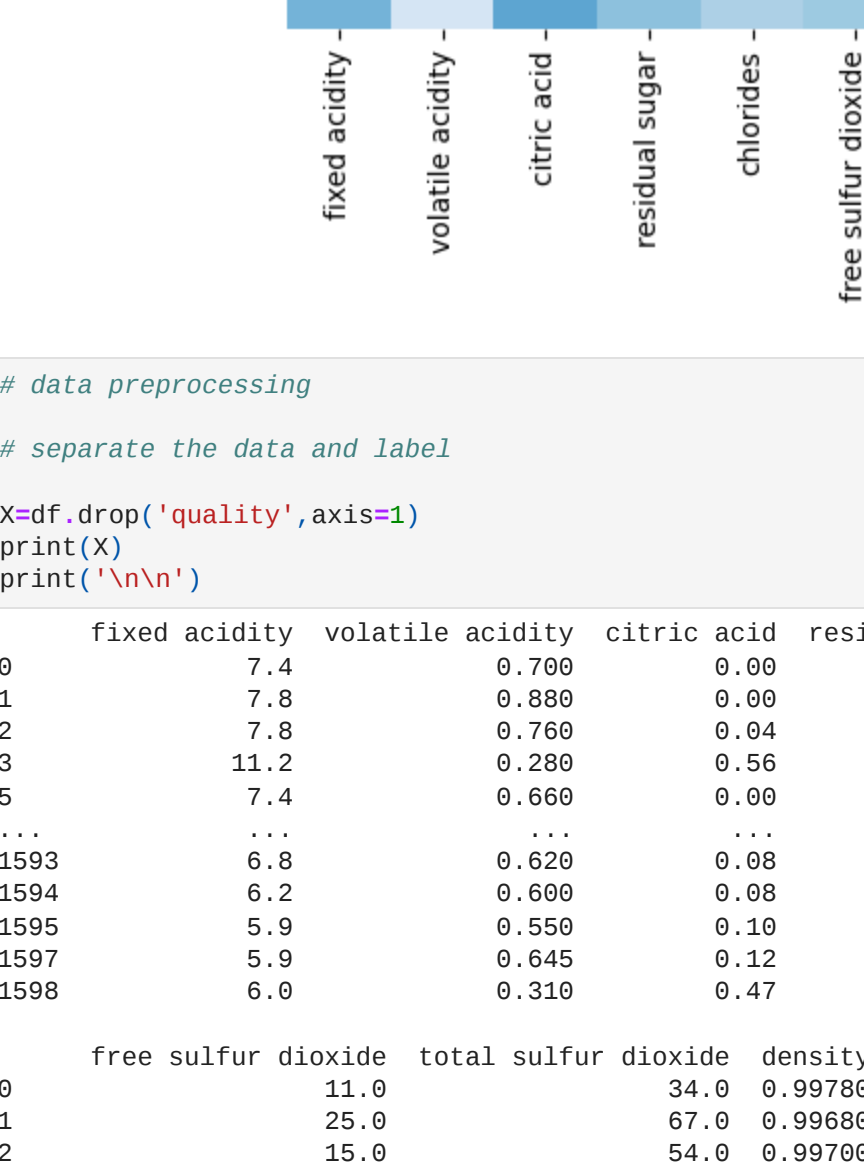
```
In [20]: #pH vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='pH', data=df)
print("\n\n")
```



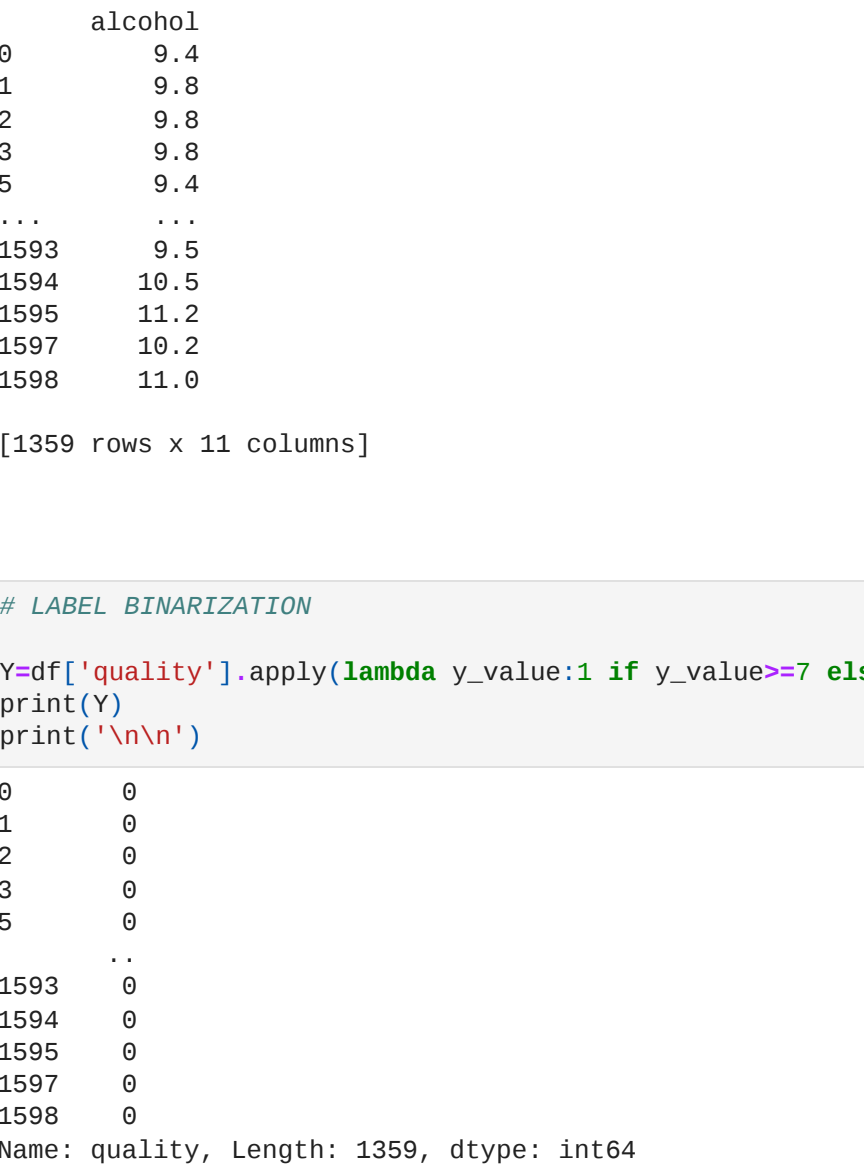
```
In [21]: #sulphates vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='sulphates', data=df)
print("\n\n")
```



```
In [22]: #alcohol vs quality

plt=plt.figure(figsize=(5,5))
sns.barplot(x='quality', y='alcohol', data=df)
print("\n\n")
```

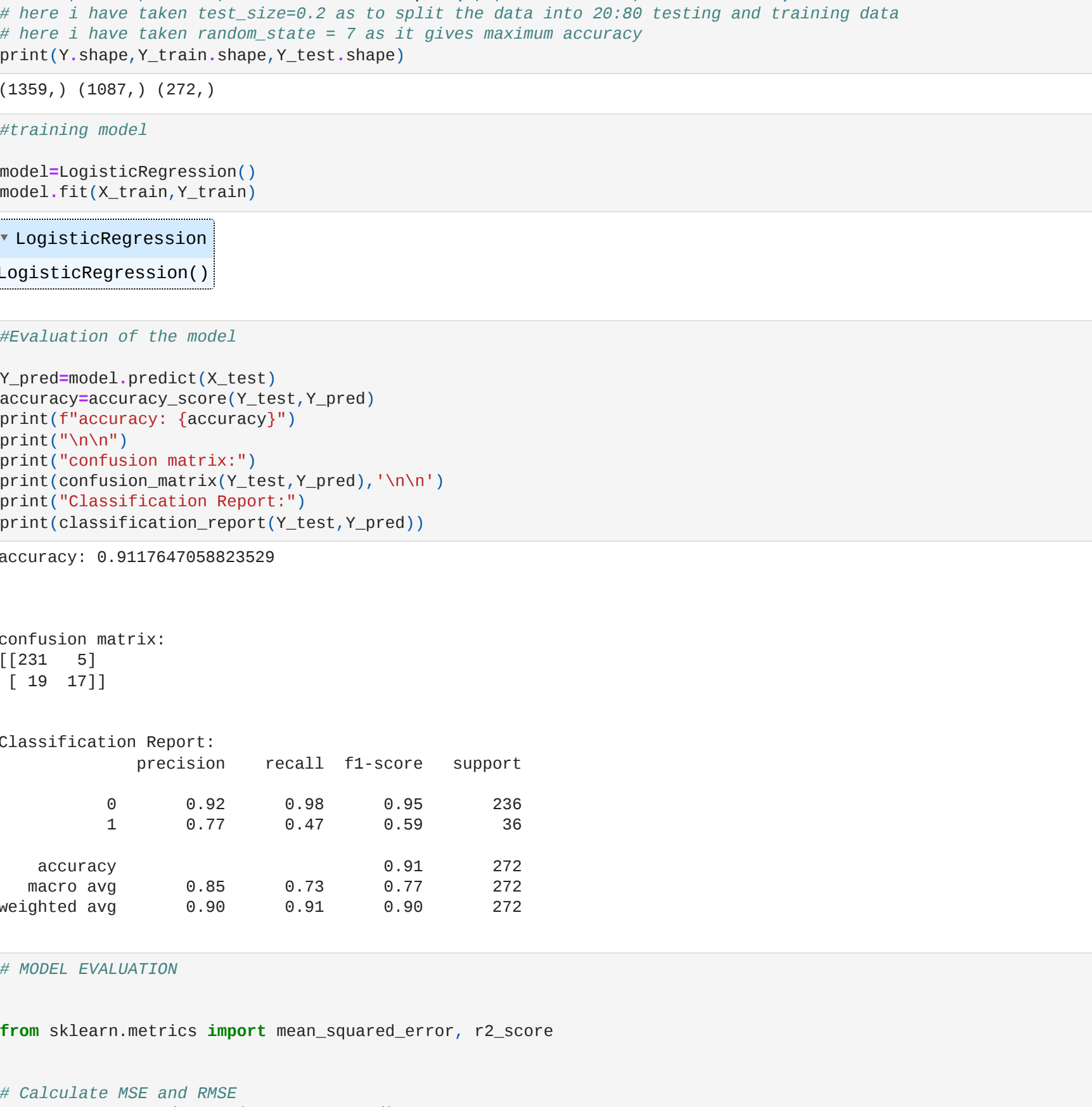


```
In [23]: # CORRELATION BETWEEN ALL THE COLUMN AND QUALITY

correlation=df.corr()
print("\n\n")

In [24]: # constructing a heatmap to understand the correlation between the column

plt.figure(figsize=(16,18))
sns.heatmap(correlation, cbar=True, square=True, fmt='.1f', annot=True, annot_kws={'size':8}, cmap='Blues')
