# EDA With Red Wine Data

**Data Set Information:**

The two datasets are related to red and white variants of the Portuguese "Vinho Verde" wine.

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 many more normal wines than excellent or poor ones).

Outlier detection algorithms could be used to detect the few excellent or poor wines. Also, we are not sure if all input variables are relevant. So it could be interesting to test feature selection methods.

**Attribute Information:**

Input variables (based on physicochemical tests): 1 - fixed acidity

2 - volatile acidity

3 - citric acid

4 - residual sugar

5 - chlorides

6 - free sulfur dioxide

7 - total sulfur dioxide

8 - density

9 - pH

10 - sulphates

11 - alcohol

Output variable (based on sensory data):

12 - quality (score between 0 and 10)

In [1]:
```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

In [2]:
```python
df=pd.read_csv('/content/sample_data/winequality-red.csv')
```

In [3]:
```python
df.head()
```

Out[3]:

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 |
| 1 | 7.8 | 0.88 | 0.00 | 2.6 | 0.098 | 25.0 | 67.0 | 0.9968 | 3.20 | 0.68 | 9.8 | 5 |
| 2 | 7.8 | 0.76 | 0.04 | 2.3 | 0.092 | 15.0 | 54.0 | 0.9970 | 3.26 | 0.65 | 9.8 | 5 |
| 3 | 11.2 | 0.28 | 0.56 | 1.9 | 0.075 | 17.0 | 60.0 | 0.9980 | 3.16 | 0.58 | 9.8 | 6 |
| 4 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 |

In [4]:
```python
## Summary of the dataset
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
```

In [5]:
```python
## descriptive summary of the dataset
df.describe()
```

Out[5]:

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol |
|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 |
| mean | 8.319637 | 0.527821 | 0.270976 | 2.538806 | 0.087467 | 15.874922 | 46.467792 | 0.996747 | 3.311113 | 0.658149 | 10.422983 |
| std | 1.741096 | 0.179060 | 0.194801 | 1.409928 | 0.047065 | 10.460157 | 32.895324 | 0.001887 | 0.154386 | 0.169507 | 1.065668 |
| min | 4.600000 | 0.120000 | 0.000000 | 0.900000 | 0.012000 | 1.000000 | 6.000000 | 0.990070 | 2.740000 | 0.330000 | 8.400000 |
| 25% | 7.100000 | 0.390000 | 0.090000 | 1.900000 | 0.070000 | 7.000000 | 22.000000 | 0.995600 | 3.210000 | 0.550000 | 9.500000 |
| 50% | 7.900000 | 0.520000 | 0.260000 | 2.200000 | 0.079000 | 14.000000 | 38.000000 | 0.996750 | 3.310000 | 0.620000 | 10.200000 |
| 75% | 9.200000 | 0.640000 | 0.420000 | 2.600000 | 0.090000 | 21.000000 | 62.000000 | 0.997835 | 3.400000 | 0.730000 | 11.100000 |
| max | 15.900000 | 1.580000 | 1.000000 | 15.500000 | 0.611000 | 72.000000 | 289.000000 | 1.003690 | 4.010000 | 2.000000 | 14.900000 |

In [6]:
```python
df.ndim
```

Out[6]: 2

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

Out[7]: (1599, 12)

In [8]: 
```python
df.size
```

Out[8]: 19188

In [9]: 
```python
## List down all the columns names
df.columns
```

Out[9]: 
```
Index(['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
       'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
       'pH', 'sulphates', 'alcohol', 'quality'],
      dtype='object')
```

In [10]: 
```python
df.nunique() # unique value in every column
```

Out[10]:

|  | **0** |
|---|---|
| **fixed acidity** | 96 |
| **volatile acidity** | 143 |
| **citric acid** | 80 |
| **residual sugar** | 91 |
| **chlorides** | 153 |
| **free sulfur dioxide** | 60 |
| **total sulfur dioxide** | 144 |
| **density** | 436 |
| **pH** | 89 |
| **sulphates** | 96 |
| **alcohol** | 65 |
| **quality** | 6 |

**dtype:** int64

In [11]:
```python
df['quality'].unique()
```

Out[11]:
```
array([5, 6, 7, 4, 8, 3])
```

In [12]:
```python
## Missing values in the dataset

