

# 葡萄酒质量预测

## 1. 导入数据

导入所需的Python库并加载白葡萄酒和红葡萄酒的两个数据集。

```
1 from IPython.core.interactiveshell import InteractiveShell
2 InteractiveShell.ast_node_interactivity = "all"
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import pandas as pd
6 import seaborn as sns
7 from matplotlib import font_manager
8
9 wine_red = pd.read_csv('winequality-red.csv', sep=';')
10 wine_white = pd.read_csv('winequality-white.csv', sep=';')
11 wine = pd.concat([wine_red, wine_white], axis=0)
12 wine.head()
```

```
In [1]: '''导入所需的Python库并加载白葡萄酒和红葡萄酒的两个数据集。'''
from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = "all"
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns

wine_red = pd.read_csv('winequality-red.csv', sep=';')
wine_white = pd.read_csv('winequality-white.csv', sep=';')
wine = pd.concat([wine_red, wine_white], axis=0)
wine.head()
```

```
Out[1]:
```

	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

葡萄酒成分列表：

列名	含义
fixed acidity	挥发性酸
volatile acidity	挥发性酸
citric acid	柠檬酸
residual sugar	剩余糖分
chlorides	氯化物
free sulfur dioxide	游离二氧化硫
total sulfur dioxide	总二氧化硫
density	密度
pH	酸碱度
sulphates	硫酸盐
alcohol	酒精
quality	质量

## 2. 检查是数据否有空值

检查是否有空值

```
1 wine.isnull().sum()
```

```
In [2]: '''检查是否有空值'''  
        wine.isnull().sum()
```

```
Out[2]: '检查是否有空值'
```

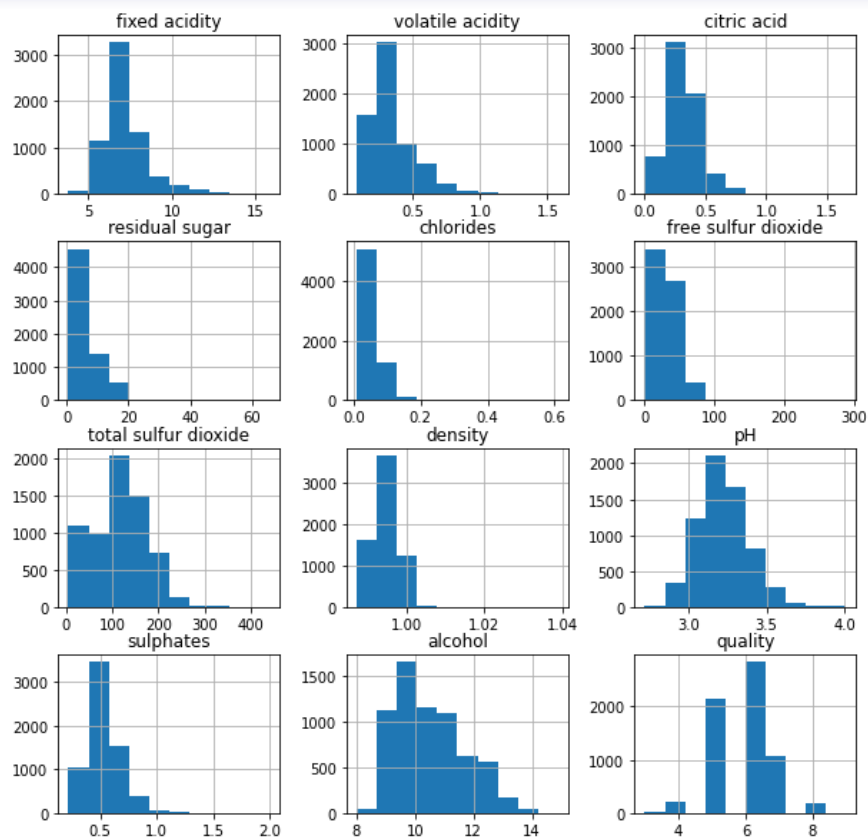
```
Out[2]: fixed acidity      0  
        volatile acidity  0  
        citric acid       0  
        residual sugar    0  
        chlorides         0  
        free sulfur dioxide 0  
        total sulfur dioxide 0  
        density          0  
        pH               0  
        sulphates        0  
        alcohol           0  
        quality           0  
        dtype: int64
```

