# Red Wine Quality Prediction Project:

```
In [2]:   # importing data sets library
          import pandas as pd
          import numpy as np
          import matplotlib.pyplot as plt
          import seaborn as sns
          from sklearn.model_selection import train_test_split
          from sklearn.ensemble import RandomForestClassifier
          from sklearn.metrics import accuracy_score
```

```
In [3]:   import warnings
          warnings.filterwarnings('ignore')
```

```
In [4]:   df = pd.read_csv('https://raw.githubusercontent.com/dsrscientist/DSData/master/wine
          df
```

Out[4]:

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcoh |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7.4 | 0.700 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.99780 | 3.51 | 0.56 | |
| 1 | 7.8 | 0.880 | 0.00 | 2.6 | 0.098 | 25.0 | 67.0 | 0.99680 | 3.20 | 0.68 | |
| 2 | 7.8 | 0.760 | 0.04 | 2.3 | 0.092 | 15.0 | 54.0 | 0.99700 | 3.26 | 0.65 | |
| 3 | 11.2 | 0.280 | 0.56 | 1.9 | 0.075 | 17.0 | 60.0 | 0.99800 | 3.16 | 0.58 | |
| 4 | 7.4 | 0.700 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.99780 | 3.51 | 0.56 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 1594 | 6.2 | 0.600 | 0.08 | 2.0 | 0.090 | 32.0 | 44.0 | 0.99490 | 3.45 | 0.58 | 1 |
| 1595 | 5.9 | 0.550 | 0.10 | 2.2 | 0.062 | 39.0 | 51.0 | 0.99512 | 3.52 | 0.76 | 1 |
| 1596 | 6.3 | 0.510 | 0.13 | 2.3 | 0.076 | 29.0 | 40.0 | 0.99574 | 3.42 | 0.75 | 1 |
| 1597 | 5.9 | 0.645 | 0.12 | 2.0 | 0.075 | 32.0 | 44.0 | 0.99547 | 3.57 | 0.71 | 1 |
| 1598 | 6.0 | 0.310 | 0.47 | 3.6 | 0.067 | 18.0 | 42.0 | 0.99549 | 3.39 | 0.66 | 1 |

1599 rows × 12 columns

```
In [5]:   df.head(15)
```

Out[5]:

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7.4 | 0.700 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 |
| 1 | 7.8 | 0.880 | 0.00 | 2.6 | 0.098 | 25.0 | 67.0 | 0.9968 | 3.20 | 0.68 | 9.8 |
| 2 | 7.8 | 0.760 | 0.04 | 2.3 | 0.092 | 15.0 | 54.0 | 0.9970 | 3.26 | 0.65 | 9.8 |
| 3 | 11.2 | 0.280 | 0.56 | 1.9 | 0.075 | 17.0 | 60.0 | 0.9980 | 3.16 | 0.58 | 9.8 |
| 4 | 7.4 | 0.700 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 |
| 5 | 7.4 | 0.660 | 0.00 | 1.8 | 0.075 | 13.0 | 40.0 | 0.9978 | 3.51 | 0.56 | 9.4 |
| 6 | 7.9 | 0.600 | 0.06 | 1.6 | 0.069 | 15.0 | 59.0 | 0.9964 | 3.30 | 0.46 | 9.4 |
| 7 | 7.3 | 0.650 | 0.00 | 1.2 | 0.065 | 15.0 | 21.0 | 0.9946 | 3.39 | 0.47 | 10.0 |
| 8 | 7.8 | 0.580 | 0.02 | 2.0 | 0.073 | 9.0 | 18.0 | 0.9968 | 3.36 | 0.57 | 9.5 |
| 9 | 7.5 | 0.500 | 0.36 | 6.1 | 0.071 | 17.0 | 102.0 | 0.9978 | 3.35 | 0.80 | 10.5 |
| 10 | 6.7 | 0.580 | 0.08 | 1.8 | 0.097 | 15.0 | 65.0 | 0.9959 | 3.28 | 0.54 | 9.2 |
| 11 | 7.5 | 0.500 | 0.36 | 6.1 | 0.071 | 17.0 | 102.0 | 0.9978 | 3.35 | 0.80 | 10.5 |
| 12 | 5.6 | 0.615 | 0.00 | 1.6 | 0.089 | 16.0 | 59.0 | 0.9943 | 3.58 | 0.52 | 9.9 |
| 13 | 7.8 | 0.610 | 0.29 | 1.6 | 0.114 | 9.0 | 29.0 | 0.9974 | 3.26 | 1.56 | 9.1 |
| 14 | 8.9 | 0.620 | 0.18 | 3.8 | 0.176 | 52.0 | 145.0 | 0.9986 | 3.16 | 0.88 | 9.2 |