df.isnull()
```

Out[12]:

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False | False | False | False | False | False |
| 1 | False | False | False | False | False | False | False | False | False | False | False | False |
| 2 | False | False | False | False | False | False | False | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False | False | False | False | False |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 1594 | False | False | False | False | False | False | False | False | False | False | False | False |
| 1595 | False | False | False | False | False | False | False | False | False | False | False | False |
| 1596 | False | False | False | False | False | False | False | False | False | False | False | False |
| 1597 | False | False | False | False | False | False | False | False | False | False | False | False |
| 1598 | False | False | False | False | False | False | False | False | False | False | False | False |

1599 rows × 12 columns

In [13]:
```python
df.isnull().sum()
```

Out[13]:

|  | 0 |
|---|---|
| **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 [14]:
```python
## Duplicate records

df[df.duplicated()]
```
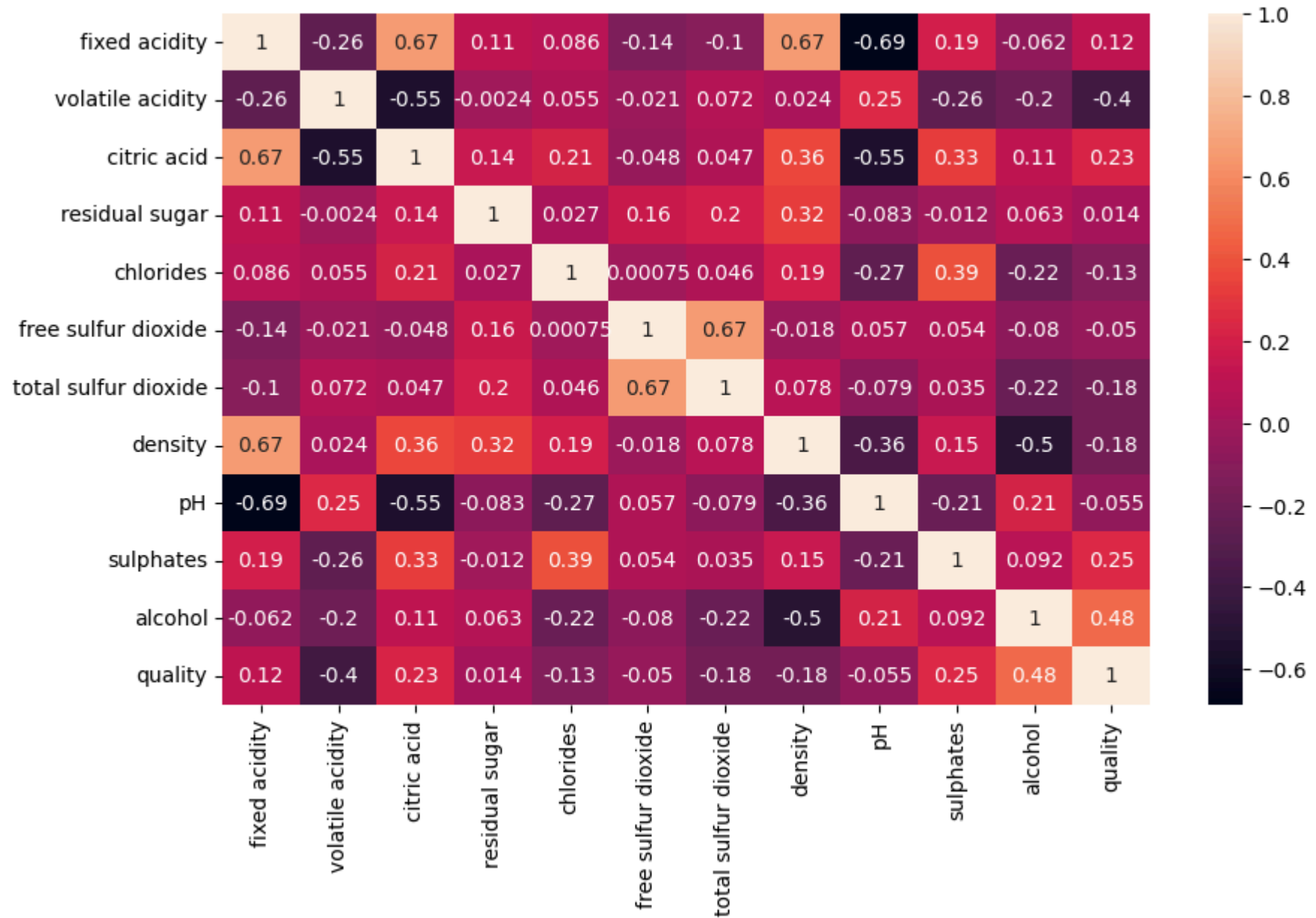
Out[14]:

| | 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 |
| **11** | 7.5 | 0.500 | 0.36 | 6.10 | 0.071 | 17.0 | 102.0 | 0.99780 | 3.35 | 0.80 | 10.5 | 5 |
| **27** | 7.9 | 0.430 | 0.21 | 1.60 | 0.106 | 10.0 | 37.0 | 0.99660 | 3.17 | 0.91 | 9.5 | 5 |
| **40** | 7.3 | 0.450 | 0.36 | 5.90 | 0.074 | 12.0 | 87.0 | 0.99780 | 3.33 | 0.83 | 10.5 | 5 |
| **65** | 7.2 | 0.725 | 0.05 | 4.65 | 0.086 | 4.0 | 11.0 | 0.99620 | 3.41 | 0.39 | 10.9 | 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 |
| **1581** | 6.2 | 0.560 | 0.09 | 1.70 | 0.053 | 24.0 | 32.0 | 0.99402 | 3.54 | 0.60 | 11.3 | 5 |
| **1596** | 6.3 | 0.510 | 0.13 | 2.30 | 0.076 | 29.0 | 40.0 | 0.99574 | 3.42 | 0.75 | 11.0 | 6 |

240 rows × 12 columns

In [16]:
```python
df.duplicated().sum()
```

Out[16]: 240

In [17]:
```python
## Remove the duplicates

