

RED WINE QUALITY ANALYSIS

```
In [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

In [2]: # Provide the full path to the CSV file
file_path = r"D:\Cognorise\Infotech\winequality-red.csv"

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

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

	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	0.9988	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]: 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: 158.0 kb

In [5]: df.isnull().sum()

Out[5]:
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 [6]: df.shape

Out[6]: (1599, 12)

In [7]: df.describe()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
count	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000
mean	8.318637	0.527821	0.270978	2.538908	0.087467	13.674922	46.467782	0.986747	3.311133	0.58149	10.422983	5.836023
std	1.741096	0.179900	0.158081	1.409028	0.047905	10.400157	32.996324	0.001087	0.154306	0.199957	1.956568	0.807959
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000	0.989070	2.740000	0.330000	8.400000	3.000000
25%	7.100000	0.300000	0.000000	1.900000	0.070000	7.000000	22.000000	0.989600	3.210000	0.500000	9.500000	5.000000
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000	0.986750	3.310000	0.620000	10.200000	6.000000
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000	0.997835	3.400000	0.730000	11.100000	6.000000
max	15.900000	1.580000	1.000000	15.500000	0.611000	72.000000	289.000000	1.003690	4.010000	2.000000	14.900000	8.000000

```
In [8]: df.corr()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
fixed acidity	1.000000	-0.256131	0.671703	0.114777	0.093705	-0.153794	-0.113181	0.668047	-0.682978	0.183006	-0.061668	0.124052
volatile acidity	-0.256131	1.000000	-0.552496	0.001918	0.061298	-0.010504	0.076470	0.022026	-0.234937	-0.260987	-0.202288	-0.300558
citric acid	0.671703	-0.552496	1.000000	0.143577	0.203823	-0.060978	0.035533	0.364947	-0.541904	0.312770	0.199903	0.226373
residual sugar	0.114777	0.001918	0.143577	1.000000	0.056610	0.187049	0.203028	0.352883	-0.089562	0.005527	0.042075	0.013732
chlorides	0.093705	0.001298	0.203823	0.056610	1.000000	0.005662	0.047400	0.200632	-0.286026	0.371260	-0.221141	-0.128907
free sulfur dioxide	-0.153794	-0.010504	-0.060978	0.187049	0.005662	1.000000	0.667666	-0.021946	0.071269	0.070377	0.051688	-0.068408
total sulfur dioxide	-0.113181	0.076470	0.035533	0.203028	0.047400	0.667666	1.000000	0.071269	-0.066495	0.042947	-0.205654	-0.185100
density	0.668047	0.022026	0.364947	0.352883	0.200632	-0.021946	0.071269	1.000000	-0.341699	0.148506	0.496180	-0.174919
pH	-0.682978	-0.234937	-0.541904	-0.089562	-0.286026	0.070377	-0.066495	-0.341699	1.000000	-0.196648	0.205633	-0.057731
sulphates	0.183006	-0.260987	0.312770	0.005527	0.371260	0.051688	0.042947	0.148506	-0.196648	1.000000	0.093995	0.251397
alcohol	-0.061668	-0.202288	0.199903	0.042075	-0.221141	-0.069408	-0.205654	-0.496180	0.205633	0.093995	1.000000	0.476166
quality	0.124052	-0.300558	0.226373	0.013732	-0.128907	-0.068408	-0.185100	-0.174919	-0.057731	0.251397	0.476166	1.000000

```
In [9]: corr_data = df.corr()

In [10]: plt.figure(figsize=(9,9))
sns.heatmap(corr_data, char = True, square = True, annot=True, fwt = '2f', cmap='Blues')
plt.title('Correlation Heatmap')

Out[10]: Text(0.5, 1.0, 'Correlation Heatmap')
```

```
In [11]: def hist_plots(df):
    plt.figure(figsize=(10,8))
    plt.hist(df)
    plt.title('Histogram Plot')
    plt.show()
    hist_plots(df['fixed acidity'])
```

```
In [12]: hist_plots(df['pH'])
```

```
In [13]: # Number of wines in each quality category
sns.catplot(x='quality', y='count', kind='count')
plt.title('Number of wines in each quality category')

Out[13]: Text(0.5, 1.0, 'Number of wines in each quality category')
```

```
In [14]: hist_plots(df['quality'])
```

```
In [15]: # plotting a barplot for quality vs volatile acidity
sns.barplot(x='quality', y='volatile acidity', data = df)

Out[15]: <AxesSubplot:xlabel='quality', ylabel='volatile acidity'>
```

```
In [16]: hist_plots(df['alcohol'])
```

```
In [17]: # plotting a barplot for quality vs alcohol
sns.barplot(x='quality', y='alcohol', data = df)

Out[17]: <AxesSubplot:xlabel='quality', ylabel='alcohol'>
```

```
In [18]: hist_plots(df['residual sugar'])
```

```
In [19]: sns.kdeplot(df['volatile acidity'])

Out[19]: <AxesSubplot:xlabel='volatile acidity', ylabel='Density'>
```

```
In [20]: sns.kdeplot(df['chlorides'])

Out[20]: <AxesSubplot:xlabel='chlorides', ylabel='Density'>
```

```
In [21]: sns.kdeplot(df['density'])

Out[21]: <AxesSubplot:xlabel='density', ylabel='Density'>
```

```
In [22]: plt.figure(figsize=(12,6))
sns.pairplot(df)
plt.show()

<Figure size 1200x600 with 0 Axes>
```

```
In [23]: df.hist(figsize=(11,11))
plt.show()
```