In [6]:
```python
# Checking missing value in data set
df.isnull().sum()
```

Out[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
```

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

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

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

Out[8]:
```
Index(['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
       'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
       'pH', 'sulphates', 'alcohol', 'quality'],
      dtype='object')
```

In [9]:
```python
df.columns.tolist()
```

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

In [10]: `df.dtypes`

```
Out[10]:  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 [11]: `df.isnull().sum()`

```
Out[11]:  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 [12]: `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
```
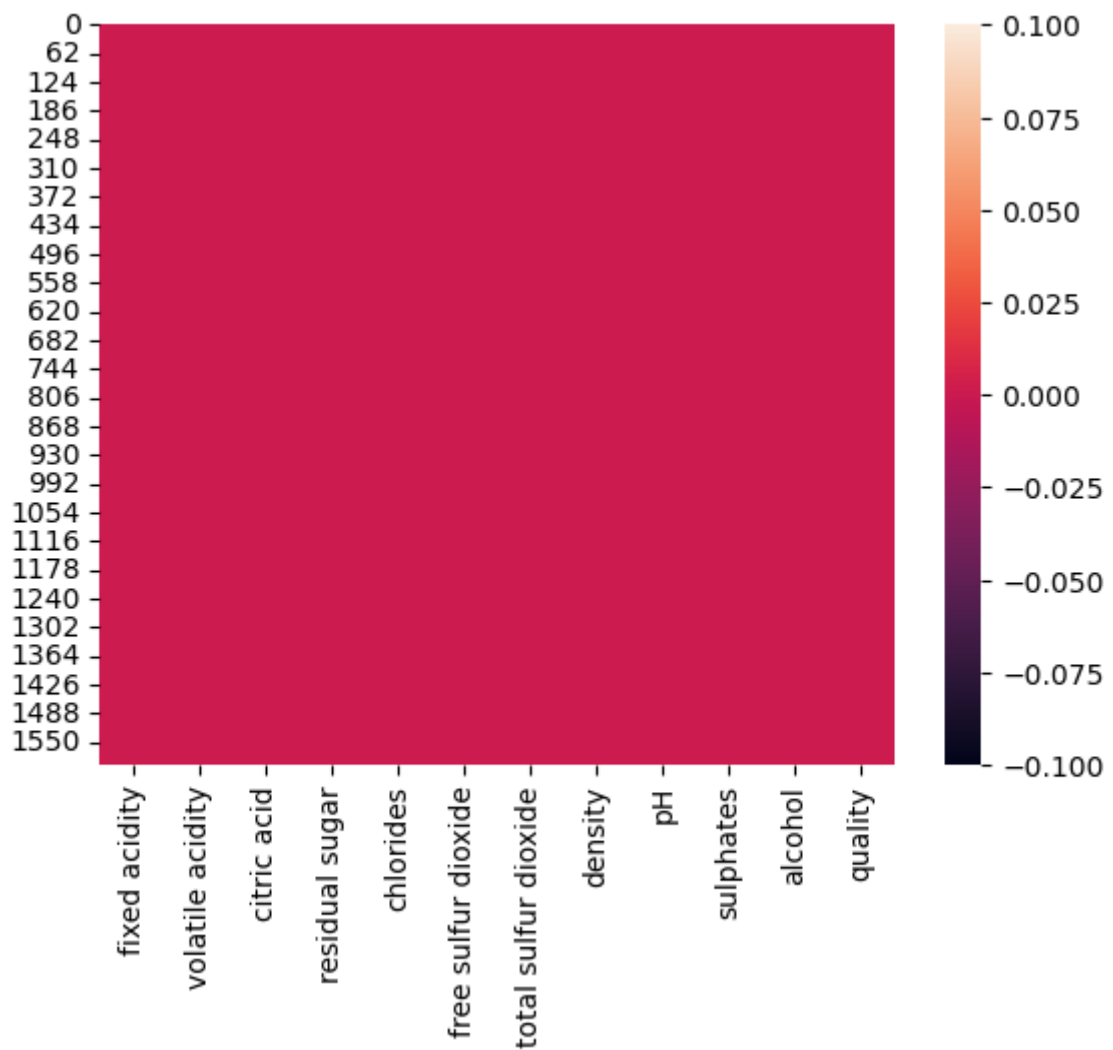
In [13]: `# Visualisation by using heatmap`
`sns.heatmap(df.isnull())`