df.drop_duplicates(inplace=True)
```

In [18]:
```python
df.shape
```

Out[18]: (1359, 12)

In [19]:
```python
## Correlation

df.corr()
```

Out[19]:

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| fixed acidity | 1.000000 | -0.255124 | 0.667437 | 0.111025 | 0.085886 | -0.140580 | -0.103777 | 0.670195 | -0.686685 | 0.190269 | -0.061596 | 0.119024 |
| volatile acidity | -0.255124 | 1.000000 | -0.551248 | -0.002449 | 0.055154 | -0.020945 | 0.071701 | 0.023943 | 0.247111 | -0.256948 | -0.197812 | -0.395214 |
| citric acid | 0.667437 | -0.551248 | 1.000000 | 0.143892 | 0.210195 | -0.048004 | 0.047358 | 0.357962 | -0.550310 | 0.326062 | 0.105108 | 0.228057 |
| residual sugar | 0.111025 | -0.002449 | 0.143892 | 1.000000 | 0.026656 | 0.160527 | 0.201038 | 0.324522 | -0.083143 | -0.011837 | 0.063281 | 0.013640 |
| chlorides | 0.085886 | 0.055154 | 0.210195 | 0.026656 | 1.000000 | 0.000749 | 0.045773 | 0.193592 | -0.270893 | 0.394557 | -0.223824 | -0.130988 |
| free sulfur dioxide | -0.140580 | -0.020945 | -0.048004 | 0.160527 | 0.000749 | 1.000000 | 0.667246 | -0.018071 | 0.056631 | 0.054126 | -0.080125 | -0.050463 |
| total sulfur dioxide | -0.103777 | 0.071701 | 0.047358 | 0.201038 | 0.045773 | 0.667246 | 1.000000 | 0.078141 | -0.079257 | 0.035291 | -0.217829 | -0.177855 |
