

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import r2_score
```

```
In [93]: red_wine = pd.read_csv('winequality-red.csv', sep=';')
white_wine = pd.read_csv('winequality-white.csv', sep=';')
```

Main objective of the analysis is to predict the quality of wine.

two datasets were created, using red and white wine samples. The inputs include objective tests (e.g. PH values) and the output is based on sensory data (median of at least 3 evaluations made by wine experts). Each expert graded the wine quality between 0 (very bad) and 10 (very excellent).

Two datasets one for red wine and one for white wine

Number of Instances: red wine - 1599; white wine - 4898.

Number of Attributes: 11 + output attribute

- 1- fixed acidity
 - 2- volatile acidity
 - 3- citric acid
 - 4- residual sugar
 - 5- chlorides
 - 6- free sulfur dioxide
 - 7- total sulfur dioxide
 - 8- density
 - 9- pH
 - 10- sulphates
 - 11- alcohol
- Output variable (based on sensory data):
- 12- quality (score between 0 and 10)

```
In [94]: red_wine.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: 156.0 KB
```

```
In [95]: red_wine.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 |

```
In [96]: white_wine.head()
```

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality |
|---|---------------|------------------|-------------|----------------|-----------|---------------------|----------------------|---------|------|-----------|---------|---------|
| 0 | 7.0 | 0.27 | 0.36 | 20.7 | 0.045 | 45.0 | 170.0 | 1.0010 | 3.00 | 0.45 | 8.8 | 6 |
| 1 | 6.3 | 0.30 | 0.34 | 1.6 | 0.049 | 14.0 | 132.0 | 0.9940 | 3.30 | 0.49 | 9.5 | 6 |
| 2 | 8.1 | 0.28 | 0.40 | 6.9 | 0.050 | 30.0 | 97.0 | 0.9951 | 3.26 | 0.44 | 10.1 | 6 |
| 3 | 7.2 | 0.23 | 0.32 | 8.5 | 0.058 | 47.0 | 186.0 | 0.9956 | 3.19 | 0.40 | 9.9 | 6 |
| 4 | 7.2 | 0.23 | 0.32 | 8.5 | 0.058 | 47.0 | 186.0 | 0.9956 | 3.19 | 0.40 | 9.9 | 6 |

check null values and data types

```
In [97]: red_wine.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: 156.0 KB
```

```
In [98]: white_wine.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4898 entries, 0 to 4897
Data columns (total 12 columns):
# Column Non-Null Count Dtype
---
0 fixed acidity 4898 non-null float64
1 volatile acidity 4898 non-null float64
2 citric acid 4898 non-null float64
3 residual sugar 4898 non-null float64
4 chlorides 4898 non-null float64
5 free sulfur dioxide 4898 non-null float64
6 total sulfur dioxide 4898 non-null float64
7 density 4898 non-null float64
8 pH 4898 non-null float64
9 sulphates 4898 non-null float64
10 alcohol 4898 non-null float64
11 quality 4898 non-null int64
dtypes: float64(11), int64(1)
memory usage: 459.3 KB
```

combine datasets in wine_df

add a new feature or column to each data set to indicate whether the wine is red or white.

```
In [99]: red_wine['color'] = 'red'
white_wine['color'] = 'white'
```

```
Out[99]:
```

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality | color |
|---|---------------|------------------|-------------|----------------|-----------|---------------------|----------------------|---------|------|-----------|---------|---------|-------|
| 0 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 | red |
| 1 | 7.8 | 0.88 | 0.00 | 2.6 | 0.098 | 25.0 | 67.0 | 0.9968 | 3.20 | 0.68 | 9.8 | 5 | red |
| 2 | 7.8 | 0.76 | 0.04 | 2.3 | 0.092 | 15.0 | 54.0 | 0.9970 | 3.26 | 0.65 | 9.8 | 5 | red |
| 3 | 11.2 | 0.28 | 0.56 | 1.9 | 0.075 | 17.0 | 60.0 | 0.9980 | 3.16 | 0.58 | 9.8 | 6 | red |
| 4 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 | red |

```
In [100]: white_wine['color'] = 'white'
white_wine.head()
```

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality | color |
|---|---------------|------------------|-------------|----------------|-----------|---------------------|----------------------|---------|------|-----------|---------|---------|-------|
| 0 | 7.0 | 0.27 | 0.36 | 20.7 | 0.045 | 45.0 | 170.0 | 1.0010 | 3.00 | 0.45 | 8.8 | 6 | white |
| 1 | 6.3 | 0.30 | 0.34 | 1.6 | 0.049 | 14.0 | 132.0 | 0.9940 | 3.30 | 0.49 | 9.5 | 6 | white |
| 2 | 8.1 | 0.28 | 0.40 | 6.9 | 0.050 | 30.0 | 97.0 | 0.9951 | 3.26 | 0.44 | 10.1 | 6 | white |
| 3 | 7.2 | 0.23 | 0.32 | 8.5 | 0.058 | 47.0 | 186.0 | 0.9956 | 3.19 | 0.40 | 9.9 | 6 | white |
| 4 | 7.2 | 0.23 | 0.32 | 8.5 | 0.058 | 47.0 | 186.0 | 0.9956 | 3.19 | 0.40 | 9.9 | 6 | white |

```
In [101]: # append dataframes
wine_df = red_wine.append(white_wine, sort=False)
# view dataframe to check for success
wine_df.head()
```

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality | color |
|---|---------------|------------------|-------------|----------------|-----------|---------------------|----------------------|---------|------|-----------|---------|---------|-------|
| 0 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 | red |
| 1 | 7.8 | 0.88 | 0.00 | 2.6 | 0.098 | 25.0 | 67.0 | 0.9968 | 3.20 | 0.68 | 9.8 | 5 | red |
| 2 | 7.8 | 0.76 | 0.04 | 2.3 | 0.092 | 15.0 | 54.0 | 0.9970 | 3.26 | 0.65 | 9.8 | 5 | red |
| 3 | 11.2 | 0.28 | 0.56 | 1.9 | 0.075 | 17.0 | 60.0 | 0.9980 | 3.16 | 0.58 | 9.8 | 6 | red |
| 4 | 7.4 | 0.70 | 0.00 | 1.9 | 0.076 | 11.0 | 34.0 | 0.9978 | 3.51 | 0.56 | 9.4 | 5 | red |

```
In [102]: wine_df.columns
```

```
Out[102]: Index(['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
        'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
        'pH', 'sulphates', 'alcohol', 'quality', 'color'],
        dtype='object')
```

```
In [103]: wine_df.to_csv('winequality_edited.csv', sep=';')
```

```
In [104]: wine_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6497 entries, 0 to 4897
Data columns (total 13 columns):
# Column Non-Null Count Dtype
---
0 fixed acidity 6497 non-null float64
1 volatile acidity 6497 non-null float64
2 citric acid 6497 non-null float64
3 residual sugar 6497 non-null float64
4 chlorides 6497 non-null float64
5 free sulfur dioxide 6497 non-null float64
6 total sulfur dioxide 6497 non-null float64
7 density 6497 non-null float64
8 pH 6497 non-null float64
9 sulphates 6497 non-null float64
10 alcohol 6497 non-null float64
11 quality 6497 non-null int64
12 color 6497 non-null object
dtypes: float64(11), int64(1), object(1)
memory usage: 710.6+ KB
```

```
In [105]: wine_df['quality'].unique()
```

```
Out[105]: array([5, 6, 7, 4, 8, 3, 9])
```

This problem could be modeled by classification or regression, we will test linear regression in this analysis

```
In [106]:
```

in this part i check the data quality, assess any issues in the data and fix it:

- null values in each column
- each column has the proper data type
- duplicate rows
- outliers
- distribution for each column (skewness)

check for nulls

```
In [105]: wine_df.isna().sum()
```

| | |
|----------------------|-------|
| 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 |
| color | 0 |
| dtype: | int64 |

check columns datatypes

```
In [106]: wine_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6497 entries, 0 to 4897
Data columns (total 13 columns):
# Column Non-Null Count Dtype
---
0 fixed acidity 6497 non-null float64
1 volatile acidity 6497 non-null float64
2 citric acid 6497 non-null float64
3 residual sugar 6497 non-null float64
4 chlorides 6497 non-null float64
5 free sulfur dioxide 6497 non-null float64
6 total sulfur dioxide 6497 non-null float64
7 density 6497 non-null float64
8 pH 6497 non-null float64
9 sulphates 6497 non-null float64
10 alcohol 6497 non-null float64
11 quality 6497 non-null int64
12 color 6497 non-null object
dtypes: float64(11), int64(1), object(1)
memory usage: 710.6+ KB
```

check duplicates and remove

```
In [107]: sum(wine_df.duplicated())
```

```
Out[107]: 1177
```

```
In [108]: wine_df.drop_duplicates(inplace=True)
```

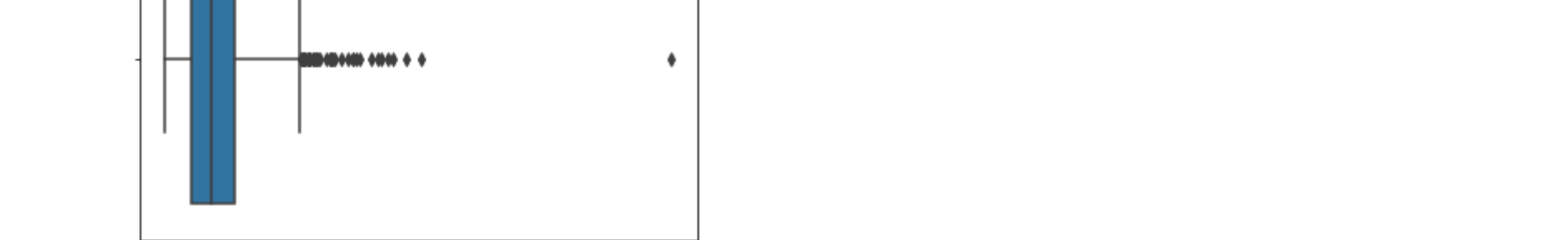
```
In [109]: sum(wine_df.duplicated())
```

```
Out[109]: 0
```

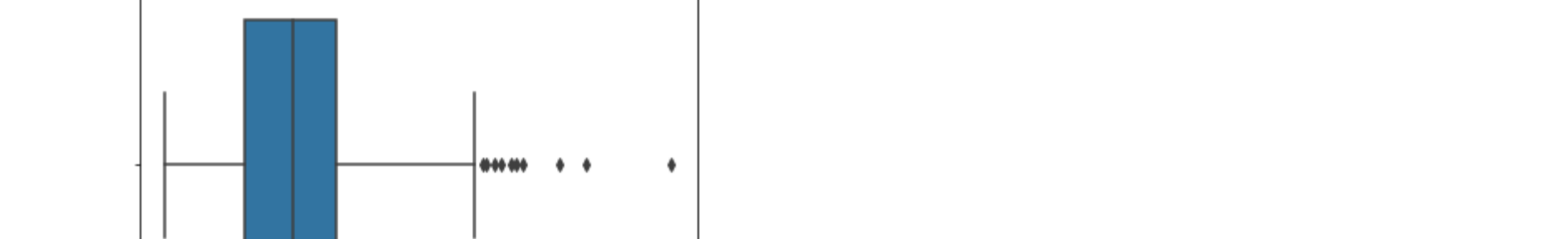
check for outliers and skewness

```
In [110]: columns_to_unskew = []
for col in wine_df.columns:
    if col != 'color':
        sns.boxplot(x = wine_df[col])
        plt.show()
```

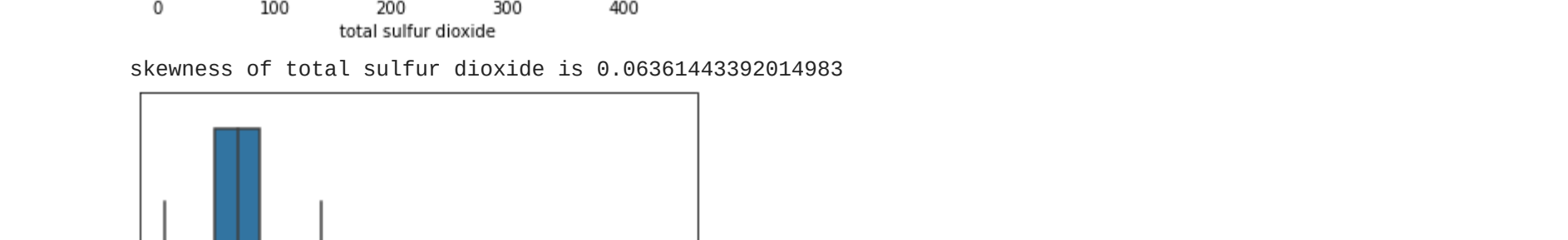
skew = wine_df[col].skew()
print("skewness of {} is {}".format(col, skew))
if skew > 3 or skew < -1:
 columns_to_unskew.append(col)



