

red-wine-quality

January 13, 2024

1 COGNORISE INFOTECH __ RED WINE QUALITY ANALYSIS __ TASK 2

```
[1]: # importing dependencies
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
```

```
[3]: # Provide the full path to the CSV file
file_path = r"C:\Users\KARTHIK\OneDrive\Desktop\CognoRise Intern\Task_2\winequality-red.csv"

# Read the CSV file into a DataFrame
df = pd.read_csv(file_path)
```

```
[4]: df.head()
```

```
[4]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	\
0	7.4	0.70	0.00	1.9	0.076	
1	7.8	0.88	0.00	2.6	0.098	
2	7.8	0.76	0.04	2.3	0.092	
3	11.2	0.28	0.56	1.9	0.075	
4	7.4	0.70	0.00	1.9	0.076	

	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	\
0	11.0	34.0	0.9978	3.51	0.56	
1	25.0	67.0	0.9968	3.20	0.68	
2	15.0	54.0	0.9970	3.26	0.65	
3	17.0	60.0	0.9980	3.16	0.58	
4	11.0	34.0	0.9978	3.51	0.56	

	alcohol	quality
0	9.4	5
1	9.8	5

2	9.8	5
3	9.8	6
4	9.4	5

```
[5]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   fixed acidity          1599 non-null   float64
1   volatile acidity       1599 non-null   float64
2   citric acid            1599 non-null   float64
3   residual sugar         1599 non-null   float64
4   chlorides              1599 non-null   float64
5   free sulfur dioxide    1599 non-null   float64
6   total sulfur dioxide   1599 non-null   float64
7   density                1599 non-null   float64
8   pH                    1599 non-null   float64
9   sulphates              1599 non-null   float64
10  alcohol                1599 non-null   float64
11  quality                1599 non-null   int64
dtypes: float64(11), int64(1)
memory usage: 150.0 KB
```

```
[6]: df.isnull().sum()
```

```
[6]: 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
```

```
[7]: df.shape
```

```
[7]: (1599, 12)
```

```
[8]: df.describe()
```

```
[8]:      fixed acidity  volatile acidity  citric acid  residual sugar \
count    1599.000000      1599.000000  1599.000000      1599.000000
mean       8.319637        0.527821    0.270976        2.538806
std        1.741096        0.179060    0.194801        1.409928
min         4.600000        0.120000    0.000000        0.900000
25%         7.100000        0.390000    0.090000        1.900000
50%         7.900000        0.520000    0.260000        2.200000
75%         9.200000        0.640000    0.420000        2.600000
max        15.900000        1.580000    1.000000       15.500000

      chlorides  free sulfur dioxide  total sulfur dioxide      density \
count    1599.000000      1599.000000      1599.000000  1599.000000
mean       0.087467      15.874922      46.467792    0.996747
std        0.047065      10.460157      32.895324    0.001887
min         0.012000        1.000000        6.000000    0.990070
25%         0.070000        7.000000      22.000000    0.995600
50%         0.079000      14.000000      38.000000    0.996750
75%         0.090000      21.000000      62.000000    0.997835
max         0.611000      72.000000     289.000000    1.003690

      pH  sulphates  alcohol  quality
count    1599.000000  1599.000000  1599.000000  1599.000000
mean       3.311113    0.658149   10.422983    5.636023
std        0.154386    0.169507    1.065668    0.807569
min         2.740000    0.330000    8.400000    3.000000
25%         3.210000    0.550000    9.500000    5.000000
50%         3.310000    0.620000   10.200000    6.000000
75%         3.400000    0.730000   11.100000    6.000000
max         4.010000    2.000000   14.900000    8.000000
```

```
[9]: df.corr()
```

```
[9]:      fixed acidity  volatile acidity  citric acid \
fixed acidity      1.000000      -0.256131    0.671703
volatile acidity  -0.256131      1.000000   -0.552496
citric acid       0.671703     -0.552496    1.000000
residual sugar    0.114777     0.001918    0.143577
chlorides         0.093705     0.061298    0.203823
free sulfur dioxide -0.153794   -0.010504   -0.060978
total sulfur dioxide -0.113181    0.076470    0.035533
density          0.668047     0.022026    0.364947
pH               -0.682978     0.234937   -0.541904
sulphates        0.183006     -0.260987    0.312770
alcohol          -0.061668     -0.202288    0.109903
quality          0.124052     -0.390558    0.226373