| density | 0.670195 | 0.023943 | 0.357962 | 0.324522 | 0.193592 | -0.018071 | 0.078141 | 1.000000 | -0.355617 | 0.146036 | -0.504995 | -0.184252 |
| pH | -0.686685 | 0.247111 | -0.550310 | -0.083143 | -0.270893 | 0.056631 | -0.079257 | -0.355617 | 1.000000 | -0.214134 | 0.213418 | -0.055245 |
| sulphates | 0.190269 | -0.256948 | 0.326062 | -0.011837 | 0.394557 | 0.054126 | 0.035291 | 0.146036 | -0.214134 | 1.000000 | 0.091621 | 0.248835 |
| alcohol | -0.061596 | -0.197812 | 0.105108 | 0.063281 | -0.223824 | -0.080125 | -0.217829 | -0.504995 | 0.213418 | 0.091621 | 1.000000 | 0.480343 |
| quality | 0.119024 | -0.395214 | 0.228057 | 0.013640 | -0.130988 | -0.050463 | -0.177855 | -0.184252 | -0.055245 | 0.248835 | 0.480343 | 1.000000 |

In [20]:
```python
plt.figure(figsize=(10,6))
sns.heatmap(df.corr(),annot=True)
```
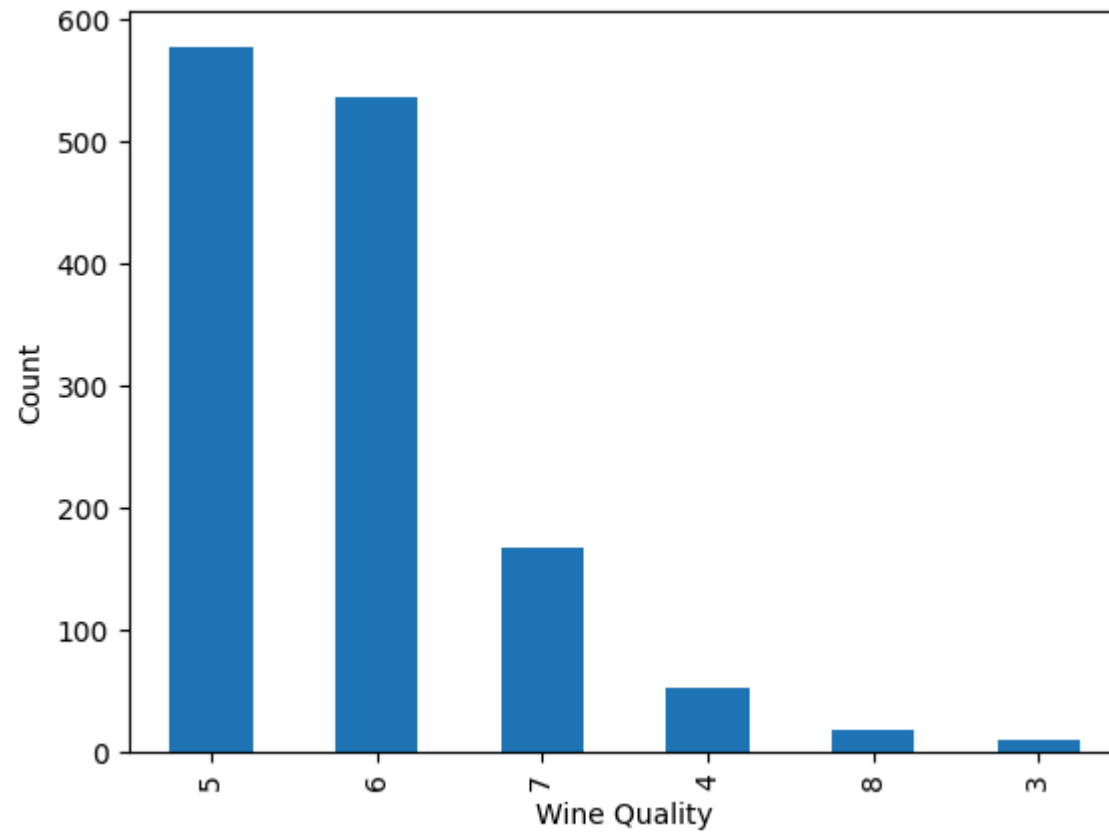
Out[20]: <Axes: >

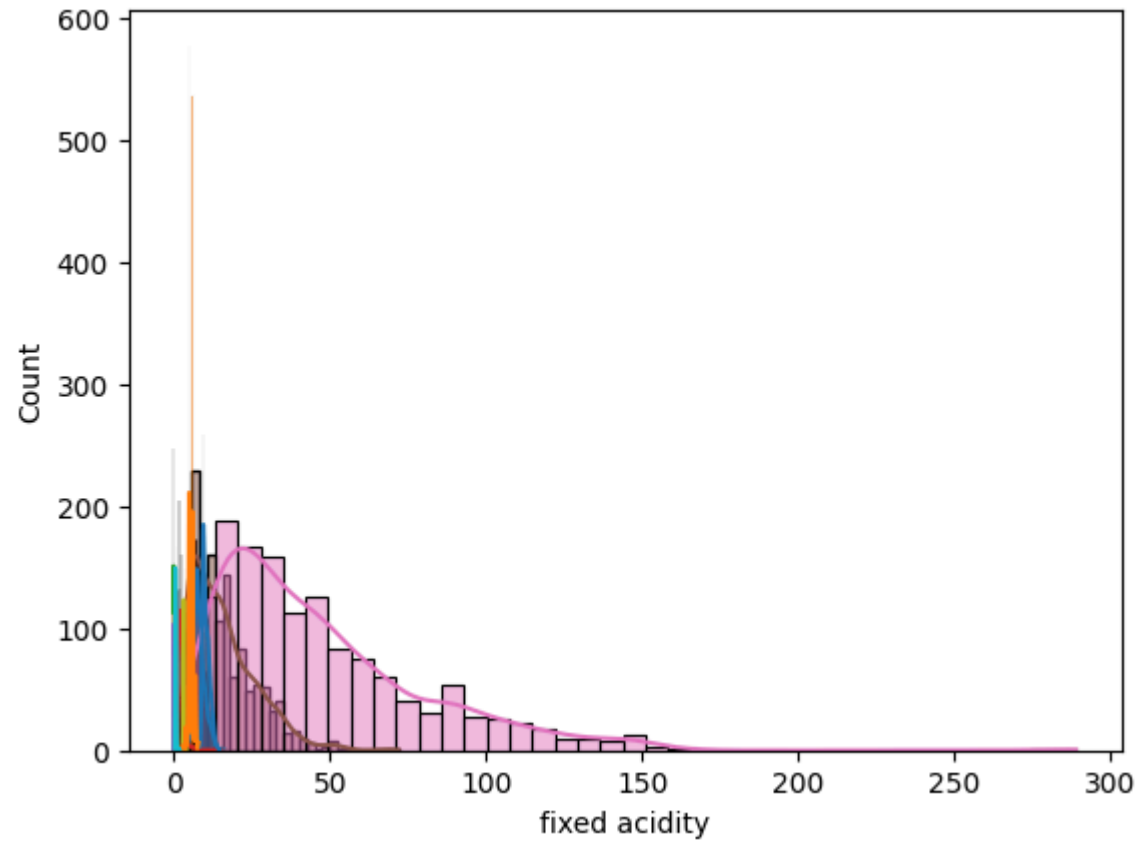## Visualization

In [21]:
```python
df.quality.value_counts().plot(kind="bar")
plt.xlabel("Wine Quality")
plt.ylabel("Count")
plt.show()
```
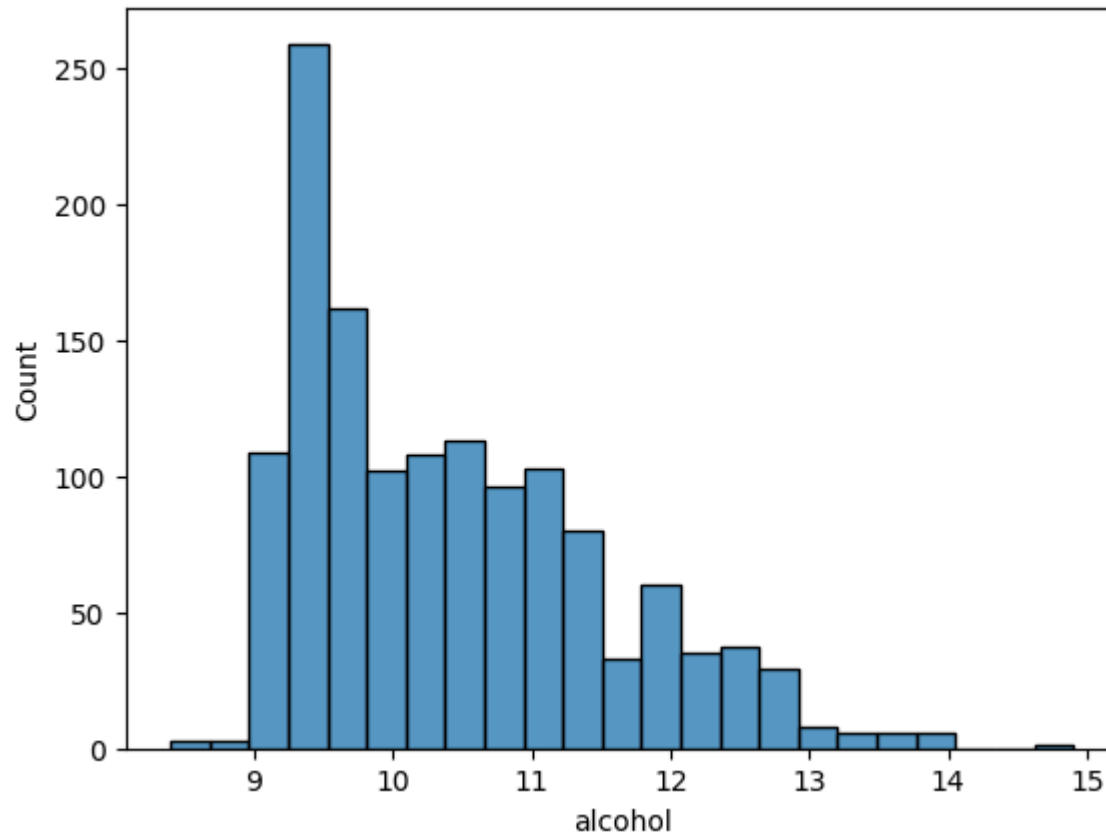


In [22]:
```python
for column in df.columns:
