

Jakie cechy jabka wpływają na postrzeganie przez ludzi jego jakości? Czy istnieją korelacje między zmienna najbardziej wpływającą na postrzeganie jakości a innymi zmiennym? Wyćwiczenie modelu wykrywającego jakość jabłka.

Sprawdzam podstawowe parametry, oczyszczam dane.

```
In [31]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import ttest_ind
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn import metrics

In [2]: df = pd.read_csv(r'C:\Users\A8797\Documents\ProjektyPython\AnalizaJabka\apple_quality.csv')
```

```
In [3]: df.head()

Out[3]:
```

| | A_id | Size | Weight | Sweetness | Crunchiness | Juiciness | Ripeness | Acidity | Quality |
|---|------|-----------|-----------|-----------|-------------|-----------|-----------|--------------|---------|
| 0 | 0.0 | -3.970049 | -2.512336 | 5.346330 | -1.012009 | 1.844900 | 0.329840 | -0.491590483 | good |
| 1 | 1.0 | -1.195217 | -2.839257 | 3.664059 | 1.588232 | 0.853286 | 0.867530 | -0.722809367 | good |
| 2 | 2.0 | -0.292024 | -1.351282 | -1.738429 | -0.342616 | 2.838636 | -0.038033 | 2.621636473 | bad |
| 3 | 3.0 | -0.657196 | -2.271627 | 1.324874 | -0.097875 | 3.637970 | -3.413761 | 0.790723217 | good |
| 4 | 4.0 | 1.364217 | -1.296612 | -0.384658 | -0.553006 | 3.030874 | -1.303849 | 0.501984036 | good |

```
In [4]: df.describe()

Out[4]:
```

| | A_id | Size | Weight | Sweetness | Crunchiness | Juiciness | Ripeness |
|-------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| count | 4000.000000 | 4000.000000 | 4000.000000 | 4000.000000 | 4000.000000 | 4000.000000 | 4000.000000 |
| mean | 1999.500000 | -0.503015 | -0.989547 | -0.470479 | 0.985478 | 0.512118 | 0.438277 |
| std | 1154.844867 | 1.928059 | 1.602507 | 1.943441 | 1.402757 | 1.930286 | 1.874427 |
| min | 0.000000 | -7.151703 | -7.149848 | -6.894485 | -6.055058 | -5.961897 | -5.864599 |
| 25% | 999.750000 | -1.816765 | -2.011770 | -1.738425 | 0.062764 | -0.801286 | -0.771677 |
| 50% | 1999.500000 | -0.513703 | -0.984736 | -0.504758 | 0.998249 | 0.534219 | 0.503445 |
| 75% | 2999.250000 | 0.805526 | 0.030976 | 0.801822 | 1.894234 | 1.835976 | 1.766212 |
| max | 3999.000000 | 6.406367 | 5.790714 | 6.374916 | 7.619852 | 7.364403 | 7.237837 |

```
In [5]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4000 entries, 0 to 4000
Data columns (total 9 columns):
#   Column      Non-Null Count  Dtype
---  --
0   A_id         4000 non-null    float64
1   Size         4000 non-null    float64
2   Weight       4000 non-null    float64
3   Sweetness    4000 non-null    float64
4   Crunchiness  4000 non-null    float64
5   Juiciness    4000 non-null    float64
6   Ripeness     4000 non-null    float64
7   Acidity      4000 non-null    object
8   Quality      4000 non-null    object
dtypes: float64(7), object(2)
memory usage: 281.4+ KB
```

```
In [6]: df = df.drop(labels = 'A_id', axis = 1)
```

```
In [7]: df.isna()
```

```
Out[7]:
```

| | Size | Weight | Sweetness | Crunchiness | Juiciness | Ripeness | Acidity | Quality |
|------|-------|--------|-----------|-------------|-----------|----------|---------|---------|
| 0 | False | False | False | False | False | False | False | False |
| 1 | False | False | False | False | False | False | False | False |
| 2 | False | False | False | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 3996 | False | False | False | False | False | False | False | False |
| 3997 | False | False | False | False | False | False | False | False |
| 3998 | False | False | False | False | False | False | False | False |
| 3999 | False | False | False | False | False | False | False | False |
| 4000 | True | True | True | True | True | True | False | True |

4001 rows x 8 columns

```
In [8]: df = df.dropna()
```

```
In [9]: df['Quality'] = df['Quality'] == 'good'.astype(int)
```

```
In [10]: df['Acidity'] = df['Acidity'].astype(float)
```

```
In [11]: df.info()

<class 'pandas.core.frame.DataFrame'>
Index: 4000 entries, 0 to 3999
Data columns (total 8 columns):
#   Column      Non-Null Count  Dtype
---  --
0   Size         4000 non-null    float64
1   Weight       4000 non-null    float64
2   Sweetness    4000 non-null    float64
3   Crunchiness  4000 non-null    float64
4   Juiciness    4000 non-null    float64
5   Ripeness     4000 non-null    float64
6   Acidity      4000 non-null    float64
7   Quality      4000 non-null    int32
dtypes: float64(7), int32(1)
memory usage: 285.6 KB

Sprawdzam korelację wszystkich zmiennych.
```

```
In [12]: df.corr()
```

```
Out[12]:
```

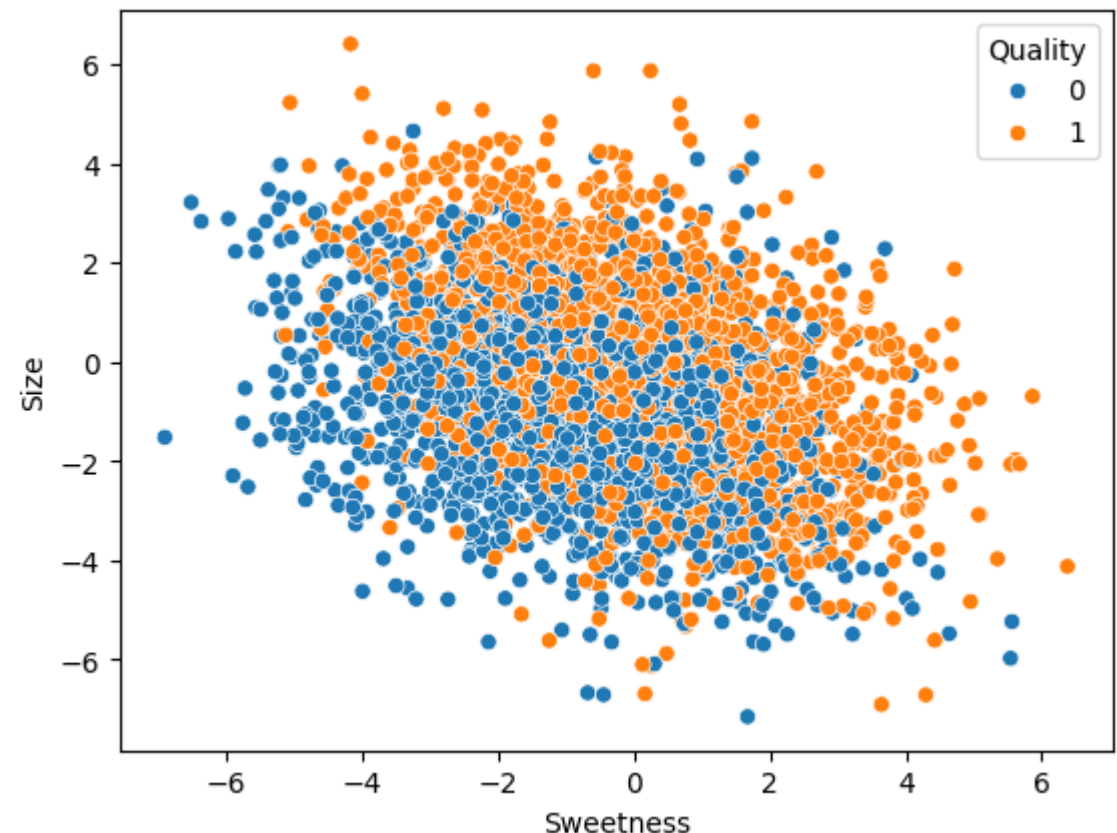
| | Size | Weight | Sweetness | Crunchiness | Juiciness | Ripeness | Acidity | Quality |
|-------------|-----------|-----------|-----------|-------------|-----------|-----------|-----------|-----------|
| Size | 1.000000 | -0.170702 | -0.324680 | 0.169868 | -0.018892 | -0.134773 | 0.196218 | 0.244007 |
| Weight | -0.170702 | 1.000000 | -0.154246 | -0.095882 | -0.092263 | -0.243824 | 0.016414 | 0.001421 |
| Sweetness | -0.324680 | -0.154246 | 1.000000 | -0.037552 | 0.095882 | 0.273800 | 0.085999 | 0.250998 |
| Crunchiness | 0.169868 | -0.095882 | -0.037552 | 1.000000 | -0.259607 | -0.201982 | 0.069943 | -0.012376 |
| Juiciness | -0.018892 | -0.092263 | 0.095882 | -0.259607 | 1.000000 | -0.097144 | 0.248714 | 0.260223 |
| Ripeness | -0.134773 | -0.243824 | -0.273800 | -0.201982 | -0.097144 | 1.000000 | -0.202669 | -0.264315 |
| Acidity | 0.196218 | 0.016414 | 0.085999 | 0.069943 | 0.248714 | -0.202669 | 1.000000 | -0.007697 |
| Quality | 0.244007 | 0.001421 | 0.250998 | -0.012376 | 0.260223 | -0.264315 | -0.007697 | 1.000000 |

Korelacja Pearsona wskazuje występowanie słabych korelacji pozytywnych między jakością jabłka a jego wielkością, słodkością i soczystością oraz słabą korelację ujemną między jakością a dojrzałością jabłka.

Wizualizację korelacji możemy narysować między słodkością a rozmiarem jabłka.

```
In [13]: sns.scatterplot(data = df, x = 'Sweetness', y = 'Size', hue = 'Quality')
```

```
Out[13]: <Axes: xlabel='Sweetness', ylabel='Size'>
```

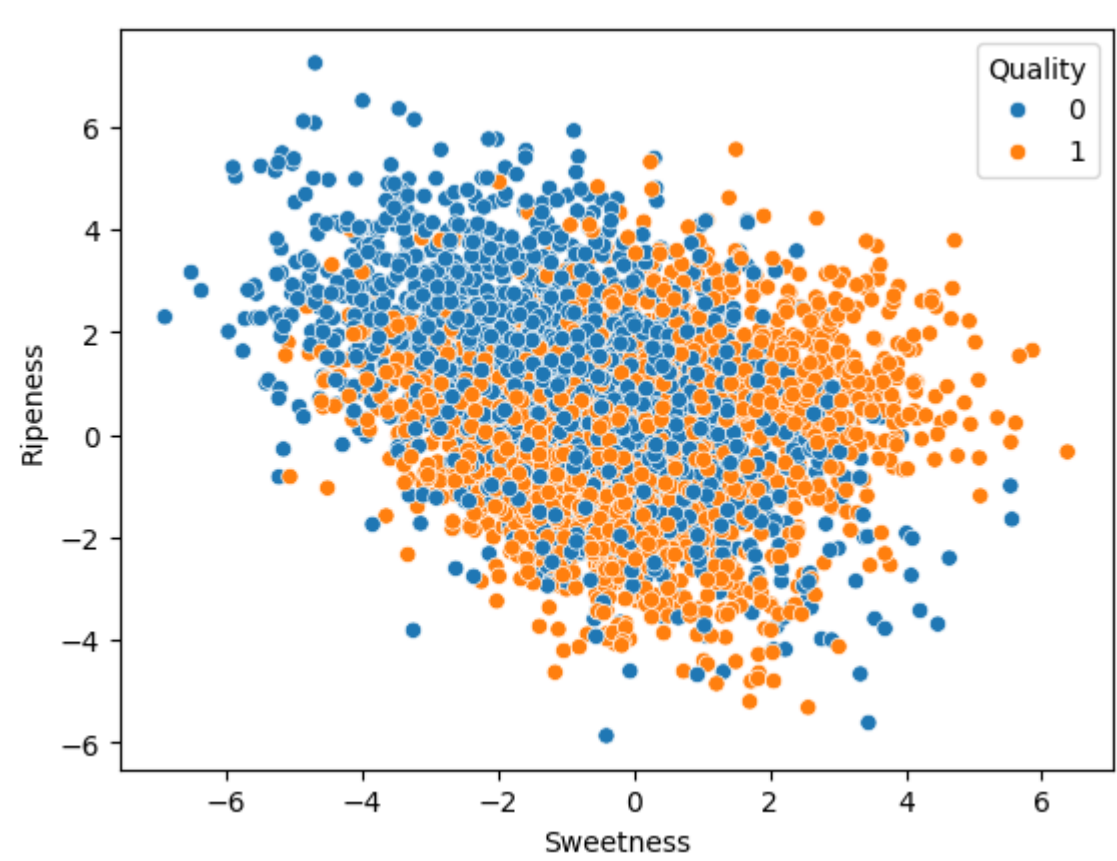


Wniosek: Im jabłko mniejsze, tym według respondentów jest słodsze oraz jest lepszej jakości

Wizualizuję korelację ujemną między słodkością a dojrzałością jabłka.

```
In [14]: sns.scatterplot(data = df, x = 'Sweetness', y = 'Ripeness', hue = 'Quality')
```

```
Out[14]: <Axes: xlabel='Sweetness', ylabel='Ripeness'>
```

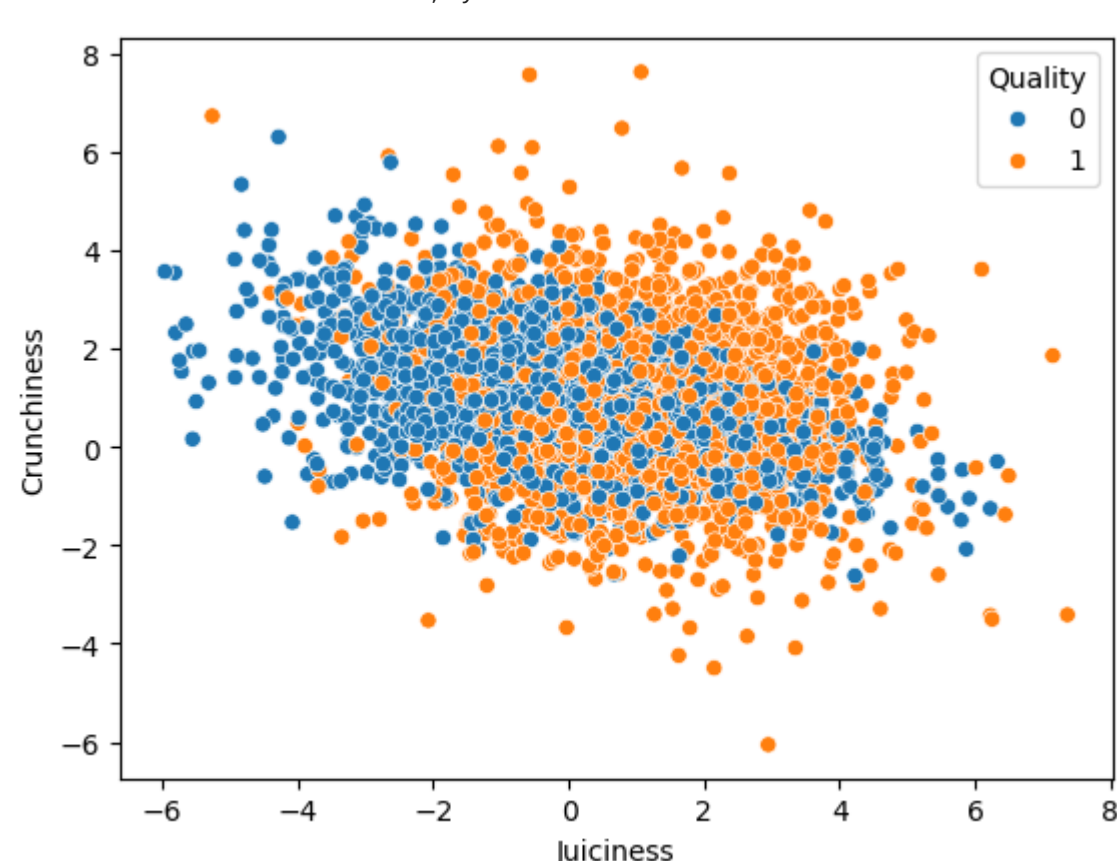


Wniosek: Im jabłko mniej dojrzałe, tym jest uznawane za słodsze

Wizualizuję korelację ujemną między soczystością a kruchością jabłka.

```
In [15]: sns.scatterplot(data = df, x = 'Juiciness', y = 'Crunchiness', hue = 'Quality')
```

```
Out[15]: <Axes: xlabel='Juiciness', ylabel='Crunchiness'>
```

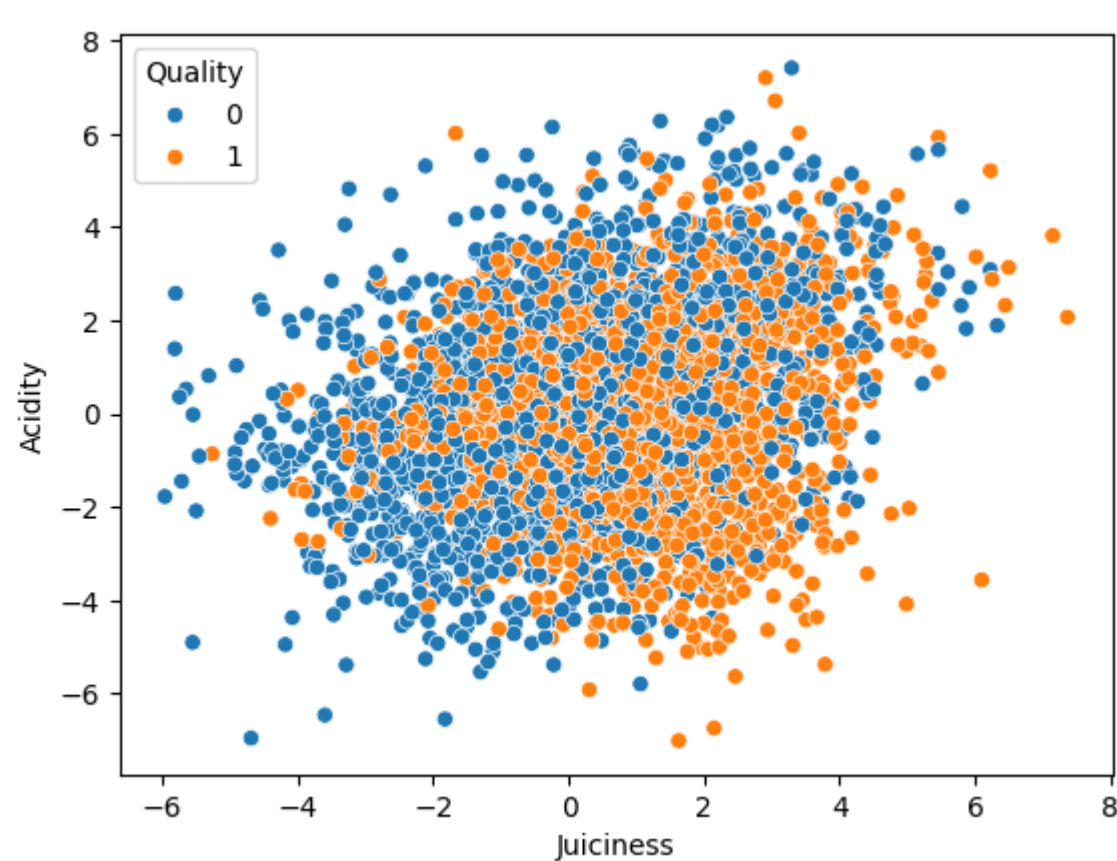


Wniosek: Im jabłko mniej chrupkie, tym jest uznawane za soczyste.

Wizualizuję korelację pozytywną między soczystością a kwasowością jabłka.

```
In [16]: sns.scatterplot(data = df, x = 'Juiciness', y = 'Acidity', hue = 'Quality')
```

```
Out[16]: <Axes: xlabel='Juiciness', ylabel='Acidity'>
```



Wniosek: Im jabłko bardziej kwaśne, tym jest uznawane za soczyste.

```
In [19]: grupa_dobre = df[df['Quality'] == 1]['Sweetness'].values
grupa_zle = df[df['Quality'] == 0]['Sweetness'].values
statystyka, p_value = ttest_ind(grupa_dobre, grupa_zle)
```

```
if p_value < 0.05:
    print('Różnice są istotne statystycznie.')
else:
    print('Brak istotnych statystycznie różnic.')
```

Różnice są istotne statystycznie.

```
In [20]: grupa_dobre = df[df['Quality'] == 1]['Juiciness'].values
grupa_zle = df[df['Quality'] == 0]['Juiciness'].values
statystyka, p_value = ttest_ind(grupa_dobre, grupa_zle)
```

```
if p_value < 0.05:
    print('Różnice są istotne statystycznie.')
else:
    print('Brak istotnych statystycznie różnic.')
```

Różnice są istotne statystycznie.

```
In [21]: grupa_dobre = df[df['Quality'] == 1]['Size'].values
grupa_zle = df[df['Quality'] == 0]['Size'].values
statystyka, p_value = ttest_ind(grupa_dobre, grupa_zle)
```

```
if p_value < 0.05:
    print('Różnice są istotne statystycznie.')
else:
    print('Brak istotnych statystycznie różnic.')
```

Różnice są istotne statystycznie.

```
In [22]: grupa_dobre = df[df['Quality'] == 1]['Ripeness'].values
grupa_zle = df[df['Quality'] == 0]['Ripeness'].values
statystyka, p_value = ttest_ind(grupa_dobre, grupa_zle)
```

```
if p_value < 0.05:
    print('Różnice są istotne statystycznie.')
else:
    print('Brak istotnych statystycznie różnic.')
```

Różnice są istotne statystycznie.

Wszystkie zmienne wpływające na jakość (posiadają korelację) jabłka są istotne statystycznie.

```
In [23]: df1 = df[df['Quality'] == 1]
```

```
In [24]: df1
```

```
Out[24]:
```

| | Size | Weight | Sweetness | Crunchiness | Juiciness | Ripeness | Acidity | Quality |
|------|-----------|-----------|-----------|-------------|-----------|-----------|-----------|---------|
| 0 | -3.970049 | -2.512336 | 5.346330 | -1.012009 | 1.844900 | 0.329840 | -0.491590 | 1 |
| 1 | -1.195217 | -2.839257 | 3.664059 | 1.588232 | 0.853286 | 0.867530 | -0.722809 | 1 |
| 3 | -0.657196 | -2.271627 | 1.324874 | -0.097875 | 3.637970 | -3.413761 | 0.790723 | 1 |
| 4 | 1.364217 | -1.296612 | -0.384658 | -0.553006 | 3.030874 | -1.303849 | 0.501984 | 1 |
| 6 | 1.331606 | 1.635956 | 0.875974 | -1.677798 | 3.106344 | -1.847417 | 2.414171 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 3992 | 1.764253 | -2.079695 | -0.083383 | -0.086724 | -1.703385 | 3.840101 | -0.338261 | 1 |
| 3994 | 1.482508 | -2.581181 | -0.306888 | 1.527877 | 1.056361 | 2.560829 | -1.232925 | 1 |
| 3996 | -0.280118 | 1.949253 | -0.204020 | -0.640196 | 0.024523 | -1.087900 | 1.854235 | 1 |
| 3998 | -4.008004 | -1.779337 | 2.366397 | -0.200329 | 2.161435 | 0.214488 | -2.229720 | 1 |
| 3999 | 0.278540 | -1.715505 | 0.121217 | -1.154075 | 1.266677 | -0.776571 | 1.599796 | 1 |

2004 rows x 8 columns

```
In [25]: df1.describe()
```

```
Out[25]:
```

| | Size | Weight | Sweetness | Crunchiness | Juiciness | Ripeness | Acidity | Quality |
|-------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|---------|
| count | 1996.000000 | 1996.000000 | 1996.000000 | 1996.000000 | 1996.000000 | 1996.000000 | 1996.000000 | 1996.0 |
| mean | -0.974358 | -0.987274 | 0.016287 | 0.968154 | 1.013356 | 0.003890 | 0.060668 | 1.0 |
| std | 1.972189 | 1.852145 | 1.910178 | 1.599191 | 1.688522 | 1.729708 | 2.006745 | 0.0 |
| min | -6.905803 | -7.149848 | -5.118948 | -6.055058 | -5.261636 | -5.313838 | -7.010538 | 1.0 |
| 25% | -1.379540 | -2.220975 | -1.281629 | -0.126736 | -0.108738 | -1.194310 | -1.332003 | 1.0 |
| 50% | 0.022554 | -0.994428 | -0.116781 | 1.049431 | 1.024190 | 0.095671 | 0.060766 | 1.0 |
| 75% | 1.344088 | 0.291035 | 1.236140 | 2.044560 | 2.148908 | 1.176599 | 1.454957 | 1.0 |
| max | 6.406367 | 5.790714 | 6.374916 | 7.619852 | 7.364403 | 5.553256 | 7.193374 | 1.0 |

```
In [26]: df2 = df[df['Quality'] == 0]
```

```
In [27]: df2.describe()
```

```
Out[27]:
```

| | Size | Weight | Sweetness | Crunchiness | Juiciness | Ripeness | Acidity | Quality |
|-------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|---------|
| count | 1996.000000 | 1996.000000 | 1996.000000 | 1996.000000 | 1996.000000 | 1996.000000 | 1996.000000 | 1996.0 |
| mean | -0.974358 | -0.987274 | -0.959195 | 1.002871 | 0.008871 | 0.994646 | 0.093151 | 0.0 |
| std | 1.761477 | 1.305158 | 1.852182 | 1.173057 | 2.024136 | 1.883318 | 2.209370 | 0.0 |
| min | -7.151703 | -7.463384 | -8.894485 | -2.620954 | -5.261897 | -5.864599 | -6.955460 | 0.0 |
| 25% | -2.169909 | -1.837784 | -2.194897 | 0.249392 | -1.374461 | -0.231485 | -1.427181 | 0.0 |
| 50% | -1.019633 | -0.975703 | -0.921758 | 0.977815 | 0.005241 | 0.965054 | -0.010807 | 0.0 |
| 75% | 0.221546 | -0.148292 | 0.312889 | 1.791937 | 1.432981 | 2.277946 | 1.595159 | 0.0 |
| max | 4.649923 | 3.081538 | 5.559624 | 6.297873 | 6.328304 | 7.237837 | 7.404736 | 0.0 |

```
In [32]: X = df[['Size','Weight','Sweetness','Crunchiness','Juiciness','Ripeness','Acidity']]
y = df['Quality']
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=16)
logreg = LogisticRegression(random_state=16)
```

```
logreg.fit(X_train, y_train)
y_pred = logreg.predict(X_test)
```

```
cnf_matrix = metrics.confusion_matrix(y_test, y_pred)
cmf_matrix
```

```
Out[32]: array([[349, 150],
               [120, 381]], dtype=int64)
```

Dokładność modelu wynosi ~73%

Próba użycia modelu

```
In [39]: new_data = [[-3, -2, -1, -2, -5, -2, -2]]
predicted_class = logreg.predict(new_data)
print(predicted_class)
```

```
[0]
C:\Users\A8797\AppData\Local\Packages\PythonSoftwareFoundation.Python.3.11.qbz2nxfrapp\LocalCache\local-packages\Python311\site-packages\sklearn\base.py:493: UserWarning: X does not have valid feature names, but LogisticRegression was fitted with feature names
warnings.warn(C
```