

WINE QUALITY PREDICTION

```
In [2]: import warnings
warnings.simplify("ignore")

In [3]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib
import scipy.stats as stats
from scipy.stats import score

In [4]: from IPython.core.display import display, HTML
display(HTML("<table>container</table>"))
```

IMPORTING DATA

```
In [4]: WD = pd.read_csv('QualityPrediction.csv')
WD.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	0.9968	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

```
Out[6]: WD.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
Dtype: float64(11), int64(1)
memory usage: 150.0 KB
```

```
In [6]: WD.shape
```

```
Out[6]: (1599, 12)
```

```
Out[7]: WD.isnull().sum()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
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 [8]: WD.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	8.319637	0.527821	0.270976	2.538066	0.087467	15.874922	46.467792	0.996747	3.311113	0.681494	9.681494	5.0
std	1.741096	0.179060	0.194801	1.409928	0.047065	10.460157	32.895324	0.001887	0.154386	0.169507	0.169507	0.0
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000	0.996070	2.740000	0.330000	8.000000	3
25%	7.100000	0.390000	0.090000	1.900000	0.070000	7.000000	22.000000	0.995600	3.210000	0.550000	9.000000	4
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000	0.996750	3.310000	0.620000	9.500000	5
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000	0.997835	3.400000	0.700000	10.000000	6
max	15.900000	1.500000	1.000000	15.500000	0.610000	72.000000	289.000000	1.003690	4.010000	2.000000	16.000000	10

Observation:

(a) There is a big gap between 75% and max values of residual sugar, free sulfur dioxide and total sulfur dioxide.

(b) This is indicative of outliers.

Renaming Columns

```
In [9]: WD = WD.rename(columns={'fixed acidity': 'fixed_acidity', 'volatile acidity': 'volatile_acidity', 'citric acid': 'citric_acid', 'residual sugar': 'residual_sugar', 'free sulfur dioxide': 'free_sulfur_dioxide', 'total sulfur dioxide': 'total_sulfur_dioxide', 'density': 'density', 'pH': 'pH', 'sulphates': 'sulphates', 'alcohol': 'alcohol', 'quality': 'quality'})
WD.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	0.9968	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

```
Out[9]: WD.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
Dtype: float64(11), int64(1)
memory usage: 150.0 KB
```

```
In [10]: WD.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	8.319637	0.527821	0.270976	2.538066	0.087467	15.874922	46.467792	0.996747	3.311113	0.681494	9.681494	5.0
std	1.741096	0.179060	0.194801	1.409928	0.047065	10.460157	32.895324	0.001887	0.154386	0.169507	0.169507	0.0
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000	0.996070	2.740000	0.330000	8.000000	3
25%	7.100000	0.390000	0.090000	1.900000	0.070000	7.000000	22.000000	0.995600	3.210000	0.550000	9.000000	4
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000	0.996750	3.310000	0.620000	9.500000	5
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000	0.997835	3.400000	0.700000	10.000000	6
max	15.900000	1.500000	1.000000	15.500000	0.610000	72.000000	289.000000	1.003690	4.010000	2.000000	16.000000	10

Observation:

(a) There is a big gap between 75% and max values of residual sugar, free sulfur dioxide and total sulfur dioxide.

(b) This is indicative of outliers.

Visualisation

```
In [10]: plt.figure(figsize=(10,10))
sns.countplot(WD['quality'], palette='viridis')
WD['quality'].value_counts()
```

quality	count
3	1
4	1
5	681
6	638
7	199
8	53
9	18
10	3

```
Out[10]: Name: quality, dtype: int64
```

From the plot and value counts it is clearly visible that there is imbalance in the dataset. Needs to be rectified.