print("\n\n")
```



```
In [25]: # data preprocessing

# separate the data and label
X=df.drop('quality',axis=1)
y=df['quality']

In [26]: # TRAIN AND TEST SPLIT

X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=7)
# here i have taken test_size=0.2 as to split the data into 20:80 (training and training data
# here i have taken random_state = 7 as to gives maximum accuracy
(1598, 11), (398, 1), (1598, 11), (398, 1)

In [27]: #training model

model=LogisticRegression()
model.fit(X_train,y_train)

Out[27]: LogisticRegression

In [28]: #Evaluation of the model

y_pred=model.predict(X_test)
accuracy=accuracy_score(y_test,y_pred)
print("accuracy: ",accuracy)

print("confusion matrix:")
print(confusion_matrix(y_test,y_pred))
print("classification_report:")
print(classification_report(y_test,y_pred))

accuracy: 0.9117647958823529

confusion matrix:
[[21  5]
 [19 17]]

Classification Report:
              precision    recall  f1-score   support

0               0.92       0.98       0.95         236
1               0.77       0.47       0.59          36

accuracy               0.85
macro avg              0.85      0.73      0.77         272
weighted avg           0.86      0.91      0.90         272

In [29]: # MODEL EVALUATION

from sklearn.metrics import mean_squared_error, r2_score

# Calculate RMSE and RMSE
mse = mean_squared_error(y_test, y_pred)
rmse = np.sqrt(mse)

# Calculate R-squared
r_squared = r2_score(y_test, y_pred)

print("RMSE: (mse)")
print("RMSE: (rmse)")
print("R-squared: (r_squared)")

RMSE: 0.13602941176478587
RMSE: 0.3682160395836603
R-squared: 0.8841612386068955

In [30]:
```