    sns.histplot(df[column],kde=True)
```
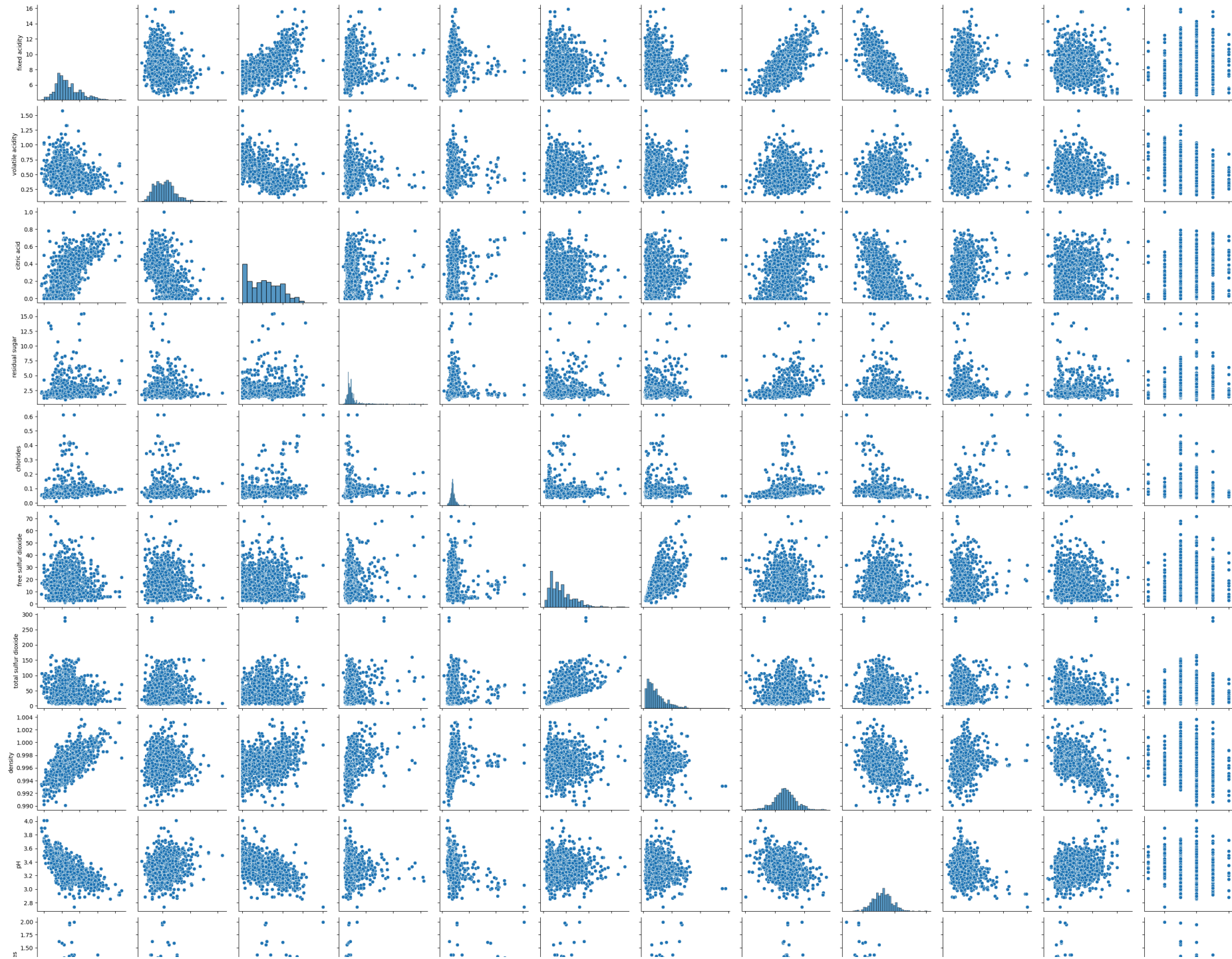
```
In [23]: sns.histplot(df['alcohol'])
```
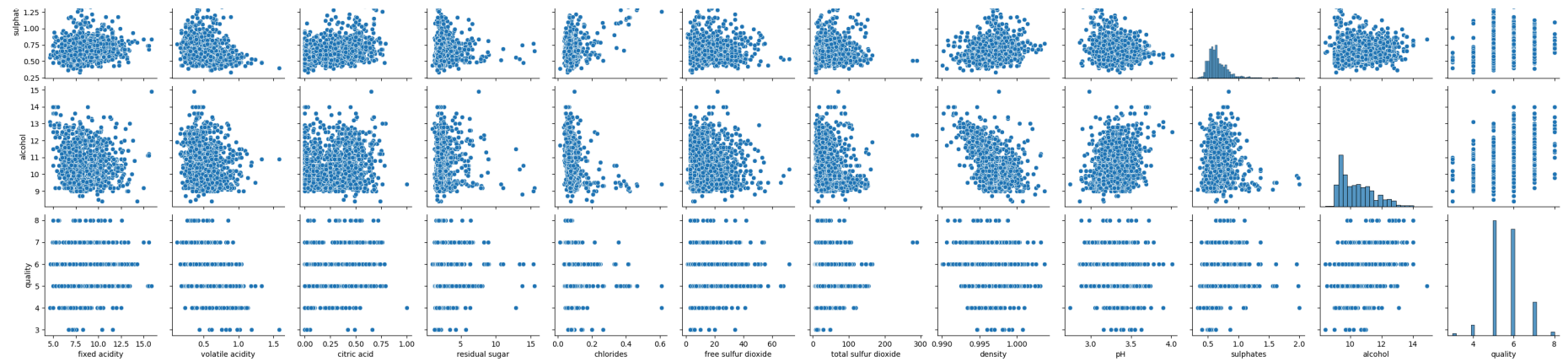
```
Out[23]: <Axes: xlabel='alcohol', ylabel='Count'>
```

In [25]: `#univariate,bivariate,multivariate analysis`
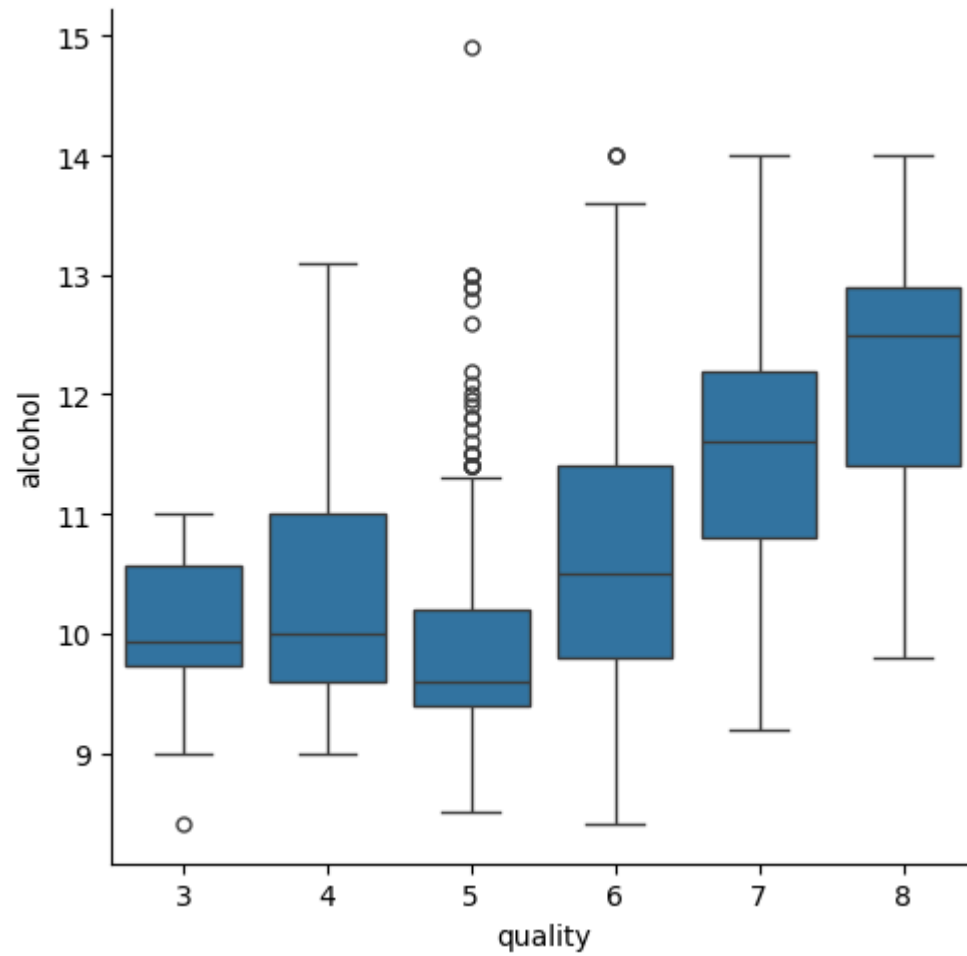`plt.figure(dpi=70)`
`sns.pairplot(df)`

Out[25]: `<seaborn.axisgrid.PairGrid at 0x7d70c811e9e0>`

`<Figure size 448x336 with 0 Axes>`

In [26]: ```python
##categorical Plot

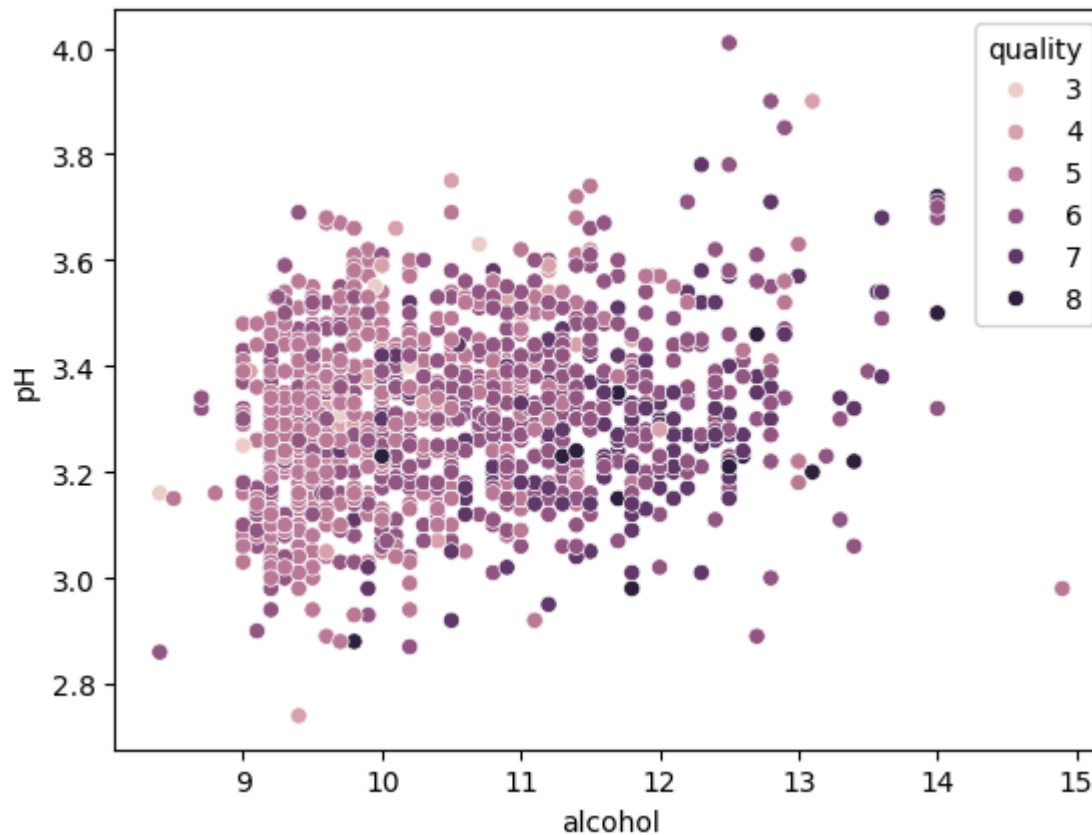sns.catplot(x='quality', y='alcohol', data=df, kind="box")
```

