

refactor_wine_quality

November 23, 2019

1 Refactor: Wine Quality Analysis

In this exercise, you'll refactor code that analyzes a wine quality dataset taken from the UCI Machine Learning Repository [here](#). Each row contains data on a wine sample, including several physicochemical properties gathered from tests, as well as a quality rating evaluated by wine experts.

The code in this notebook first renames the columns of the dataset and then calculates some statistics on how some features may be related to quality ratings. Can you refactor this code to make it more clean and modular?

```
In [10]: import pandas as pd
         df = pd.read_csv('winequality-red.csv', sep=';')
         df.head()
```

```
Out[10]:
```

	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

1.0.1 Renaming Columns

You want to replace the spaces in the column labels with underscores to be able to reference columns with dot notation. Here's one way you could've done it.

```
In [11]: new_df = df.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'
                                     })

new_df.head()
```

```
Out[11]:
```

	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

And here's a slightly better way you could do it. You can avoid making naming errors due to typos caused by manual typing. However, this looks a little repetitive. Can you make it better?

```
In [23]: labels = list(df.columns)
         #print(labels)
         # labels[0] = labels[0].replace(' ', '_')
         # labels[1] = labels[1].replace(' ', '_')
         # labels[2] = labels[2].replace(' ', '_')
         # labels[3] = labels[3].replace(' ', '_')
         # labels[5] = labels[5].replace(' ', '_')
         # labels[6] = labels[6].replace(' ', '_')
         # df.columns = labels

         # df.head()
```

```

i = 0
for each in labels:
    labels[i] = each.replace(' ', '_')
    i+=1
df.columns = labels
#print(labels)
df.head()

```

```

Out[23]:
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

```

1.0.2 Analyzing Features

Now that your columns are ready, you want to see how different features of this dataset relate to the quality rating of the wine. A very simple way you could do this is by observing the mean quality rating for the top and bottom half of each feature. The code below does this for four features. It looks pretty repetitive right now. Can you make this more concise?

You might challenge yourself to figure out how to make this code more efficient! But you don't need to worry too much about efficiency right now - we will cover that more in the next section.

```

In [24]: # def groupby_quality(df_col,col_name):
#         median = df_col.median()
#         for i, col in enumerate(df_col):
#             if col >= median:
#                 df.loc[i,col_name] = 'high'
#             else:
#                 df.loc[i,col_name] = 'low'
#         df.groupby(col_name).quality.mean()

```

```

In [25]: median_alcohol = df.alcohol.median()
for i, alcohol in enumerate(df.alcohol):

```

```

        if alcohol >= median_alcohol:
            df.loc[i, 'alcohol'] = 'high'
        else:
            df.loc[i, 'alcohol'] = 'low'
    df.groupby('alcohol').quality.mean()

Out[25]: alcohol
high      5.958904
low       5.310302
Name: quality, dtype: float64

In [26]: median_pH = df.pH.median()
for i, pH in enumerate(df.pH):
    if pH >= median_pH:
        df.loc[i, 'pH'] = 'high'
    else:
        df.loc[i, 'pH'] = 'low'
df.groupby('pH').quality.mean()

Out[26]: pH
high      5.598039
low       5.675607
Name: quality, dtype: float64

In [27]: median_sugar = df.residual_sugar.median()
for i, sugar in enumerate(df.residual_sugar):
    if sugar >= median_sugar:
        df.loc[i, 'residual_sugar'] = 'high'
    else:
        df.loc[i, 'residual_sugar'] = 'low'
df.groupby('residual_sugar').quality.mean()

Out[27]: residual_sugar
high      5.665880
low       5.602394
Name: quality, dtype: float64

In [28]: median_citric_acid = df.citric_acid.median()
for i, citric_acid in enumerate(df.citric_acid):
    if citric_acid >= median_citric_acid:
        df.loc[i, 'citric_acid'] = 'high'
    else:
        df.loc[i, 'citric_acid'] = 'low'
df.groupby('citric_acid').quality.mean()

Out[28]: citric_acid
high      5.822360
low       5.447103
Name: quality, dtype: float64

In [ ]:

```