Out[13]: `<Axes: >`



In [14]: `df['fixed acidity'].unique()`

```
Out[14]: array([ 7.4,  7.8, 11.2,  7.9,  7.3,  7.5,  6.7,  5.6,  8.9,  8.5,  8.1,
                7.6,  6.9,  6.3,  7.1,  8.3,  5.2,  5.7,  8.8,  6.8,  4.6,  7.7,
                8.7,  6.4,  6.6,  8.6, 10.2,  7. ,  7.2,  9.3,  8. ,  9.7,  6.2,
                5. ,  4.7,  8.4, 10.1,  9.4,  9. ,  8.2,  6.1,  5.8,  9.2, 11.5,
                5.4,  9.6, 12.8, 11. , 11.6, 12. , 15. , 10.8, 11.1, 10. , 12.5,
               11.8, 10.9, 10.3, 11.4,  9.9, 10.4, 13.3, 10.6,  9.8, 13.4, 10.7,
               11.9, 12.4, 12.2, 13.8,  9.1, 13.5, 10.5, 12.6, 14. , 13.7,  9.5,
               12.7, 12.3, 15.6,  5.3, 11.3, 13. ,  6.5, 12.9, 14.3, 15.5, 11.7,
               13.2, 15.9, 12.1,  5.1,  4.9,  5.9,  6. ,  5.5])
```