Plotting Feature variables against Target Label

```
In [11]: sns.pairplot(WD, kind='scatter', diag_kind='kde', diag_sharey=True, figsize=(25,10))
```

Observations regarding features compared to the target variable(quality)increasing are:

- fixed acidity - no pattern
- volatile acidity - negative trend
- citric - acid positive trend
- residual sugar - no pattern
- chlorides - negative trend
- free sulfur dioxide - no fixed pattern
- total sulfur dioxide - no fixed pattern
- density - more or less constant
- pH - more or less constant
- sulphates - positive trend
- alcohol - positive trend

Correlations

```
In [12]: # External Corr
correlations = WD.corr()['quality'].sort_values(ascending=False)
print(correlations)
```

	quality
quality	1.000000
alcohol	0.476166
sulphates	0.231397
citric_acid	0.226373
fixed_acidity	0.144052
residual_sugar	0.013732
free_sulfur_dioxide	-0.050656
pH	-0.057731
chlorides	-0.128807
density	-0.174919
total_sulfur_dioxide	-0.185100
volatile_acidity	-0.390558
Name: quality, dtype: float64	

```
In [13]: correlations.plot(kind='bar', color='r')
```

```
Out[13]: <AxesSubplot:~>
```

```
In [14]: # Internal Corr
plt.figure(figsize=(16,10))
sns.heatmap(WD.corr(), vmin=-1, vmax=1, annot=True, cmap='Spectral', annot_kws={'size': 12})
plt.style.use('dark')
plt.xticks(fontsize=14)
plt.yticks(fontsize=14)
plt.show()
```

- Fixed acidity and citric acid are positively correlated with a value of 0.672 which is close to 1. Similarly, fixed acidity and density are positively correlated with a value of 0.668. The other 2 column that's positively correlated are free sulfur dioxide and total sulfur dioxide with a value of 0.668. The only negatively correlated columns are fixed acidity and pH with a value -0.683 being close to the value -1. Quality column is least correlated with residual sugar showing a coefficient value of 0.014.
- Need to check for multicollinearity.

PRE PROCESSING

Dropping Columns

Free_sulfur_dioxide can be dropped - Also fixed acidity can be dropped to obviate multicollinearity as it has relatively strong corr with citric_acid,density and pH.

```
In [15]: WD = WD.drop(['free_sulfur_dioxide', 'fixed_acidity'], axis=1)
WD
```

	volatile acidity	citric acid	residual_sugar	chlorides	total_sulfur_dioxide	density	pH	sulphates	alcohol	quality
0	0.700	0.00	1.9	0.076	34.0	0.99780	3.51	0.56	9.4	5
1	0.880	0.00	2.6	0.098	67.0	0.99680	3.20	0.68	9.8	5
2	0.760	0.04	2.3	0.092	54.0	0.99700	3.26	0.65	9.8	5
3	0.280	0.56	1.9	0.075	60.0	0.99800	3.16	0.58	9.8	6
4	0.700	0.00	1.9	0.076	34.0	0.99780	3.51	0.56	9.4	5
...
1594	0.600	0.08	2.0	0.090	44.0	0.99490	3.45	0.58	10.5	5
1595	0.550	0.10	2.2	0.062	51.0	0.99512	3.52	0.76	11.2	6
1596	0.510	0.13	2.3	0.076	40.0	0.99574	3.42	0.75	11.0	6
1597	0.645	0.12	2.0	0.075	44.0	0.99547	3.57	0.71	10.2	5
1598	0.310	0.47	3.6	0.067	42.0	0.99549	3.39	0.66	11.0	6

1599 rows x 10 columns

Outlier Removal using z score

- We can use the fact that if z score of a datapoint is greater than 3, it can be treated as an outlier

```
In [16]: zmap, abs(zscore)
threshold=3
np.where(z>3)

WD[WD[(z<3).all(axis=1)]]
WD
```