检查发现数据集无空值，则不必进行空值处理

## 3. 绘制直方图查看数据的分布

查看各个特征值的分布数量

```
1 wine.hist(figsize=(10, 10))
```

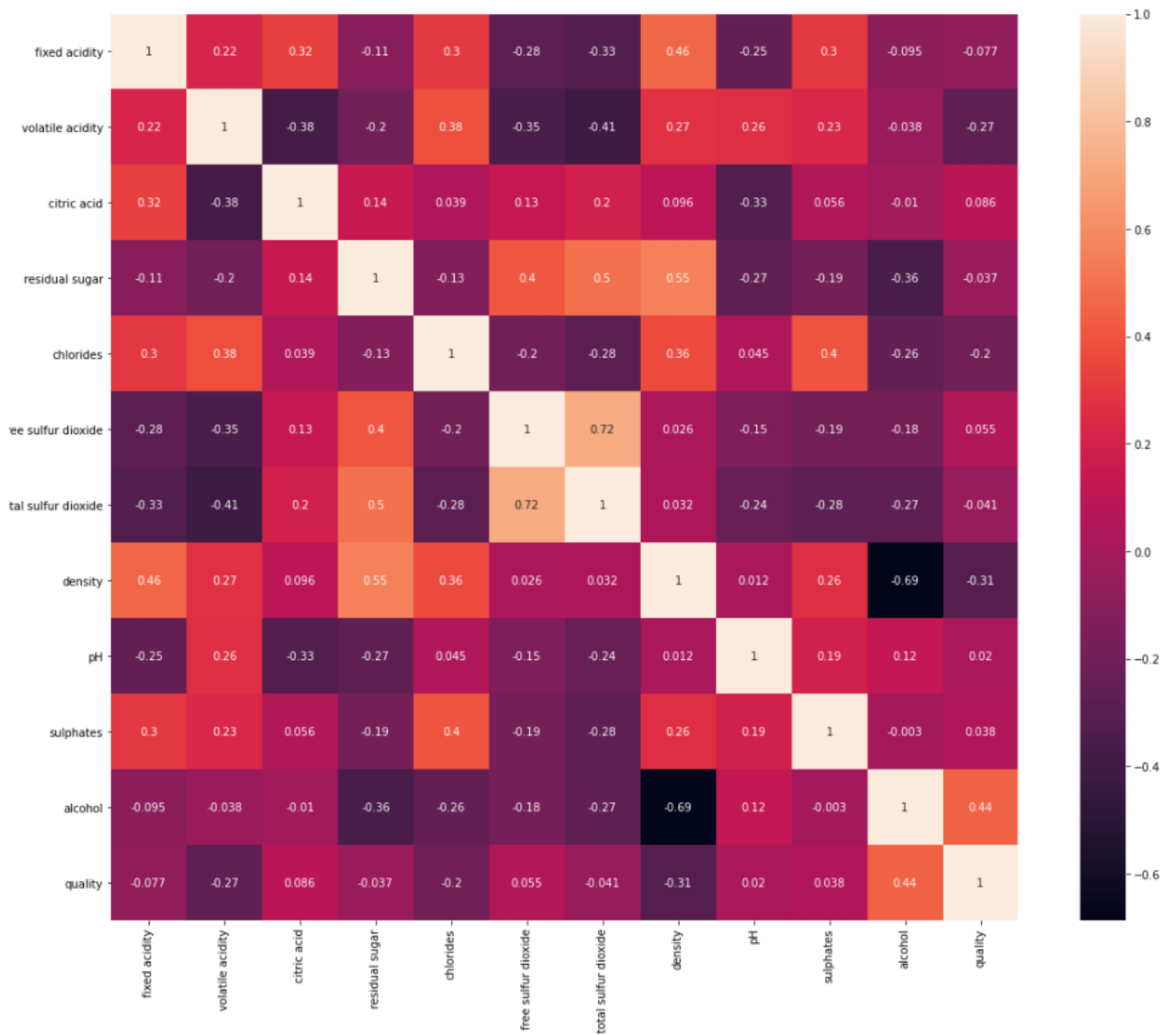


横坐标为特征值，纵坐标为数量

## 4. 绘制各个特征量与葡萄酒质量的相关性

找出输出(质量)变量与所有输入变量之间的相关性:

```
1 plt.subplots(figsize=(20,15)) #定义热力图大小为 (20, 15)  
2 corr = wine.corr() #相关系数  
3 sns.heatmap(corr,square=True, annot=True) #以corr为数据创建热力图
```



图中发现citric acid(柠檬酸)、free sulfur dioxide(游离二氧化硫)、pH(酸碱度)、sulphates(硫酸盐)、alcohol(酒精)与quality(质量)为正相关，于是我们选择这六个特征量绘制相关性图

```
1 cols = corr.nlargest(6, 'quality')['quality'].index #找到与目标值相关性最大的6个特征，而这几个特征之间的相关性要低。
2 corrcoeff = np.corrcoef(wine[cols].values.T)
3 plt.subplots(figsize=(20,15)) #设置热力图大小
4 sns.heatmap(corrcoeff,square=True, annot=True, xticklabels= cols.values, yticklabels=cols.values) #创建热力图
```



## 5. 用各种算法计算拟合训练数据，并根据测试值确定预测输出的准确性

### 准备工作

使用机器学习中的sklearn库，将数据集拆分为测试和训练数据集，我使用了20%的数据作为测试数据集：

```
1 y = wine["quality"]    #获取quality列
2 x = wine.drop(columns=["quality", "fixed acidity", "volatile acidity", "residual sugar", "chlorides",
3                       "total sulfur dioxide", "density"], axis=1)    #删去quality列和负相关的列
4 from sklearn.model_selection import train_test_split
5 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)
```

由于不同的列具有不同的值，因此使用StandardScaler库归一化值以获得准确的预测结果

```
1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 x_train = scaler.fit_transform(x_train)
4 x_test = scaler.fit_transform(x_test)
```

### 一、线性回归

```

1 from sklearn.metrics import accuracy_score, confusion_matrix
2 from sklearn.linear_model import LogisticRegression
3 logreg = LogisticRegression()
4 logreg.fit(x_train, y_train)
5 pred_logreg = logreg.predict(x_test)    #逻辑回归预测葡萄酒的质量
6 accuracy = accuracy_score(pred_logreg, y_test)    #计算分类准确率分数
7 print("Logreg Accuracy Score %.2f" % accuracy)

```

结果图：

Out[8]: '现在，我将根据各种算法拟合我的训练数据，并根据测试值确定预测输出的准确性。Python实现如下：'

Out[8]: '线性算法：'

Out[8]: LogisticRegression()

Logreg Accuracy Score 0.49

## 二、KNN

```

1 from sklearn.neighbors import KNeighborsClassifier
2 cm = confusion_matrix(pred_logreg, y_test)    #混淆矩阵
3 knn = KNeighborsClassifier(n_neighbors=1)    #K近邻算法，选取最近的点的个数为1
4 knn.fit(x_train, y_train)
5 pred_knn = knn.predict(x_test)    #knn算法预测葡萄酒的质量
6 accuracy = accuracy_score(pred_knn, y_test)    #计算分类准确率分数
7 print("Knn Accuracy Score %.2f" % accuracy)

```

结果图：

Out[9]: 'KNN'

Out[9]: KNeighborsClassifier(n\_neighbors=1)

Knn Accuracy Score 0.59

## 三、支持向量机SVC

```

1 from sklearn.svm import SVC
2 svc = SVC()
3 svc.fit(x_train, y_train)    #训练数据集
4 pred_svc =svc.predict(x_test)    #SVC预测葡萄酒的质量
5 accuracy = accuracy_score(pred_svc, y_test)    #计算分类准确率分数
6 print("SVC Accuracy Score %.2f" % accuracy)

```

结果图：

Out[10]: '支持向量机SVC'

Out[10]: SVC()

SVC Accuracy Score 0.52

## 四、随机森林

```

1 from sklearn.ensemble import RandomForestClassifier
2 rf = RandomForestClassifier()
3 rf.fit(x_train, y_train)    #训练数据集
4 pred_rf =rf.predict(x_test)    #随机森林预测葡萄酒的质量
5 accuracy = accuracy_score(pred_rf, y_test)    #计算分类准确率分数
6 print("Random Forest Accuracy Score %.2f" % accuracy)

```

结果图：

```
Out[11]: '随机森林'
```

```
Out[11]: RandomForestClassifier()
```

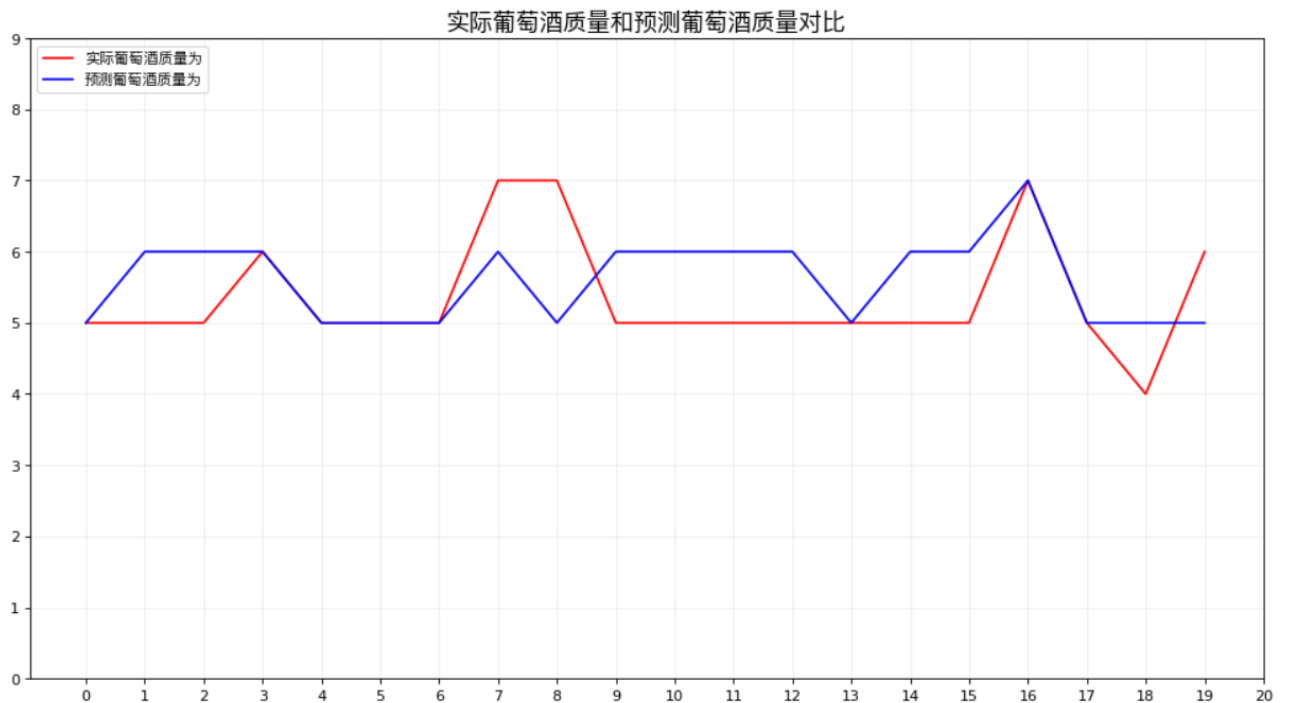
```
Random Forest Accuracy Score 0.66
```

在以上四种算法中，我们发现随机森林算法的预测准确率最高，所以接下来我们选用随机森林算法来比较实际葡萄酒质量和预测葡萄酒质量。

我们从x\_test测试数据集中徐州呢前20行创建新的观测集

```
1 new_observations = x.head(20) # 用数据集中的前20行创建新观测集
2 new_observations = scaler.fit_transform(new_observations) #对new_observations数据进行归一化
3 y_predicted = rf.predict(new_observations) # 用随机森林法对新观测集进行质量预测
4
5 #设置字体
6 my_font = font_manager.FontProperties(fname=r"C:/Users/xjqhre/PingFang SC.ttf")
7 #设置图片大小
8 plt.figure(figsize=(15, 8), dpi=80)
9 #画图—折线图
10 plt.plot(range(len(y.head(20).values)), y.head(20).values, label='实际葡萄酒质量为',color="r")
11 plt.plot(range(len(y_predicted)), y_predicted, label='预测葡萄酒质量为',color="b")
12 #设置x,y坐标
13 plt.xticks(range(len(y_predicted)+1))
14 plt.yticks(range(10))
15 #设置网格线
16 plt.grid(alpha=0.2)
17 plt.legend(prop=my_font,loc="upper left")
18 plt.title('实际葡萄酒质量和预测葡萄酒质量对比',fontproperties=my_font, size = 16)
19 plt.show()
```

结果图：



从折线图上我们可以看到，预测结果和实际结果有些许不同，但总体上还是一致的