In [15]: `df['fixed acidity'].nunique()`

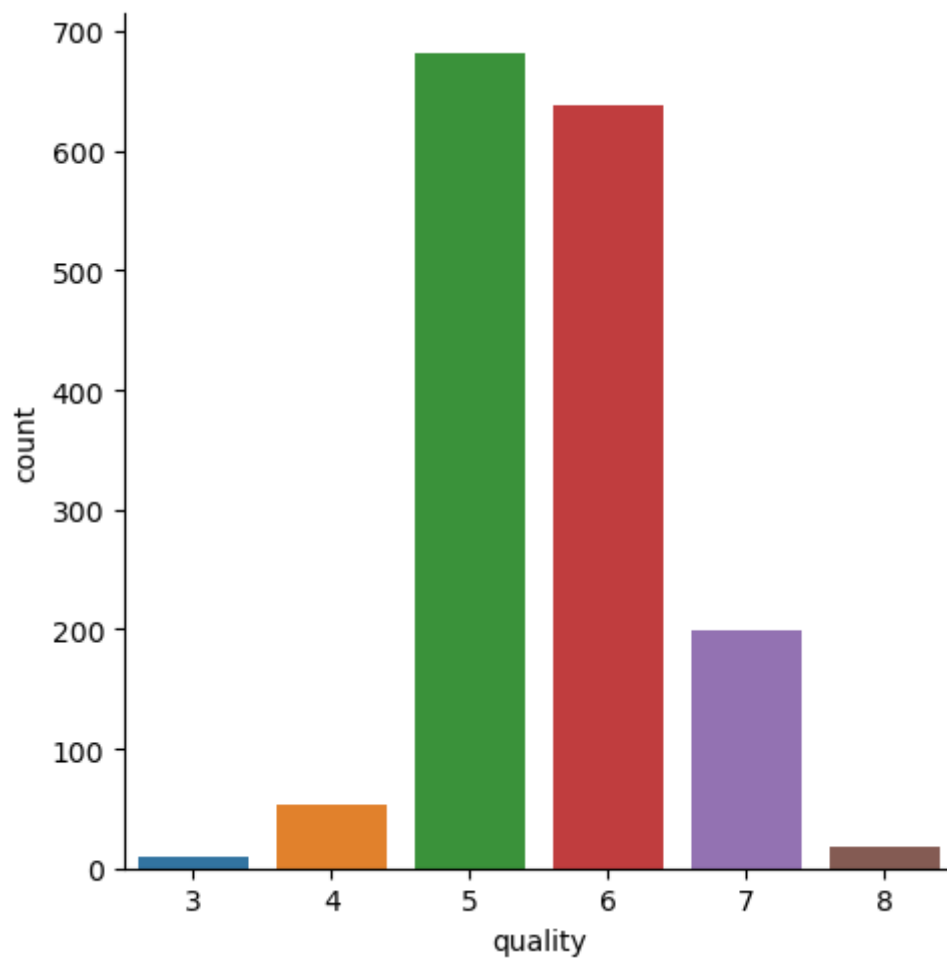Out[15]: 96

In [16]:
```python
# data analysis and visualization
df.describe()
```

Out[16]:

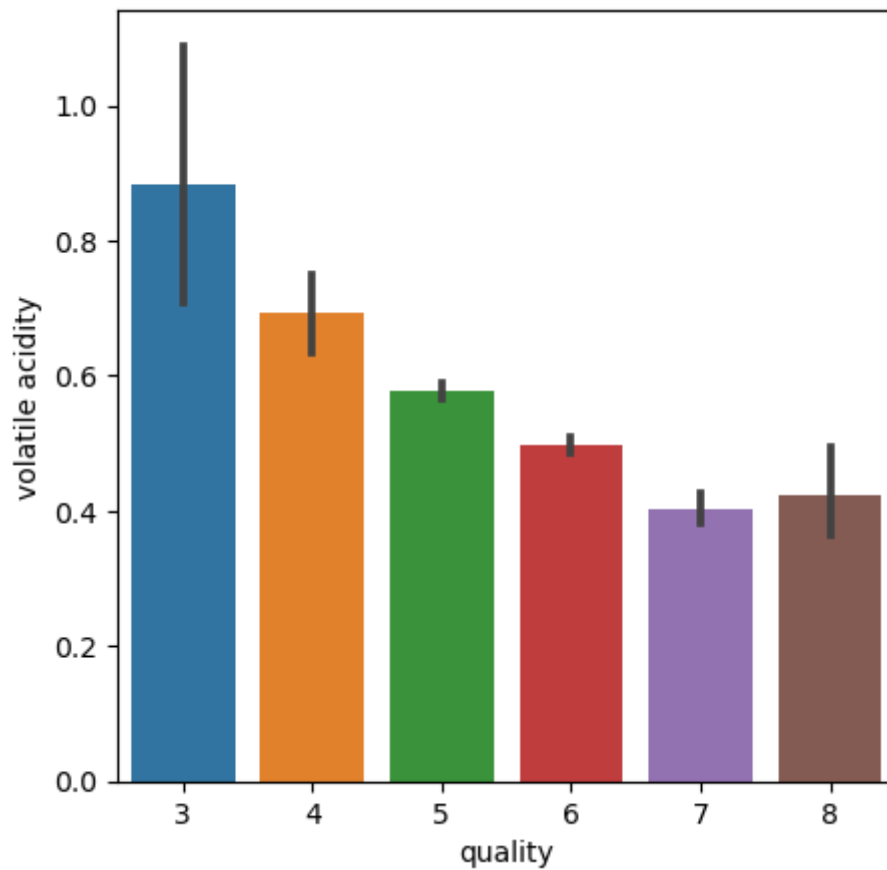| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide |
|---|---|---|---|---|---|---|---|
| count | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 | 1599.000000 |
| mean | 8.319637 | 0.527821 | 0.270976 | 2.538806 | 0.087467 | 15.874922 | 46.467792 |
| std | 1.741096 | 0.179060 | 0.194801 | 1.409928 | 0.047065 | 10.460157 | 32.895324 |
| min | 4.600000 | 0.120000 | 0.000000 | 0.900000 | 0.012000 | 1.000000 | 6.000000 |
| 25% | 7.100000 | 0.390000 | 0.090000 | 1.900000 | 0.070000 | 7.000000 | 22.000000 |
| 50% | 7.900000 | 0.520000 | 0.260000 | 2.200000 | 0.079000 | 14.000000 | 38.000000 |
| 75% | 9.200000 | 0.640000 | 0.420000 | 2.600000 | 0.090000 | 21.000000 | 62.000000 |
| max | 15.900000 | 1.580000 | 1.000000 | 15.500000 | 0.611000 | 72.000000 | 289.000000 |

In [17]:
```python
# Number of values of each quality
sns.catplot(x='quality', data = df, kind = "count")
```

Out[17]: `<seaborn.axisgrid.FacetGrid at 0x1e5f2df6dd0>`

In [18]:
```python
# volatile Acidity v/s Quality
plot = plt.figure(figsize=(5,5))
sns.barplot(x= 'quality', y = 'volatile acidity', data = df)
print('volatile acidity is inversely proptional to quality')
```
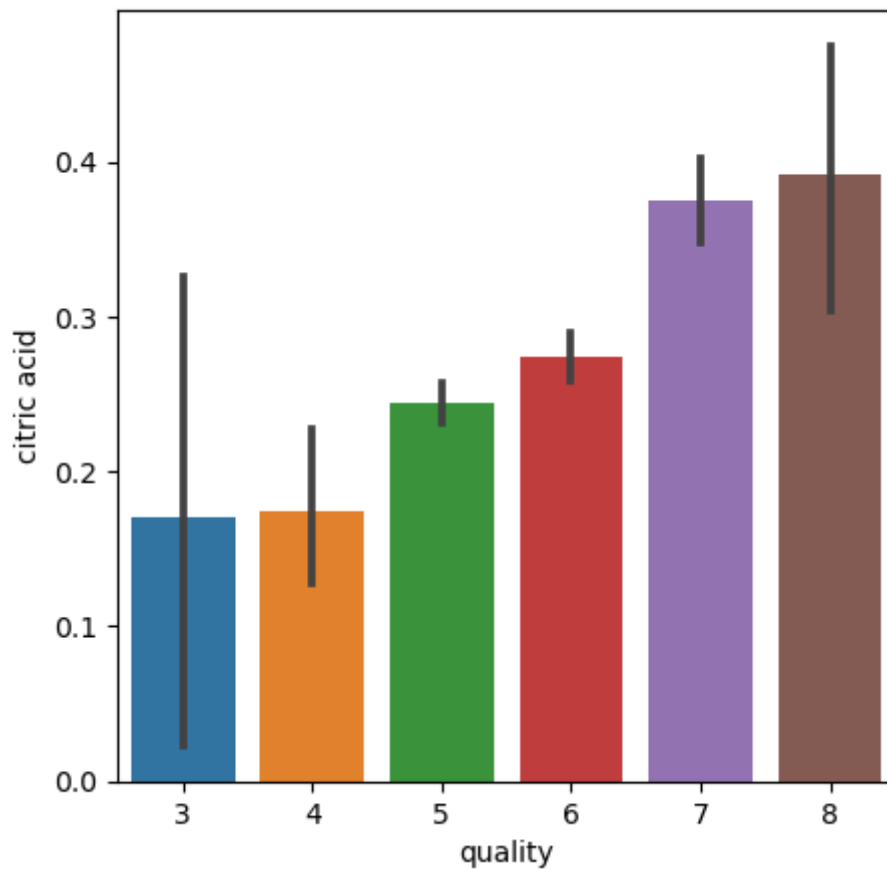
volatile acidity is inversely proptional to quality

In [19]:
```python
# citric acid v/s Quality
plot = plt.figure(figsize=(5,5))
sns.barplot(x= 'quality', y = 'citric acid', data = df)
print('volatile acidity is directly proptional to quality')
```

volatile acidity is directly proptional to quality

```
In [20]:  # residual sugar v/s Quality
          plot = plt.figure(figsize=(5,5))
          sns.barplot(x= 'quality', y = 'residual sugar', data = df)
```
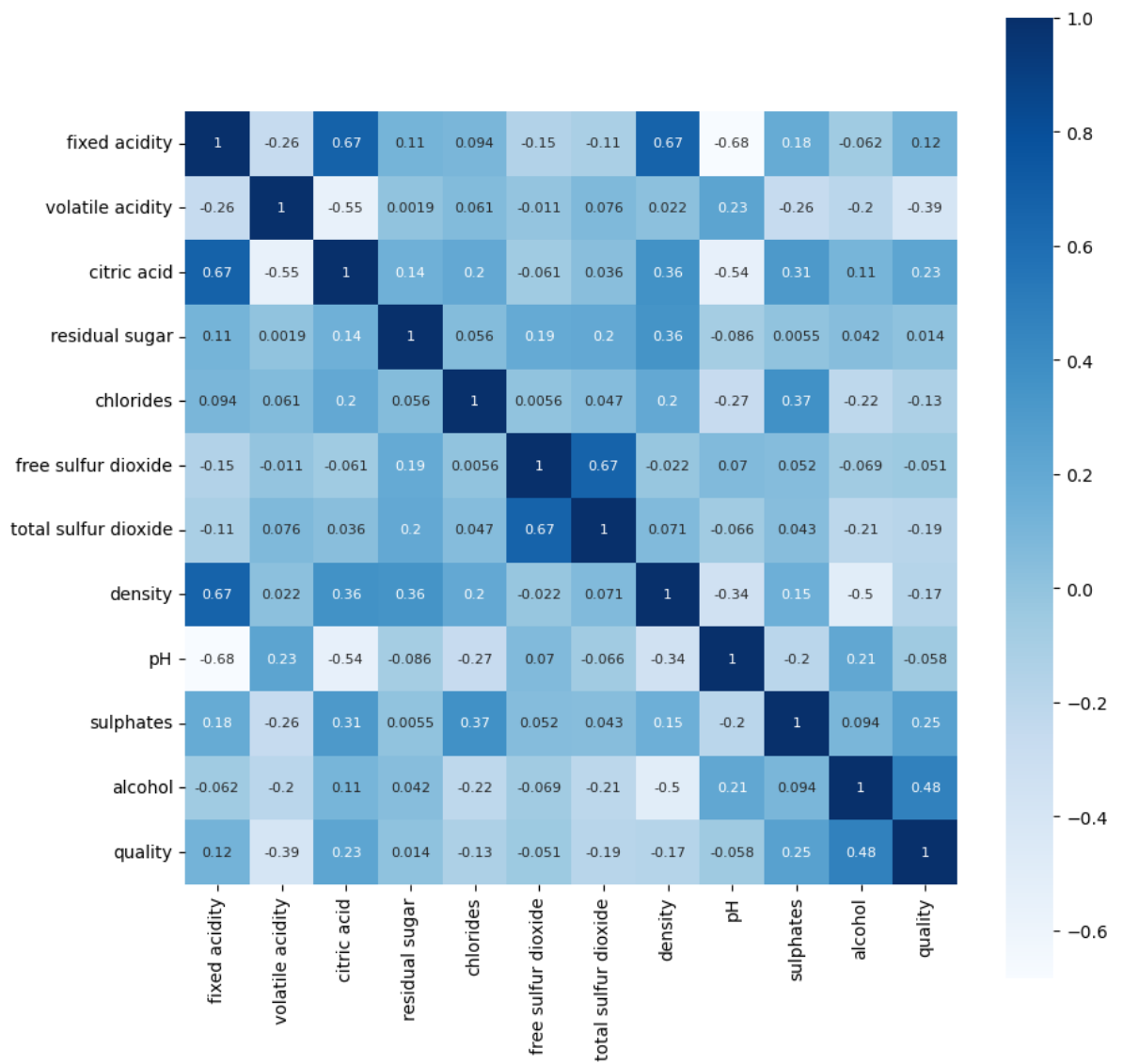
Out[20]:  `<Axes: xlabel='quality', ylabel='residual sugar'>`



```
In [21]:  correlation = df.corr()
          correlation
```

|  | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density |  |
|---|---|---|---|---|---|---|---|---|---|
| fixed acidity | 1.000000 | -0.256131 | 0.671703 | 0.114777 | 0.093705 | -0.153794 | -0.113181 | 0.668047 | -( |
| volatile acidity | -0.256131 | 1.000000 | -0.552496 | 0.001918 | 0.061298 | -0.010504 | 0.076470 | 0.022026 | ( |
| citric acid | 0.671703 | -0.552496 | 1.000000 | 0.143577 | 0.203823 | -0.060978 | 0.035533 | 0.364947 | -( |
| residual sugar | 0.114777 | 0.001918 | 0.143577 | 1.000000 | 0.055610 | 0.187049 | 0.203028 | 0.355283 | -( |
| chlorides | 0.093705 | 0.061298 | 0.203823 | 0.055610 | 1.000000 | 0.005562 | 0.047400 | 0.200632 | -( |
| free sulfur dioxide | -0.153794 | -0.010504 | -0.060978 | 0.187049 | 0.005562 | 1.000000 | 0.667666 | -0.021946 | ( |
| total sulfur dioxide | -0.113181 | 0.076470 | 0.035533 | 0.203028 | 0.047400 | 0.667666 | 1.000000 | 0.071269 | -( |
| density | 0.668047 | 0.022026 | 0.364947 | 0.355283 | 0.200632 | -0.021946 | 0.071269 | 1.000000 | -( |
| pH | -0.682978 | 0.234937 | -0.541904 | -0.085652 | -0.265026 | 0.070377 | -0.066495 | -0.341699 | |
| sulphates | 0.183006 | -0.260987 | 0.312770 | 0.005527 | 0.371260 | 0.051658 | 0.042947 | 0.148506 | -( |
| alcohol | -0.061668 | -0.202288 | 0.109903 | 0.042075 | -0.221141 | -0.069408 | -0.205654 | -0.496180 | ( |
| quality | 0.124052 | -0.390558 | 0.226373 | 0.013732 | -0.128907 | -0.050656 | -0.185100 | -0.174919 | -( |

In [22]:
```python
plt.figure(figsize=(10,10))
sns.heatmap(correlation, cbar = True, square = True, annot = True, annot_kws={'siz
```

Out[22]: <Axes: >

In [23]: 
```python
#data prepocessing
x = df.drop('quality',axis=1)
```

In [24]: 
```python
print(x)
```

|      | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | \ |
|------|---------------|------------------|-------------|----------------|-----------|---|
| 0    | 7.4           | 0.700            | 0.00        | 1.9            | 0.076     |   |
| 1    | 7.8           | 0.880            | 0.00        | 2.6            | 0.098     |   |
| 2    | 7.8           | 0.760            | 0.04        | 2.3            | 0.092     |   |
| 3    | 11.2          | 0.280            | 0.56        | 1.9            | 0.075     |   |
| 4    | 7.4           | 0.700            | 0.00        | 1.9            | 0.076     |   |
| ...  | ...           | ...              | ...         | ...            | ...       |   |
| 1594 | 6.2           | 0.600            | 0.08        | 2.0            | 0.090     |   |
| 1595 | 5.9           | 0.550            | 0.10        | 2.2            | 0.062     |   |
| 1596 | 6.3           | 0.510            | 0.13        | 2.3            | 0.076     |   |
| 1597 | 5.9           | 0.645            | 0.12        | 2.0            | 0.075     |   |
| 1598 | 6.0           | 0.310            | 0.47        | 3.6            | 0.067     |   |

|      | free sulfur dioxide | total sulfur dioxide | density | pH   | sulphates | \ |
|------|---------------------|----------------------|---------|------|-----------|---|
| 0    | 11.0                | 34.0                 | 0.99780 | 3.51 | 0.56      |   |
| 1    | 25.0                | 67.0                 | 0.99680 | 3.20 | 0.68      |   |
| 2    | 15.0                | 54.0                 | 0.99700 | 3.26 | 0.65      |   |
| 3    | 17.0                | 60.0                 | 0.99800 | 3.16 | 0.58      |   |
| 4    | 11.0                | 34.0                 | 0.99780 | 3.51 | 0.56      |   |
| ...  | ...                 | ...                  | ...     | ...  | ...       |   |
| 1594 | 32.0                | 44.0                 | 0.99490 | 3.45 | 0.58      |   |
| 1595 | 39.0                | 51.0                 | 0.99512 | 3.52 | 0.76      |   |
| 1596 | 29.0                | 40.0                 | 0.99574 | 3.42 | 0.75      |   |
| 1597 | 32.0                | 44.0                 | 0.99547 | 3.57 | 0.71      |   |
| 1598 | 18.0                | 42.0                 | 0.99549 | 3.39 | 0.66      |   |

|      | alcohol |
|------|---------|
| 0    | 9.4     |
| 1    | 9.8     |
| 2    | 9.8     |
| 3    | 9.8     |
| 4    | 9.4     |
| ...  | ...     |
| 1594 | 10.5    |
| 1595 | 11.2    |
| 1596 | 11.0    |
| 1597 | 10.2    |
| 1598 | 11.0    |

[1599 rows x 11 columns]

In [26]:
```python
#label binarization
Y = df['quality'].apply(lambda y_value: 1 if y_value>= 7 else 0)

print(Y)
```

```
0       0
1       0
2       0
3       0
4       0
       ..
1594    0
1595    0
1596    0
1597    0
1598    0
Name: quality, Length: 1599, dtype: int64
```

In [27]:
```python
#train & Test Split

X_train, X_test, Y_train, Y_test = train_test_split(x, Y, test_size=0.2, random_sta
```

```python
In [28]: print(Y.shape, Y_train.shape, Y_test.shape)

(1599,) (1279,) (320,)

In [29]: #Model Training: Random Forest Classifier
         model = RandomForestClassifier()

In [30]: model.fit(X_train, Y_train)

Out[30]:  ▾ RandomForestClassifier

         RandomForestClassifier()


In [31]: #Accuracy Score
         # accuracy on test data
         X_test_prediction = model.predict(X_test)
         test_data_accuracy = accuracy_score(X_test_prediction, Y_test)

In [32]: print('Accuracy : ', test_data_accuracy)

Accuracy :  0.925

In [33]: #Building a Predictive System

In [34]: input_data = ('7.3,0.65,0.0,1.2,0.065,15.0,21.0,0.9946,3.39,0.47,10.0,')

In [35]: input_data = (7.5,0.5,0.36,6.1,0.071,17.0,102.0,0.9978,3.35,0.8,10.5)

In [36]: input_data_as_numpy_array = np.asarray(input_data)

In [37]: input_data_reshaped = input_data_as_numpy_array.reshape(1,-1)

In [38]: prediction = model.predict(input_data_reshaped)
         print(prediction)

         if (prediction[0]==1):
           print('Good Quality Wine')
         else:
           print('Bad Quality Wine')

[0]
Bad Quality Wine

In [ ]:
```