	volatile acidity	citric acid	residual_sugar	chlorides	total_sulfur_dioxide	density	pH	sulphates	alcohol	quality
0	0.700	0.00	1.9	0.076	34.0	0.99780	3.51	0.56	9.4	5
1	0.880	0.00	2.6	0.098	67.0	0.99680	3.20	0.68	9.8	5
2	0.760	0.04	2.3	0.092	54.0	0.99700	3.26	0.65	9.8	5
3	0.280	0.56	1.9	0.075	60.0	0.99800	3.16	0.58	9.8	6
4	0.700	0.00	1.9	0.076	34.0	0.99780	3.51	0.56	9.4	5
...
1594	0.600	0.08	2.0	0.090	44.0	0.99490	3.45	0.58	10.5	5
1595	0.550	0.10	2.2	0.062	51.0	0.99512	3.52	0.76	11.2	6
1596	0.510	0.13	2.3	0.076	40.0	0.99574	3.42	0.75	11.0	6
1597	0.645	0.12	2.0	0.075	44.0	0.99547	3.57	0.71	10.2	5
1598	0.310	0.47	3.6	0.067	42.0	0.99549	3.39	0.66	11.0	6

1471 rows x 10 columns

From 1599 rowsof dataset, the same now reduced to 1471 rows.

Rectification of Class Imbalance in Dataset using SMOTE(Synthetic Minority Oversampling Technique)

```
In [17]: from imblearn.over_sampling import SMOTE

In [18]: X = WD.drop('quality',axis=1)
Y = WD['quality']

In [19]: Y.value_counts()
```

quality	count
3	1
4	1
5	624
6	624
7	199
8	47
9	16

```
Out[19]: Name: quality, dtype: int64
```

```
In [20]: os = SMOTE()
Xs, Ys = os.fit_resample(X, Y)
```

```
Out[20]: Y.value_counts()
```

quality	count
3	1
4	1
5	624
6	624
7	199
8	47
9	16

```
Out[21]: os = SMOTE()
Xs, Ys = os.fit_resample(X, Y)
```

quality	count
3	1
4	1
5	624
6	624
7	199
8	47
9	16

```
Out[22]: Y
```

quality	count
3	1
4	1
5	624
6	624
7	199
8	47
9	16

```
Out[23]: Name: quality, Length: 3120, dtype: int64
```

Converting Output Label to Binary

```
In [23]: Y = Y.apply(lambda y_value: 1 if y_value>7 else 0)
Y
```

quality	count
0	0
1	0
2	0
3	0
4	0
5	1
6	1
7	1
8	1
9	1

```
Out[23]: Name: quality, Length: 3120, dtype: int64
```

```
In [24]: X
```

	volatile acidity	citric acid	residual_sugar	chlorides	total_sulfur_dioxide	density	pH	sulphates	alcohol
0	0.700000	0.000000	1.900000	0.076000	34.0				

Text(0.5, 1.0, 'OOB Error Rate Across Various Forest sizes \n(From 100 to 1000 trees)')



```
In [56]: print('OOB Error rate for 160 trees is: {}'.format(oob_series_RFC[160]))
```

OOB Error rate for 160 trees is: 0.05048

```
In [67]: # Refine the tree via OOB Output
RFCF.set_params(n_estimators=160,
                bootstrap = True,
                oob_score=False, warm_start=False)
```

```
Out[67]: RandomForestClassifier(criterion='entropy', max_depth=20, max_features='sqrt',
                               min_sample_split=5, n_estimators=160, random_state=42)
```

```
In [68]: RFCF.fit(X_train, Y_train)
          Y_pred2 = RFCF.predict(X_test)
          RFCF_acc = (accuracy_score(Y_test, Y_pred2))*100
          print("Accuracy score for the Best Model is", RFCF_acc)
          print(classification_report(Y_test, Y_pred2))
```

Accuracy score for the Best Model is: 93.58974358974359

	precision	recall	f1-score	support
0	0.96	0.93	0.95	370
1	0.91	0.94	0.92	254
accuracy			0.94	624
macro avg	0.93	0.94	0.93	624
weighted avg	0.94	0.94	0.94	624

ROC Curve

```
In [72]: disp = metrics.plot_roc_curve(RFCF, X_test, Y_test)
          disp.figure_.suptitle("ROC Curve")
          plt.show()
```



In []: