

assessing

October 17, 2017

1 Assessing

Use the space below to explore `winequality-red.csv` and `winequality-white.csv` to answer the quiz questions below.

```
In [5]: import pandas as pd
```

```
df_red = pd.read_csv('winequality-red.csv', sep=';')
df_red.head()
```

```
Out[5]:
```

	fixed_acidity	volatile_acidity	citric_acid	residual_sugar	chlorides	\
0	7.4	0.70	0.00	1.9	0.076	
1	7.8	0.88	0.00	2.6	0.098	
2	7.8	0.76	0.04	2.3	0.092	
3	11.2	0.28	0.56	1.9	0.075	
4	7.4	0.70	0.00	1.9	0.076	

	free_sulfur_dioxide	total_sulfur-dioxide	density	pH	sulphates	\
0	11.0	34.0	0.9978	3.51	0.56	
1	25.0	67.0	0.9968	3.20	0.68	
2	15.0	54.0	0.9970	3.26	0.65	
3	17.0	60.0	0.9980	3.16	0.58	
4	11.0	34.0	0.9978	3.51	0.56	

	alcohol	quality
0	9.4	5
1	9.8	5
2	9.8	5
3	9.8	6
4	9.4	5

```
In [8]: df_red.shape
```

```
Out[8]: (1599, 12)
```

```
In [13]: df_red.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
fixed_acidity      1599 non-null float64
volatile_acidity   1599 non-null float64
citric_acid        1599 non-null float64
residual_sugar     1599 non-null float64
chlorides          1599 non-null float64
free_sulfur_dioxide 1599 non-null float64
total_sulfur-dioxide 1599 non-null float64
density            1599 non-null float64
pH                 1599 non-null float64
sulphates          1599 non-null float64
alcohol            1599 non-null float64
quality            1599 non-null int64
dtypes: float64(11), int64(1)
memory usage: 150.0 KB

```

```

In [6]: df_white = pd.read_csv('winequality-white.csv', sep=';')
df_white.head()

```

```

Out[6]:
  fixed_acidity  volatile_acidity  citric_acid  residual_sugar  chlorides \
0             7.0              0.27         0.36             20.7      0.045
1             6.3              0.30         0.34              1.6      0.049
2             8.1              0.28         0.40              6.9      0.050
3             7.2              0.23         0.32              8.5      0.058
4             7.2              0.23         0.32              8.5      0.058

  free_sulfur_dioxide  total_sulfur_dioxide  density  pH  sulphates \
0                 45.0                 170.0   1.0010  3.00      0.45
1                 14.0                 132.0   0.9940  3.30      0.49
2                 30.0                 97.0   0.9951  3.26      0.44
3                 47.0                 186.0   0.9956  3.19      0.40
4                 47.0                 186.0   0.9956  3.19      0.40

  alcohol  quality
0       8.8       6
1       9.5       6
2      10.1       6
3       9.9       6
4       9.9       6

```

```

In [9]: df_white.shape

```

```

Out[9]: (4898, 12)

```

```

In [14]: df_white.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4898 entries, 0 to 4897
Data columns (total 12 columns):
fixed_acidity      4898 non-null float64
volatile_acidity   4898 non-null float64
citric_acid        4898 non-null float64
residual_sugar     4898 non-null float64
chlorides          4898 non-null float64
free_sulfur_dioxide 4898 non-null float64
total_sulfur_dioxide 4898 non-null float64
density            4898 non-null float64
pH                 4898 non-null float64
sulphates          4898 non-null float64
alcohol            4898 non-null float64
quality            4898 non-null int64
dtypes: float64(11), int64(1)
memory usage: 459.3 KB

```

```
In [15]: sum(df_red.duplicated())
```

```
Out[15]: 240
```

```
In [16]: sum(df_white.duplicated())
```

```
Out[16]: 937
```

```
In [17]: df_red.nunique()
```

```

Out[17]: 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 [18]: df_white.nunique()
```

```

Out[18]: fixed_acidity      68
         volatile_acidity   125
         citric_acid        87
         residual_sugar     310

```

```

chlorides          160
free_sulfur_dioxide 132
total_sulfur_dioxide 251
density            890
pH                 103
sulphates          79
alcohol            103
quality            7
dtype: int64

```

```
In [19]: df_red.describe()
```

```

Out[19]:
      fixed_acidity  volatile_acidity  citric_acid  residual_sugar  \
count      1599.000000      1599.000000  1599.000000      1599.000000
mean         8.319637         0.527821    0.270976         2.538806
std          1.741096         0.179060    0.194801         1.409928
min           4.600000         0.120000    0.000000         0.900000
25%           7.100000         0.390000    0.090000         1.900000
50%           7.900000         0.520000    0.260000         2.200000
75%           9.200000         0.640000    0.420000         2.600000
max          15.900000         1.580000    1.000000        15.500000

      chlorides  free_sulfur_dioxide  total_sulfur-dioxide  density  \
count      1599.000000      1599.000000      1599.000000  1599.000000
mean         0.087467        15.874922        46.467792    0.996747
std          0.047065        10.460157        32.895324    0.001887
min           0.012000         1.000000         6.000000    0.990070
25%           0.070000         7.000000        22.000000    0.995600
50%           0.079000        14.000000        38.000000    0.996750
75%           0.090000        21.000000        62.000000    0.997835
max           0.611000        72.000000       289.000000    1.003690

      pH  sulphates  alcohol  quality
count      1599.000000  1599.000000  1599.000000  1599.000000
mean         3.311113    0.658149    10.422983    5.636023
std          0.154386    0.169507    1.065668    0.807569
min           2.740000    0.330000    8.400000    3.000000
25%           3.210000    0.550000    9.500000    5.000000
50%           3.310000    0.620000   10.200000    6.000000
75%           3.400000    0.730000   11.100000    6.000000
max           4.010000    2.000000   14.900000    8.000000

```

```
In [ ]:
```