      residual sugar  chlorides  free sulfur dioxide \
```

fixed acidity	0.114777	0.093705	-0.153794
volatile acidity	0.001918	0.061298	-0.010504
citric acid	0.143577	0.203823	-0.060978
residual sugar	1.000000	0.055610	0.187049
chlorides	0.055610	1.000000	0.005562
free sulfur dioxide	0.187049	0.005562	1.000000
total sulfur dioxide	0.203028	0.047400	0.667666
density	0.355283	0.200632	-0.021946
pH	-0.085652	-0.265026	0.070377
sulphates	0.005527	0.371260	0.051658
alcohol	0.042075	-0.221141	-0.069408
quality	0.013732	-0.128907	-0.050656

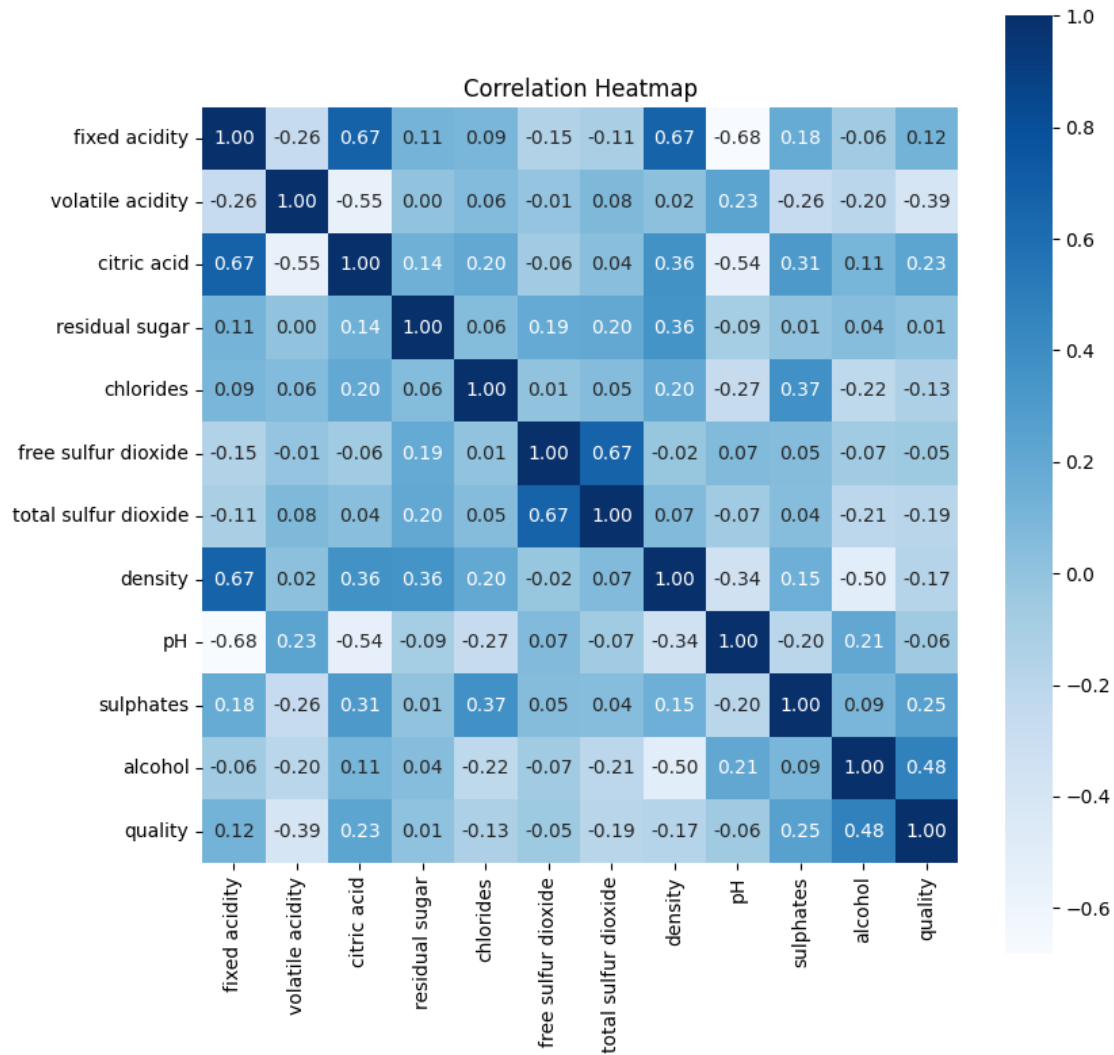
	total sulfur dioxide	density	pH	sulphates \
fixed acidity	-0.113181	0.668047	-0.682978	0.183006
volatile acidity	0.076470	0.022026	0.234937	-0.260987
citric acid	0.035533	0.364947	-0.541904	0.312770
residual sugar	0.203028	0.355283	-0.085652	0.005527
chlorides	0.047400	0.200632	-0.265026	0.371260
free sulfur dioxide	0.667666	-0.021946	0.070377	0.051658
total sulfur dioxide	1.000000	0.071269	-0.066495	0.042947
density	0.071269	1.000000	-0.341699	0.148506
pH	-0.066495	-0.341699	1.000000	-0.196648
sulphates	0.042947	0.148506	-0.196648	1.000000
alcohol	-0.205654	-0.496180	0.205633	0.093595
quality	-0.185100	-0.174919	-0.057731	0.251397

	alcohol	quality
fixed acidity	-0.061668	0.124052
volatile acidity	-0.202288	-0.390558
citric acid	0.109903	0.226373
residual sugar	0.042075	0.013732
chlorides	-0.221141	-0.128907
free sulfur dioxide	-0.069408	-0.050656
total sulfur dioxide	-0.205654	-0.185100
density	-0.496180	-0.174919
pH	0.205633	-0.057731
sulphates	0.093595	0.251397
alcohol	1.000000	0.476166
quality	0.476166	1.000000

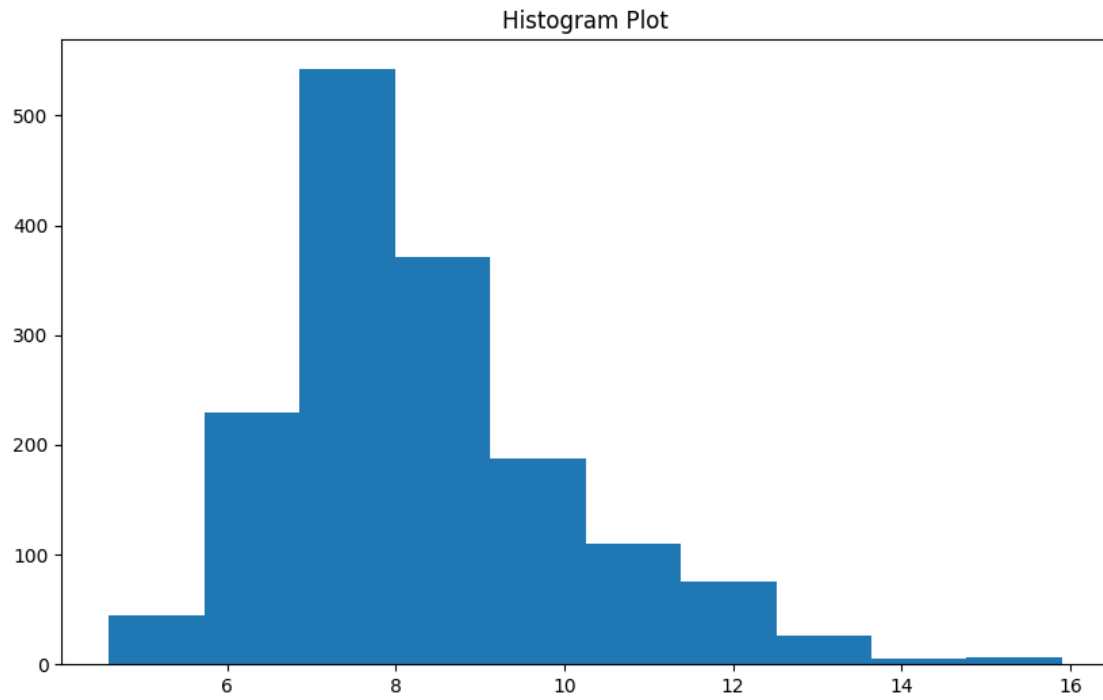
```
[10]: corr_data = df.corr()
```

```
[11]: plt.figure(figsize = (9, 9))
sns.heatmap(corr_data, cbar = True, square= True, annot=True, fmt= '.2f',
            cmap='Blues')
plt.title('Correlation Heatmap')
```

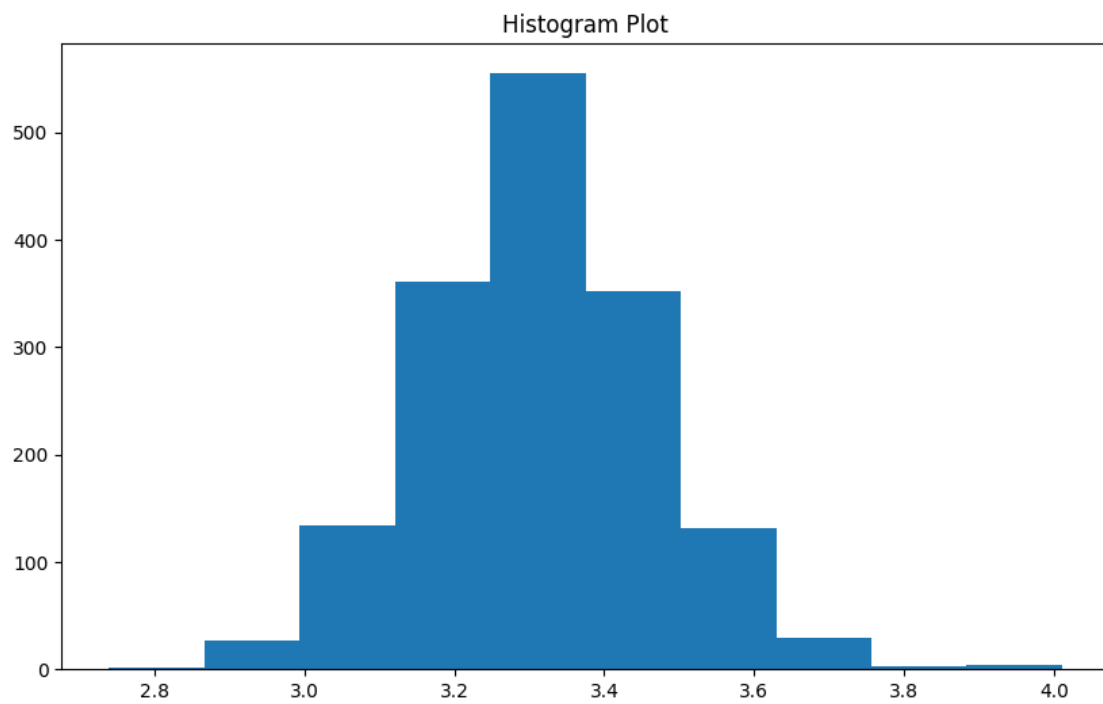
```
[11]: Text(0.5, 1.0, 'Correlation Heatmap')
```



```
[26]: def hist_plots(df):
        plt.figure(figsize=(10, 6))
        plt.hist(df)
        plt.title("Histogram Plot")
        plt.show()
        hist_plots(df['fixed acidity'])
```

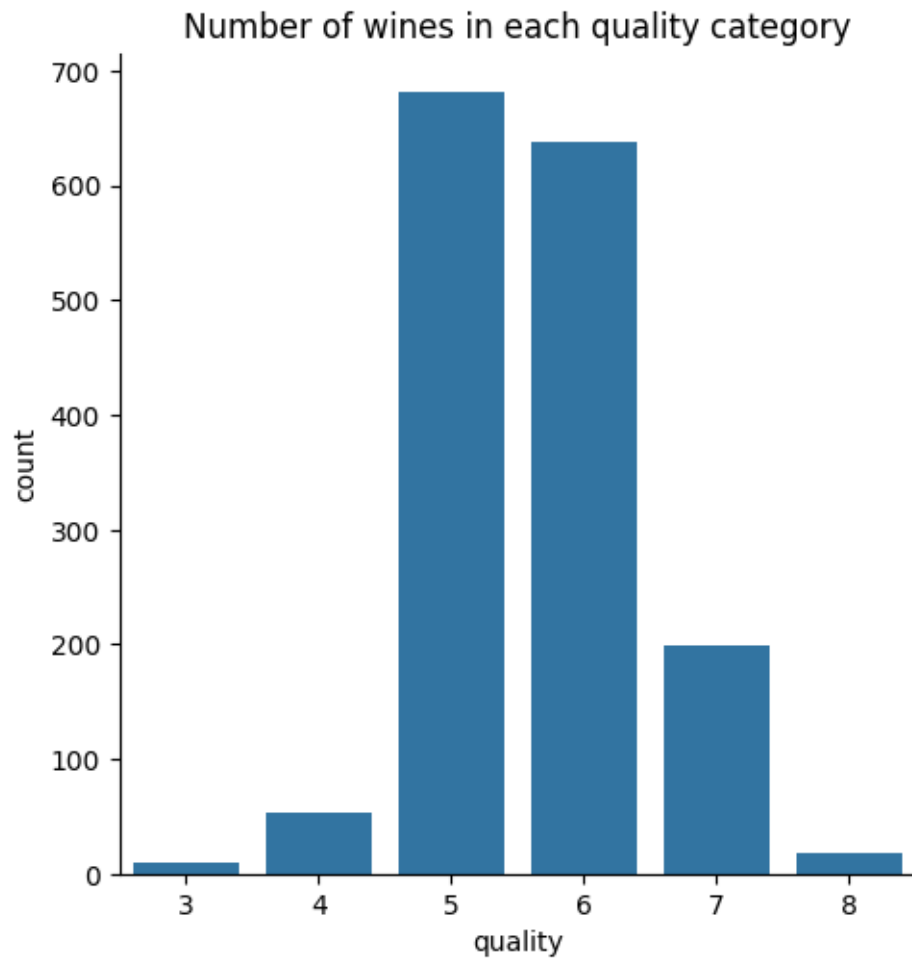


```
[27]: hist_plots(df['pH'])
```

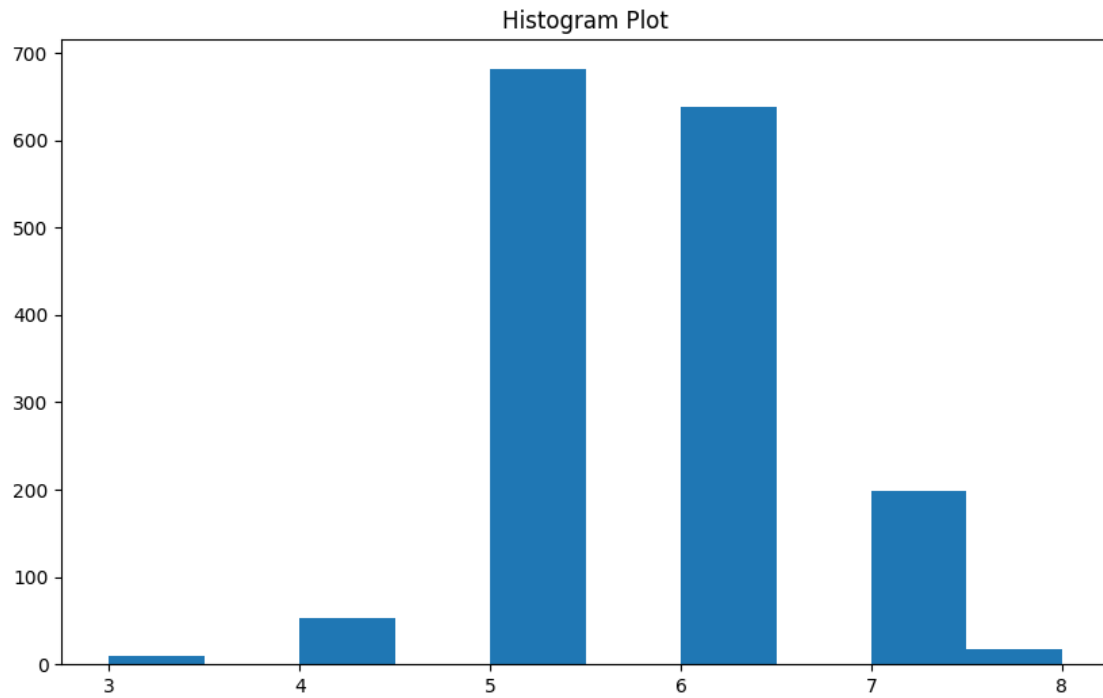


```
[12]: # Number of wines in each quality category
sns.catplot(x='quality', data=df, kind='count')
plt.title('Number of wines in each quality category')
```

```
[12]: Text(0.5, 1.0, 'Number of wines in each quality category')
```

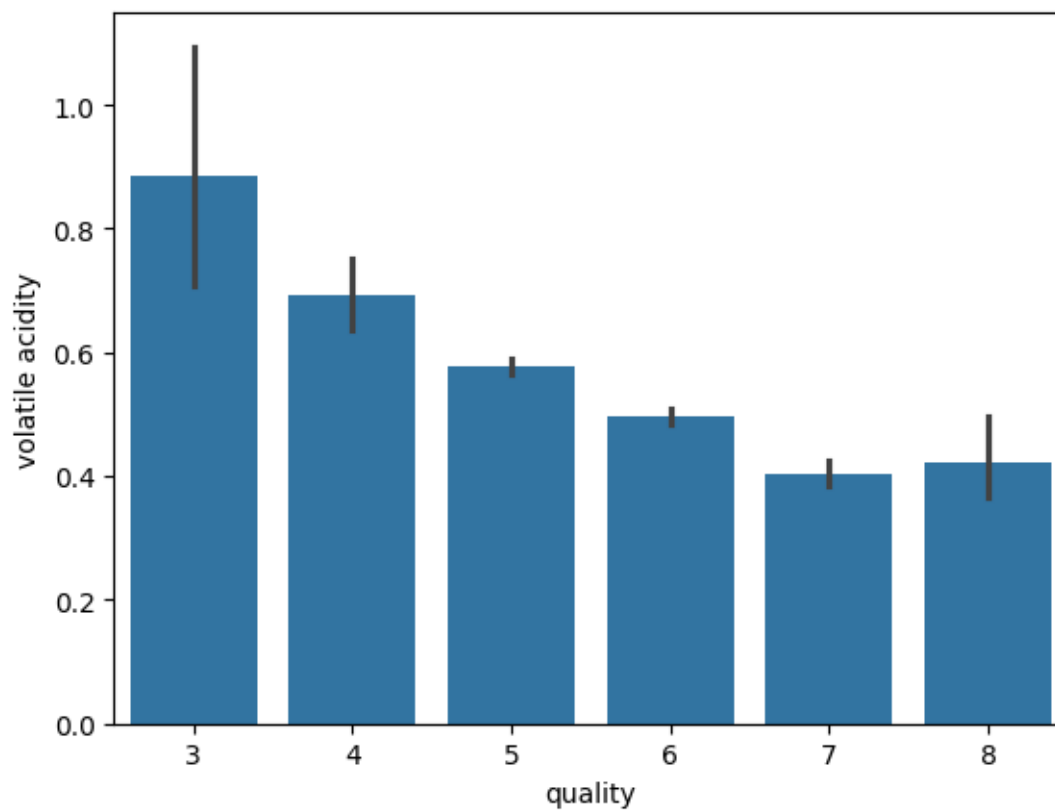


```
[28]: hist_plots(df['quality'])
```

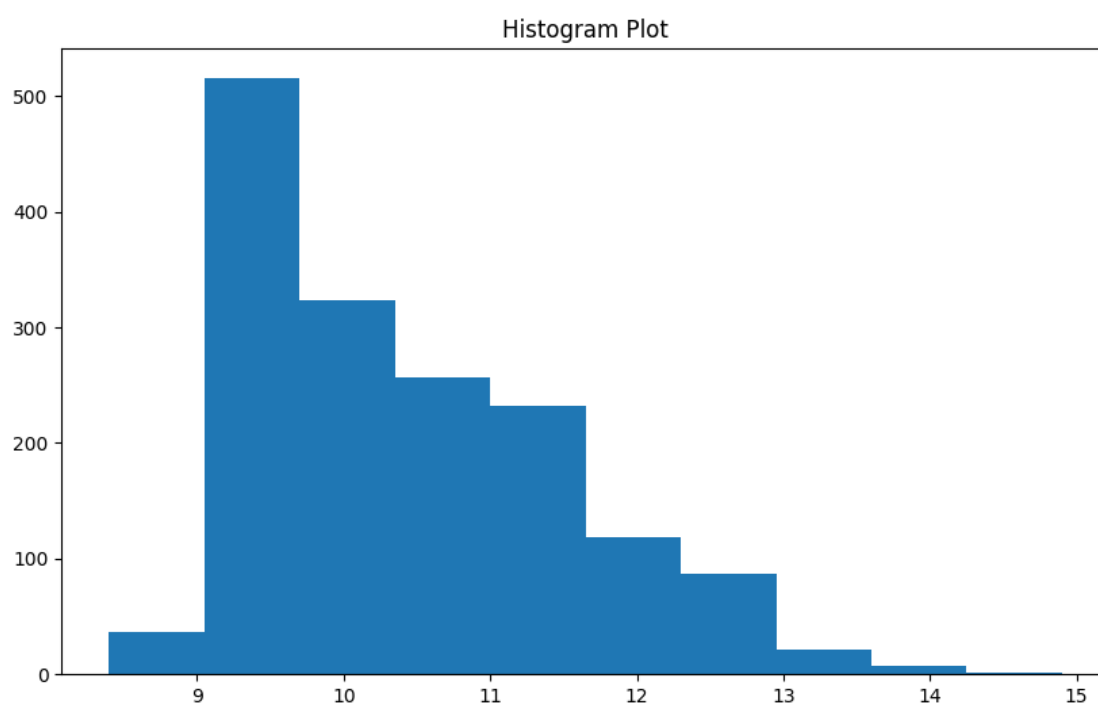


```
[13]: # plotting a barplot for quality vs volatile acidity  
sns.barplot(x = 'quality', y= 'volatile acidity', data = df)
```

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

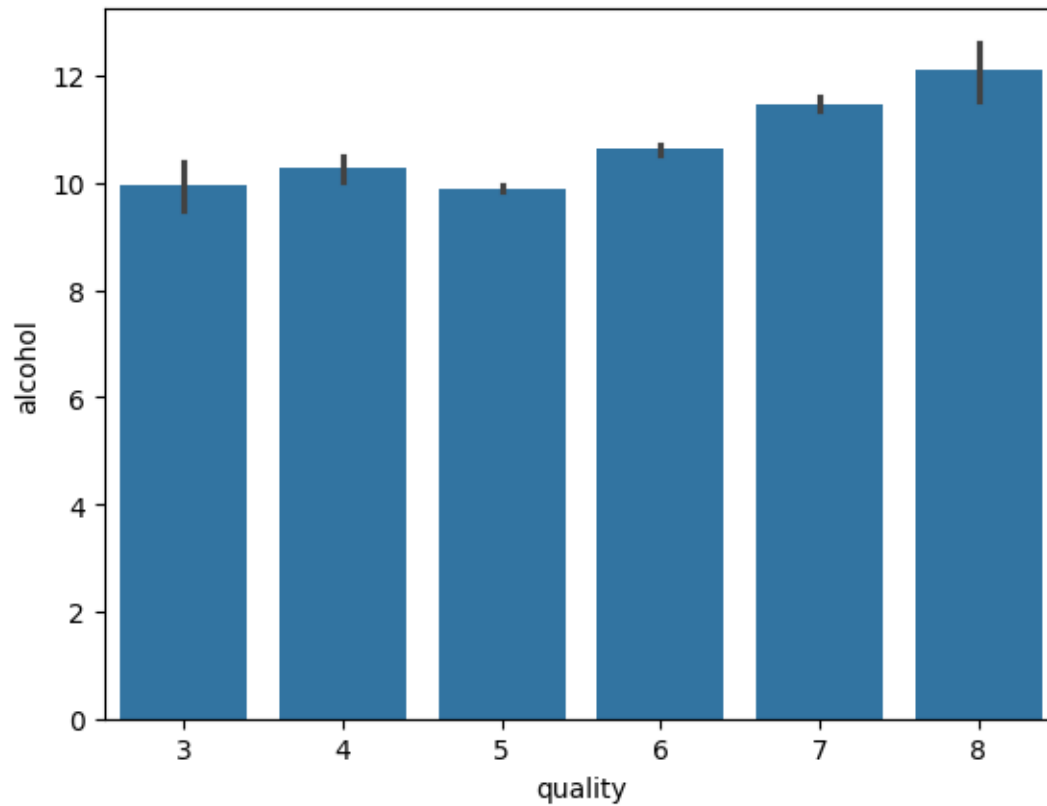



```
[29]: hist_plots(df['alcohol'])
```

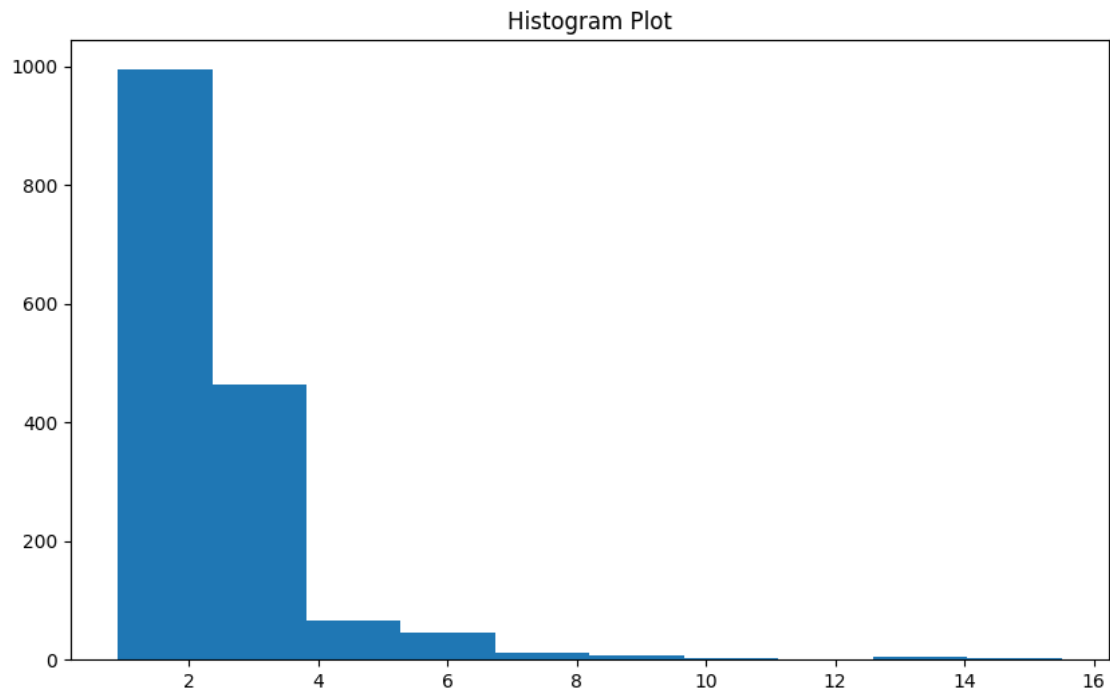


```
[14]: # plotting a barplot for quality vs alcohol
sns.barplot(x = 'quality', y = 'alcohol', data = df)
```

```
[14]: <Axes: xlabel='quality', ylabel='alcohol'>
```

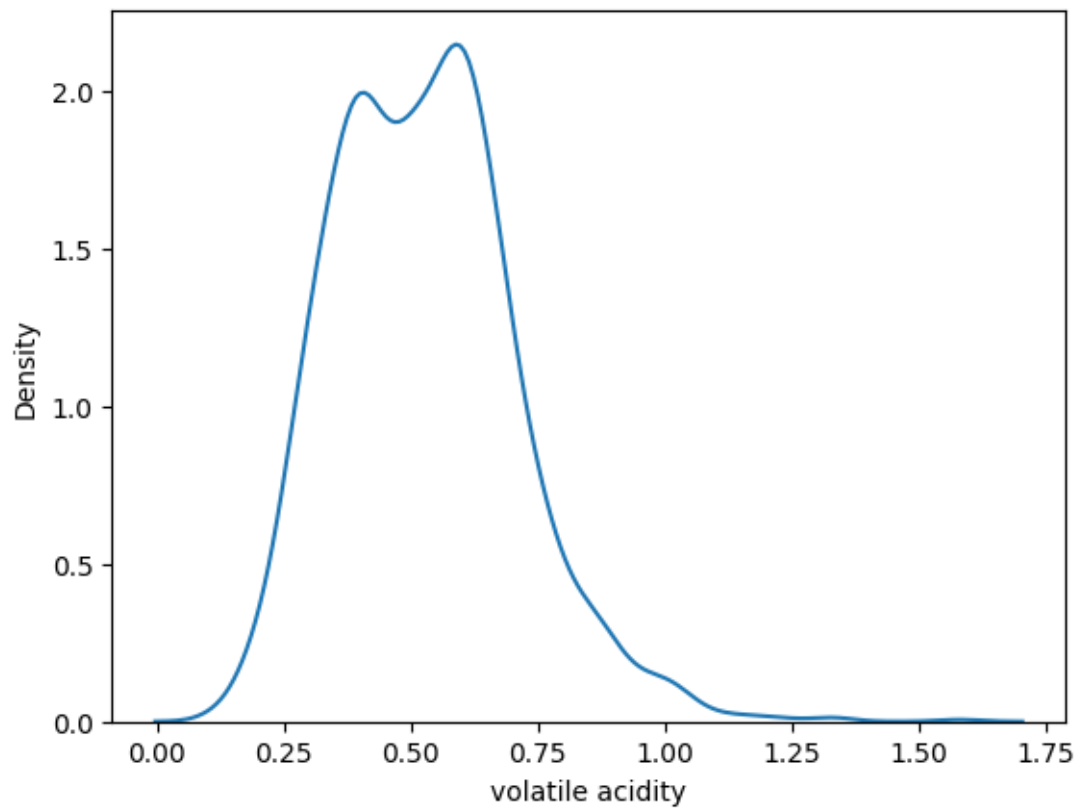


```
[30]: hist_plots(df['residual sugar'])
```



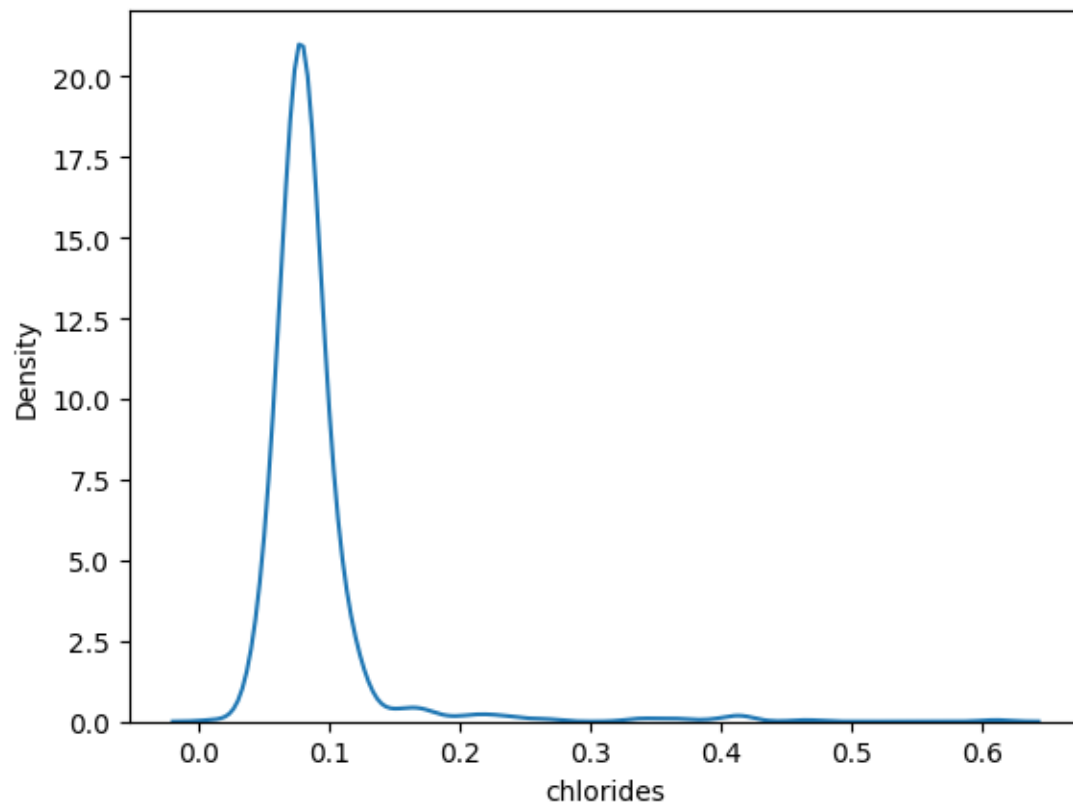
```
[31]: sns.kdeplot(df['volatile acidity'])
```

```
[31]: <Axes: xlabel='volatile acidity', ylabel='Density'>
```



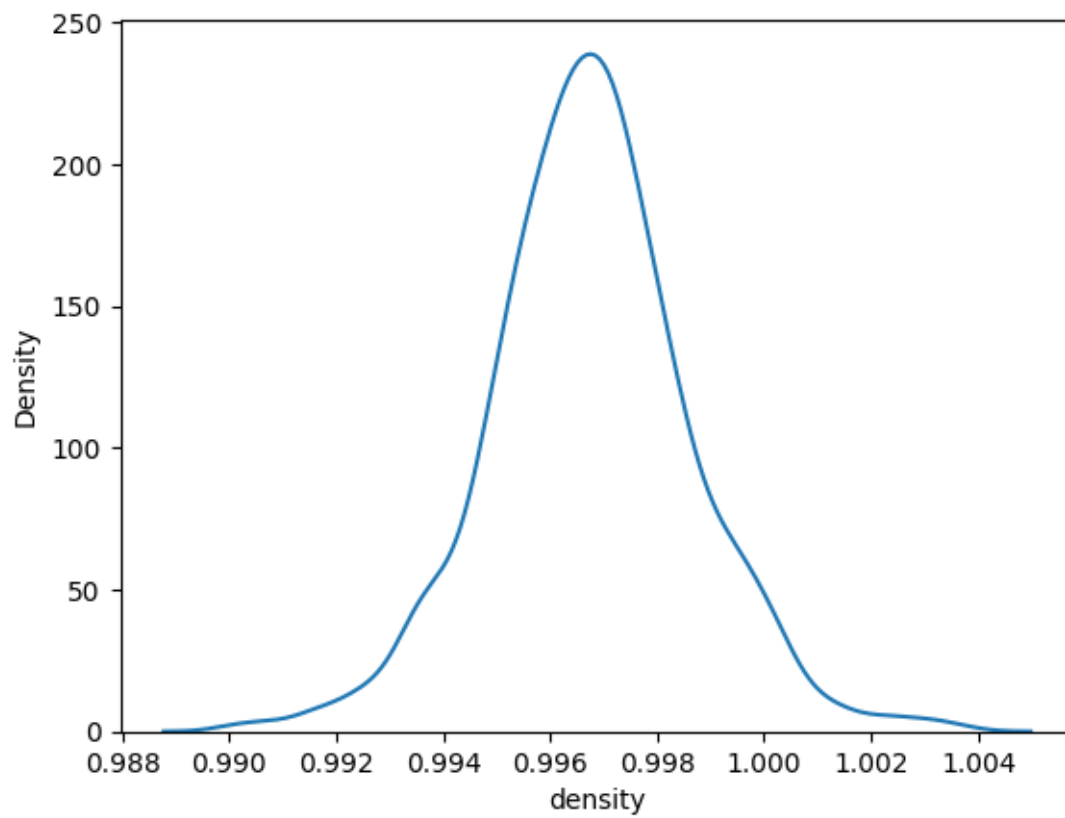
```
[33]: sns.kdeplot(df['chlorides'])
```

```
[33]: <Axes: xlabel='chlorides', ylabel='Density'>
```



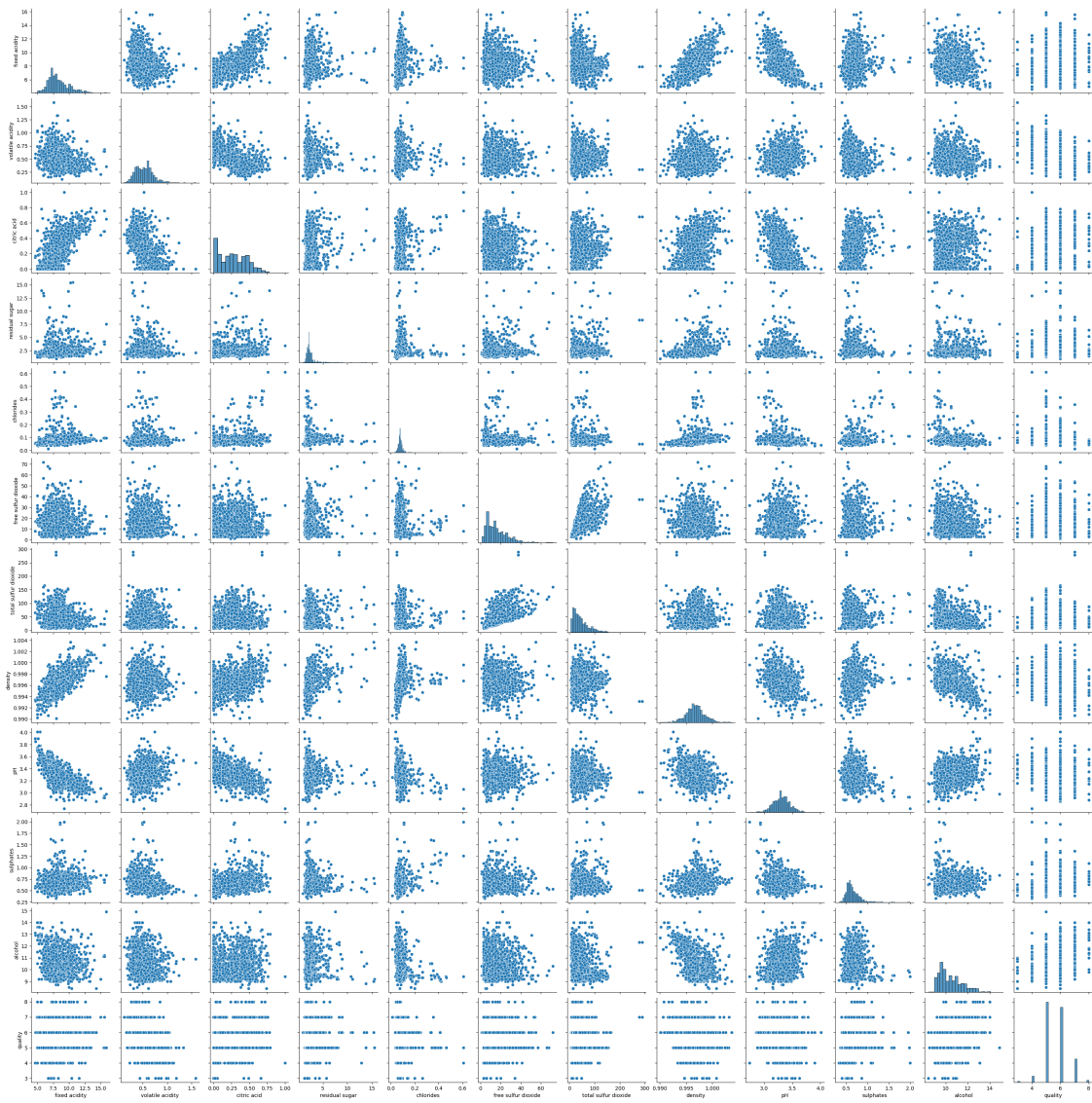
```
[34]: sns.kdeplot(df['density'])
```

```
[34]: <Axes: xlabel='density', ylabel='Density'>
```

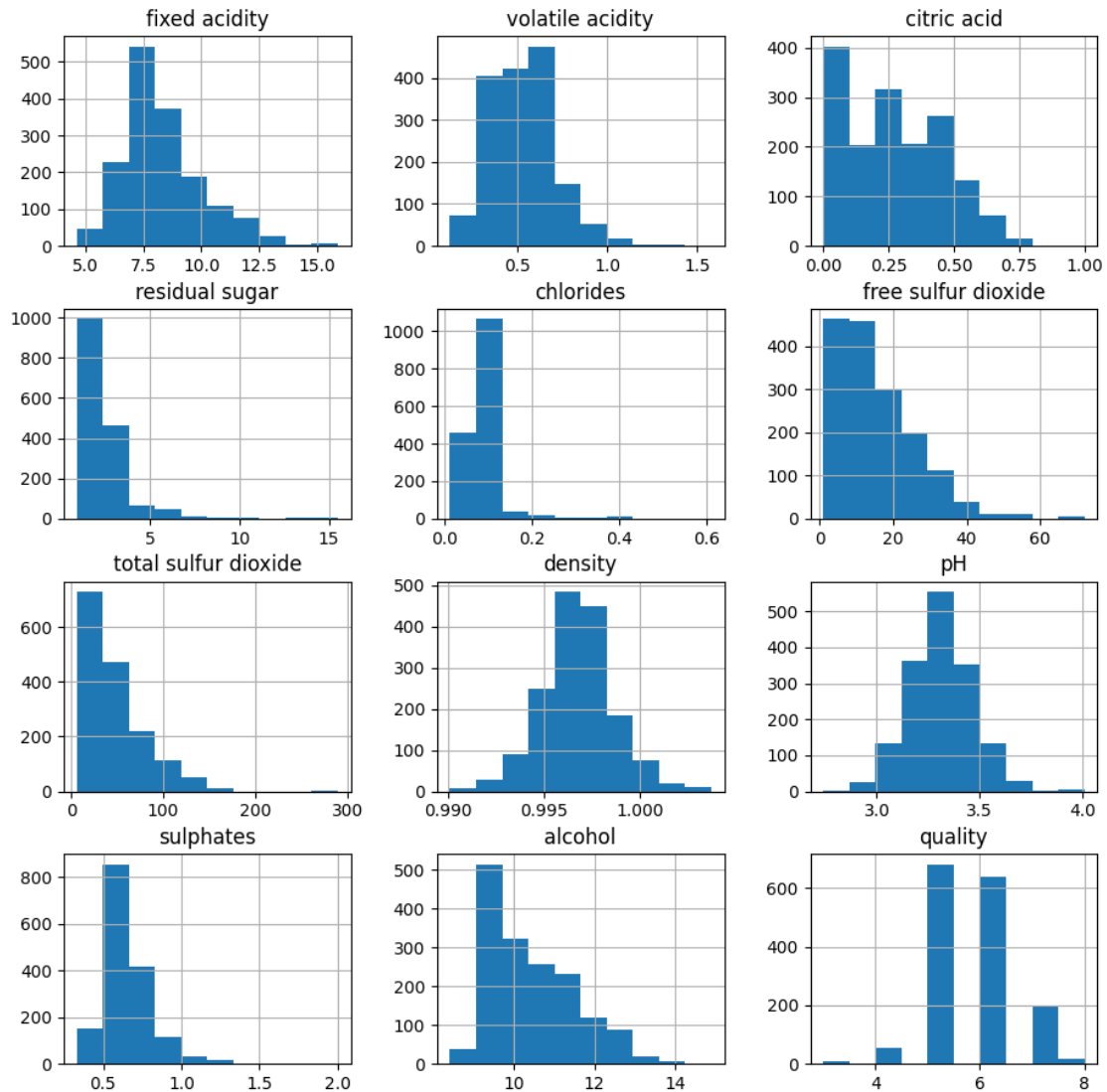


```
[35]: plt.figure(figsize = (12,6))  
sns.pairplot(df)  
plt.show()
```

<Figure size 1200x600 with 0 Axes>



```
[37]: df.hist(figsize=(11, 11))
plt.show()
```



```
[ ]:
```

```
[15]: # Separating the features and labels
```

```
X = df.drop('quality', axis= 1)
y = df['quality'].apply(lambda y_value: 1 if y_value >= 7 else 0)
```

```
[16]: y.value_counts()
```

```
[16]: quality
0    1382
1     217
Name: count, dtype: int64
```



```
[17]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
↳ random_state= 3)
```

```
[18]: # Training The Model
```

```
model = RandomForestClassifier(n_estimators=100, max_depth=5, random_state=1)
```

```
[19]: model.fit(X_train, y_train)
```

```
[19]: RandomForestClassifier(max_depth=5, random_state=1)
```

```
[20]: from sklearn.metrics import accuracy_score
```

```
[21]: #accuracy on test data
```

```
X_test_preds = model.predict(X_test)
```

```
test_accuracy = accuracy_score(y_test, X_test_preds)
```

```
[22]: print("Test accuracy: {:.2f}%".format(test_accuracy * 100))
```

```
Test accuracy: 91.56%
```

```
[ ]:
```