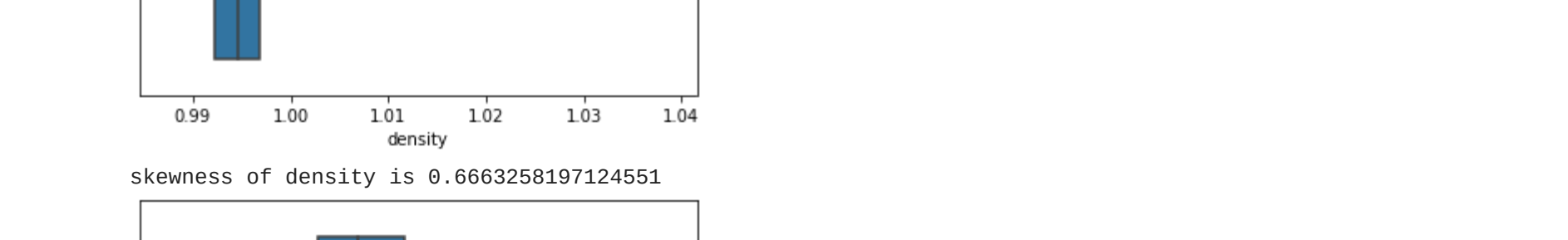
skewness of fixed acidity is 1.6504171812404804



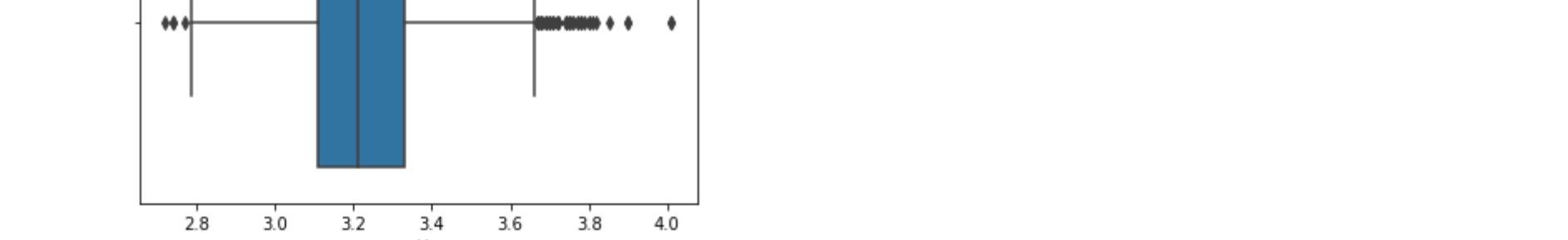
skewness of volatile acidity is 1.5045572014401647



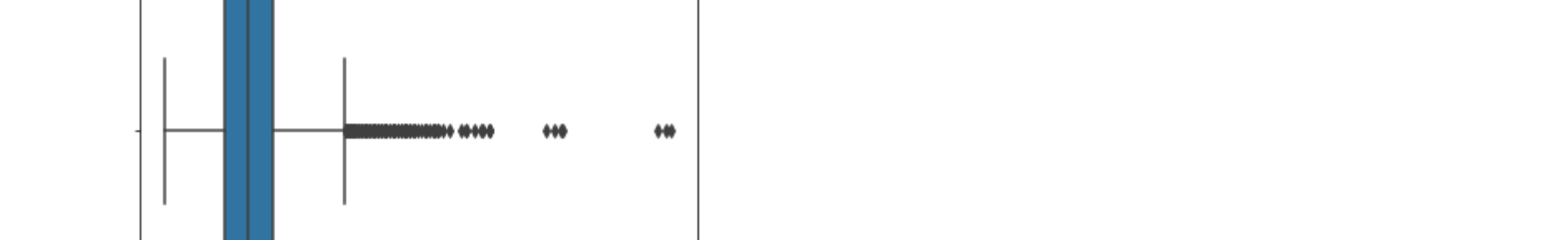
skewness of citric acid is 0.48545902797982095



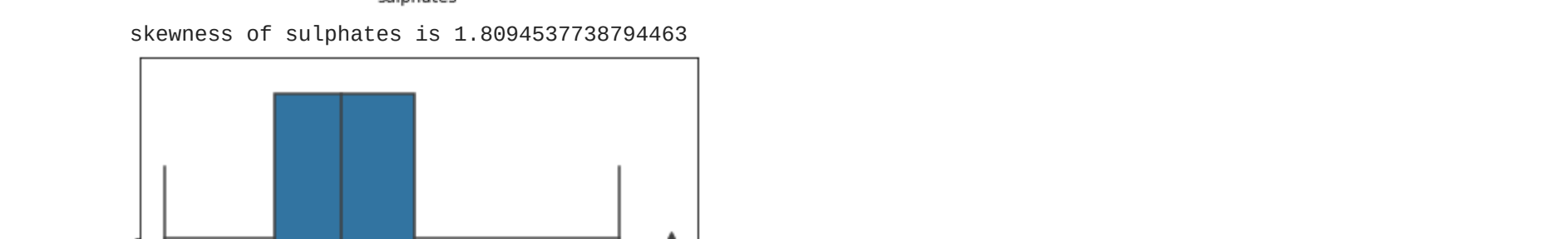
skewness of residual sugar is 1.7065502704869113



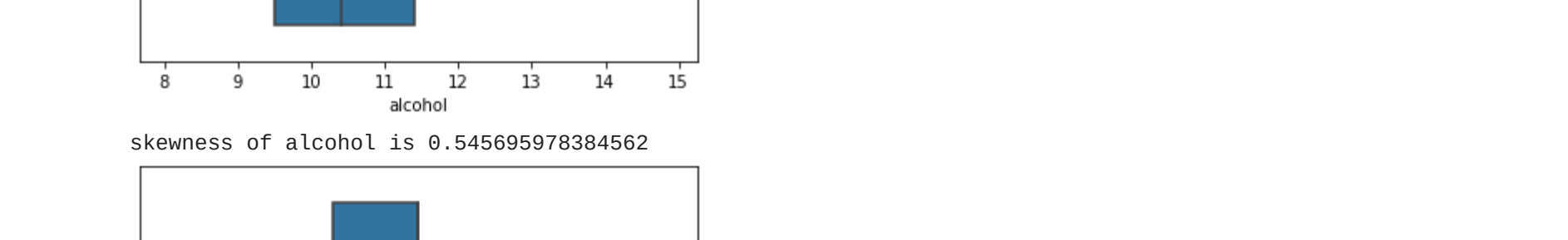
skewness of chlorides is 5.33823660791629



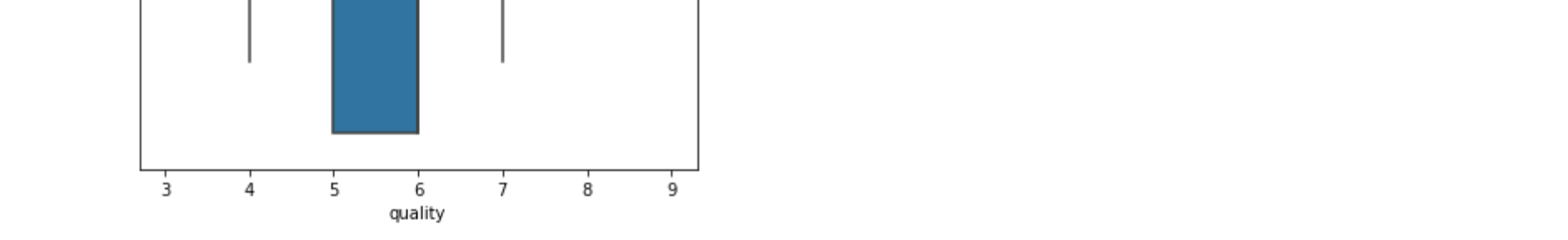
skewness of free sulfur dioxide is 1.362194619464467



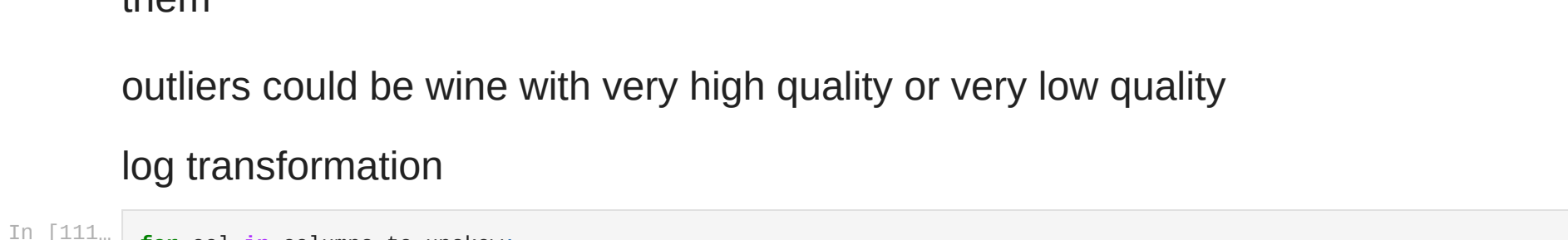
skewness of total sulfur dioxide is 0.8631443302014983



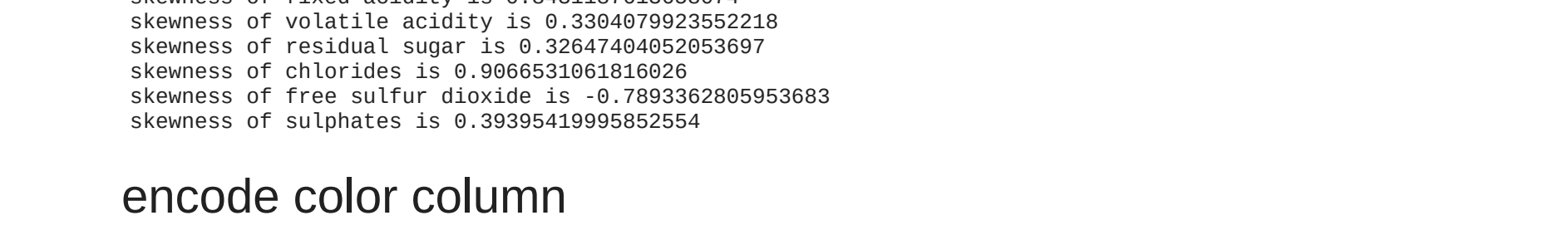
skewness of density is 0.6663258197124551



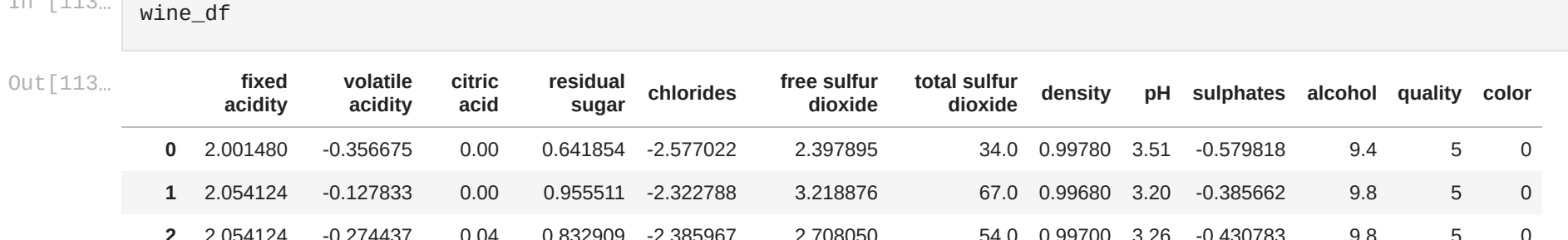
skewness of pH is 0.38996921434552967



skewness of sulphates is 1.8894537738794463



skewness of alcohol is 0.545695778384562



skewness of quality is 0.1474673665121148

most of the features are skewed and have lots of outliers so we can't remove them

logits transformation

```
In [111]: for col in columns_to_unskew:
wine_df[col] = np.log(wine_df[col])
print("skewness of {} is {}".format(col, wine_df[col].skew()))
```

skewness of fixed acidity is 0.8431387613658074
skewness of volatile acidity is 0.3304879923552218
skewness of citric acid is 0.3304879923552218
skewness of residual sugar is 0.966531063816026
skewness of chlorides is 0.789332895953683
skewness of free sulfur dioxide is 0.789332895953683
skewness of sulphates is 0.3939541999552554

encode color column

```
In [112]: wine_df['color'] = pd.get_dummies(wine_df['color'], drop_first=True)
```

```
In [113]: wine_df
```

| | fixed acidity | volatile acidity | citric acid | residual sugar | chlorides | free sulfur dioxide | total sulfur dioxide | density | pH | sulphates | alcohol | quality | color |
|---|---------------|------------------|-------------|----------------|-----------|---------------------|----------------------|---------|------|-----------|---------|---------|-------|
| 0 | 2.001480 | -0.356875 | 0.00 | 0.641854 | -2.577022 | 2.397895 | 34.0 | 0.99780 | 3.51 | -0.579818 | 9.4 | 5 | 0 |
| 1 | 2.054124 | -0.127833 | 0.00 | 0.955511 | -2.322788 | 3.218876 | 67.0 | 0.99680 | 3.20 | -0.385662 | 9.8 | 5 | 0 |
| 2 | 2.054124 | -0.274437 | 0.04 | 0.832900 | -2.385967 | 2.708050 | 54.0 | 0.99700 | 3.26 | -0.430783 | 9.8 | 5 | 0 |
| 3 | 2.415914 | -0.272966 | 0.56 | 0.641854 | -2.590267 | 2.833213 | 60.0 | 0.99800 | 3.16 | -0.544727 | 9.8 | 6 | 0 |
| 4 | 2.001480 | -0.415515 | 0.00 | 0.587787 | -2.590267 | 2.564949 | 40.0 | 0.99780 | 3.51 | -0.579818 | 9.4 | 5 | 0 |

5320 rows x 13 columns

```
In [160]: wine_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 | 0.220109 | 0.320073 | -0.068722 | 0.376426 | -0.326154 | -0.307004 | 0.483131 | -0.291425 | 0.281534 | -0.123079 | -0.091052 |
| volatile acidity | 0.220109 | 1.000000 | -0.340334 | -0.086402 | 0.443129 | -0.380076 | -0.365291 | 0.307104 | 0.201273 | 0.245456 | -0.048363 | -0.257036 |
| citric acid | 0.320073 | -0.340334 | 1.000000 | 0.122839 | -0.021140 | 0.116065 | 0.194835 | 0.094758 | -0.344735 | 0.037581 | -0.005496 | 0.097765 |
| residual sugar | -0.068722 | -0.086402 | 0.122839 | 1.000000 | -0.097723 | 0.365977 | 0.669693 | 0.480986 | -0.214883 | -0.154597 | -0.258628 | -0.027043 |
| chlorides | 0.376426 | 0.443129 | -0.021140 | -0.097723 | 1.000000 | -0.286262 | -0.310019 | 0.533019 | 0.109704 | 0.395847 | -0.382938 | -0.272067 |
| free sulfur dioxide | -0.326154 | -0.380076 | 0.116065 | 0.365977 | -0.286262 | 1.000000 | 0.740785 | 0.407085 | -0.139170 | -0.229993 | -0.124016 | 0.115441 |
| total sulfur dioxide | -0.307004 | -0.365291 | 0.194835 | 0.669693 | -0.310019 | 0.740785 | 1.000000 | 1.000000 | -0.222956 | -0.265884 | -0.490004 | 0.050296 |
| density | 0.483131 | 0.307104 | 0.094758 | 0.480986 | 0.533019 | -0.077548 | 0.006711 | 0.304273 | 0.034273 | -0.067811 | -0.326434 | 0.039733 |
| pH | -0.291425 | 0.201273 | -0.344735 | -0.021140 | -0.139170 | -0.222956 | 0.034273 | 1.000000 | 0.209510 | 0.097314 | -0.039733 | 0.039733 |
| sulphates | 0.281534 | 0.245456 | 0.037581 | -0.154597 | 0.395847 | -0.229993 | -0.265884 | 0.302393 | 0.209510 | 1.000000 | -0.038807 | 0.031693 |
| alcohol | -0.123079 | -0.048363 | -0.005496 | -0.258628 | -0.382938 | -0.124016 | -0.490004 | -0.667811 | 0.097314 | -0.038807 | 1.000000 | 0.469422 |
| quality | -0.091052 | -0.257036 | 0.097765 | -0.027043 | -0.272067 | 0.115441 | -0.050296 | -0.326434 | 0.039733 | 0.031693 | 0.469422 | 1.000000 |
| color | -0 | | | | | | | | | | | |


```
Out[130]: /usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 5.158e+01, tolerance: 1.018e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 2.367e+00, tolerance: 1.838e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 6.055e+01, tolerance: 1.799e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 1.241e+01, tolerance: 1.819e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 6.072e+00, tolerance: 1.038e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 6.967e+00, tolerance: 1.799e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 8.672e+00, tolerance: 1.038e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 6.967e+00, tolerance: 1.838e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 2.363e+00, tolerance: 1.838e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 3.614e+00, tolerance: 1.038e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 4.361e+00, tolerance: 1.838e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 9.764e+00, tolerance: 1.799e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive
/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_coordinate_descent.py:648: ConvergenceWarning: Ob
ective did not converge. You might want to increase the number of iterations, check the scale of the features o
r consider increasing regularisation. Duality gap: 7.448e+00, tolerance: 2.728e-01
  coef_, l1_reg, l2_reg, X, y, max_iter, tol, rng, random, positive

Out[132]: GridSearchCV(cv=3,
  estimator=Pipeline(steps=[('polynomial', PolynomialFeatures()),
                             ('scaler', StandardScaler()),
                             ('ridge', Ridge())]),
  param_grid=[('polynomial_degree': [1, 2, 3, 4, 5, 6],
               ('ridge__alpha': [0.01, 0.001, 0.0001, 1, 5, 10, 20],
                'ridge__max_iter': [1000, 2000, 3000, 4000])]),
               ('lasso_alpha': [0.01, 0.001, 0.0001, 1, 5, 10, 20],
                'lasso_max_iter': [2000, 3000],
                'polynomial_degree': [1, 2, 3, 4, 5])]),
  scoring='neg_mean_squared_error',
  verbose=0)
```

```
In [132]: grid_cv.best_params_

Out[132]: {'lasso_alpha': 0.001, 'lasso_max_iter': 2000, 'polynomial_degree': 3}

In [133]: print("Train score: ", grid_cv.score(X_train, y_train))

Train score: 0.3978064888995494

In [134]: ypred = grid_cv.predict(X_test)
ridge_score = r2_score(y_test, ypred)
scores_dict['lasso_poly'] = lasso_score

In [135]: scores_dict

Out[135]: {'LR': 0.3210036115533088,
  'lasso_poly': 0.364958402889341846,
  'poly': 0.32100361155335937,
  'ridge': 0.3223181818486342853}

let's check some of the output

In [158]: np.append(y_pred, y_test, axis=1)[:10]

Out[158]: array([[6.51201689, 7.        ],
  [5.52979997, 6.        ],
  [5.36506128, 6.        ],
  [4.887849   , 4.        ],
  [5.08026871, 5.        ],
  [5.88433399, 6.        ],
  [5.49176547, 5.        ],
  [5.32022178, 5.        ],
  [4.9881621 , 5.        ],
  [5.89402777, 5.        ]])
```

opservation and flaws

- 1- several of the attributes are correlated
- 2- Not all features are relevant that's why feature selection gave better results
- 3- Data had lots of outliers

Best model

the best score on the testset was achieved by ridge model with a polynomial degree of 3. but still the score was only 37% which is not good.

Linear regression might not be the best fit for this dataset. We can Test classification techniuges or a different regression model.

future steps

- 1- Outlier detection algorithms could be used to detect the few excellent or poor wines.
- 2- use classification algorithms
- 3- other regression models

```
In [ ]:

In [ ]:

In [ ]:
```