Out[26]: `<seaborn.axisgrid.FacetGrid at 0x7d70c74c9480>`

In [27]: `df.head()`

Out[27]:

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 |
| 1 | 7.8 | 0.88 | 0.00 | 2.6 | 0.098 | 25.0 | 67.0 | 0.9968 | 3.20 | 0.68 | 9.8 | 5 |
| 2 | 7.8 | 0.76 | 0.04 | 2.3 | 0.092 | 15.0 | 54.0 | 0.9970 | 3.26 | 0.65 | 9.8 | 5 |
| 3 | 11.2 | 0.28 | 0.56 | 1.9 | 0.075 | 17.0 | 60.0 | 0.9980 | 3.16 | 0.58 | 9.8 | 6 |
| 5 | 7.4 | 0.66 | 0.00 | 1.8 | 0.075 | 13.0 | 40.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 |

In [28]: 
```python
sns.scatterplot(x='alcohol',y='pH',hue='quality',data=df)
```

Out[28]: `<Axes: xlabel='alcohol', ylabel='pH'>`

**Conclusion:**

- Alcohol content, volatile acidity, and sulfur dioxide levels were the most significant factors impacting wine quality.

- Sweetness (residual sugar) and pH were not as crucial in determining the quality score, based on this dataset.

- Wines with higher alcohol content, lower volatile acidity, and balanced sulfur dioxide levels generally received better ratings.

In [29]:
```
!jupyter nbconvert --to html /content/Red_Wine_EDA.ipynb
```

```
[NbConvertApp] Converting notebook /content/Red_Wine_EDA.ipynb to html
[NbConvertApp] Writing 3831319 bytes to /content/Red_Wine_EDA.html
```

In [ ]: