Лабортаторная работа №1 по курсо ТМО

Бекетов Роман

ИУ5-62Б

# Разведочный анализ данных. Исследование и визуализация данных.

1) Текстовое описание набора данных

Этот датасет содержит информацию о различных атрибутах набора фруктов - яблоков, позволяющую получить представление об их характеристиках. Набор данных включает такие сведения, как идентификатор фрукта, размер, вес, сладость, хрусткость, сочность, спелость, кислотность и качество.

```python
# !pip install numpy pandas seaborn matplotlib

# !pip install scipy

import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(style="ticks")

df_data = pd.read_csv("apple_quality.csv")

df_data.sample(5)
```

```
        A_id      Size    Weight  Sweetness  Crunchiness  Juiciness
Ripeness  \
3144  3144.0 -2.248008 -2.310955  -0.430174     1.798926  -0.690108
0.448416
1552  1552.0 -0.196499 -1.534853   0.828281     1.717767   0.838220
1.396164
2336  2336.0 -1.286048 -2.465491  -0.889402     0.661327  -0.833476
1.743180
2488  2488.0  2.094886 -2.762798  -3.236776     2.443926   0.935719 -
0.044279
1853  1853.0 -1.205505  2.423058   3.473933    -1.469719  -2.792099 -
2.546413

          Acidity Quality
3144  -2.872017555     bad
1552   0.643928312    good
2336  -2.830523857     bad
```

```
2488   -1.008545765      good
1853    2.871475348      good
```

df_data.shape

(4001, 9)

df_data.columns

```
Index(['A_id', 'Size', 'Weight', 'Sweetness', 'Crunchiness',
'Juiciness',
        'Ripeness', 'Acidity', 'Quality'],
       dtype='object')
```

df_data.dtypes

```
A_id           float64
Size           float64
Weight         float64
Sweetness      float64
Crunchiness    float64
Juiciness      float64
Ripeness       float64
Acidity         object
Quality         object
dtype: object
```

```python
print("Количесво пропусков")
for col in df_data:
    print(f"{col} = {df_data[df_data[col].isnull()].shape[0]}")
```

```
Количесво пропусков
A_id = 1
Size = 1
Weight = 1
Sweetness = 1
Crunchiness = 1
Juiciness = 1
Ripeness = 1
Acidity = 0
Quality = 1
```

df_data.describe()

|       | A_id | Size | Weight | Sweetness | Crunchiness |
|-------|------|------|--------|-----------|-------------|
| count | 4000.000000 | 4000.000000 | 4000.000000 | 4000.000000 | 4000.000000 |
| mean  | 1999.500000 | -0.503015 | -0.989547 | -0.470479 | 0.985478 |
| std   | 1154.844867 | 1.928059 | 1.602507 | 1.943441 | 1.402757 |

|       |              |           |           |           |           |
|-------|--------------|-----------|-----------|-----------|-----------|
| min   | 0.000000     | -7.151703 | -7.149848 | -6.894485 | -6.055058 |
| 25%   | 999.750000   | -1.816765 | -2.011770 | -1.738425 | 0.062764  |
| 50%   | 1999.500000  | -0.513703 | -0.984736 | -0.504758 | 0.998249  |
| 75%   | 2999.250000  | 0.805526  | 0.030976  | 0.801922  | 1.894234  |
| max   | 3999.000000  | 6.406367  | 5.790714  | 6.374916  | 7.619852  |

```
          Juiciness      Ripeness
count   4000.000000   4000.000000
mean       0.512118      0.498277
std        1.930286      1.874427
min       -5.961897     -5.864599
25%       -0.801286     -0.771677
50%        0.534219      0.503445
75%        1.835976      1.766212
max        7.364403      7.237837
```

```python
df_data.Quality.unique()
```

```
array(['good', 'bad', nan], dtype=object)
```

```python
fig, ax = plt.subplots(figsize=(10,10))
sns.scatterplot(ax=ax, x='Sweetness', y='Size', data=df_data)
```

```
<Axes: xlabel='Sweetness', ylabel='Size'>
```

```
fig, ax = plt.subplots(figsize=(10,10))
sns.distplot(df_data['Sweetness'])
```

/var/folders/bg/1zs8qp8d26v62zscyhsgmff40000gq/T/
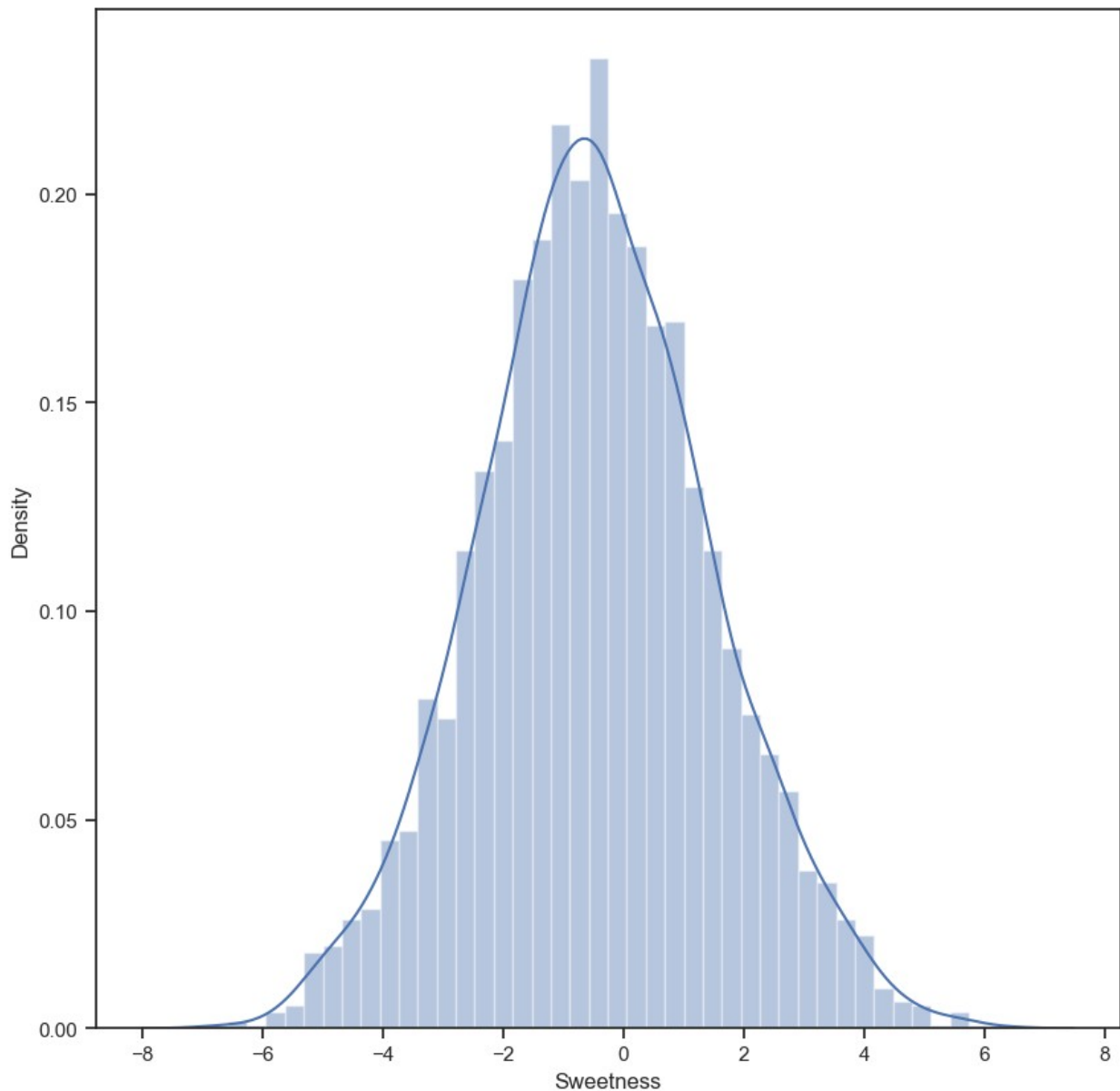ipykernel_2543/3326567540.py:2: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn
v0.14.0.

Please adapt your code to use either `displot` (a figure-level
function with
similar flexibility) or `histplot` (an axes-level function for
histograms).
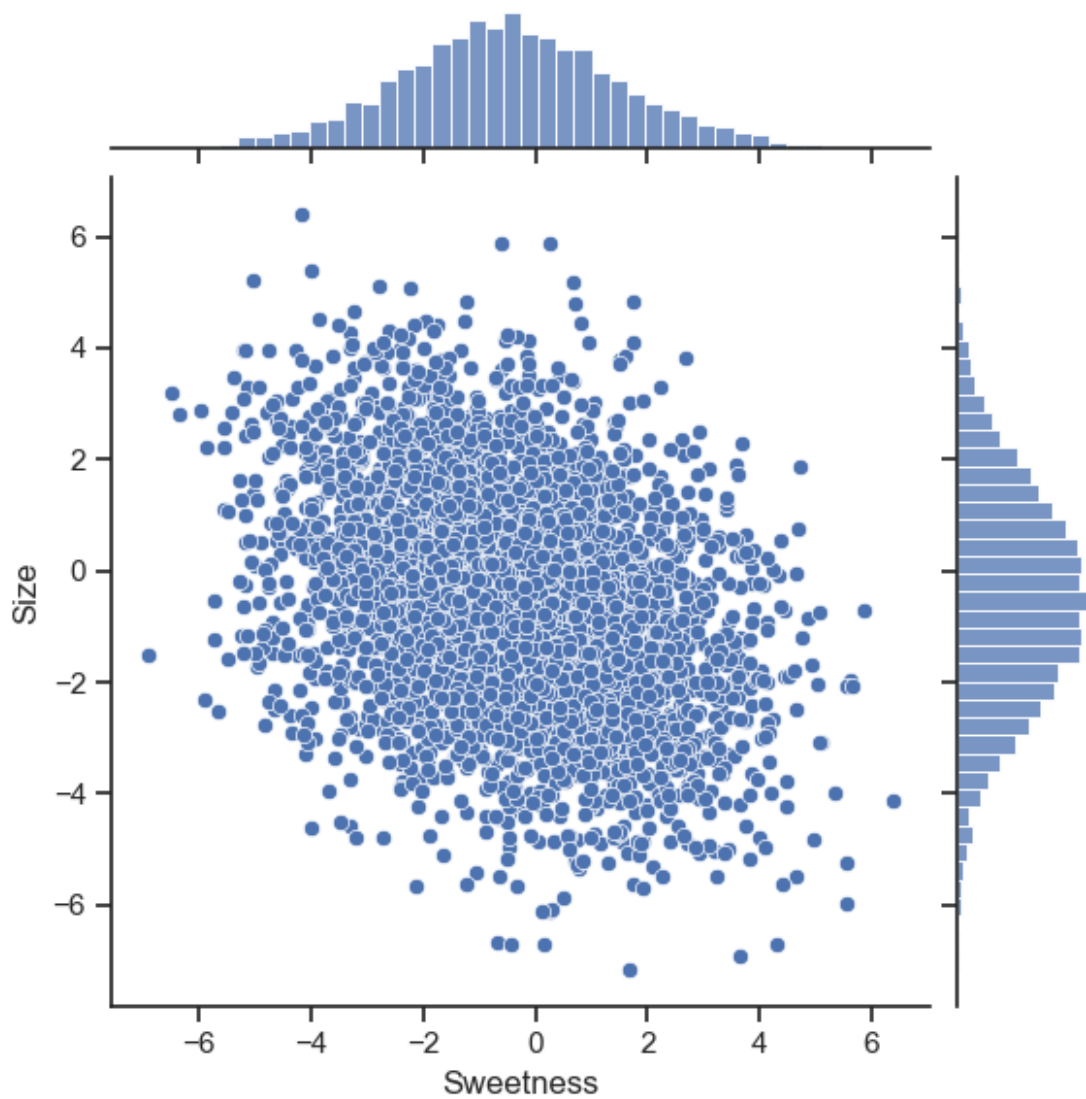
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

```
    sns.distplot(df_data['Sweetness'])
```

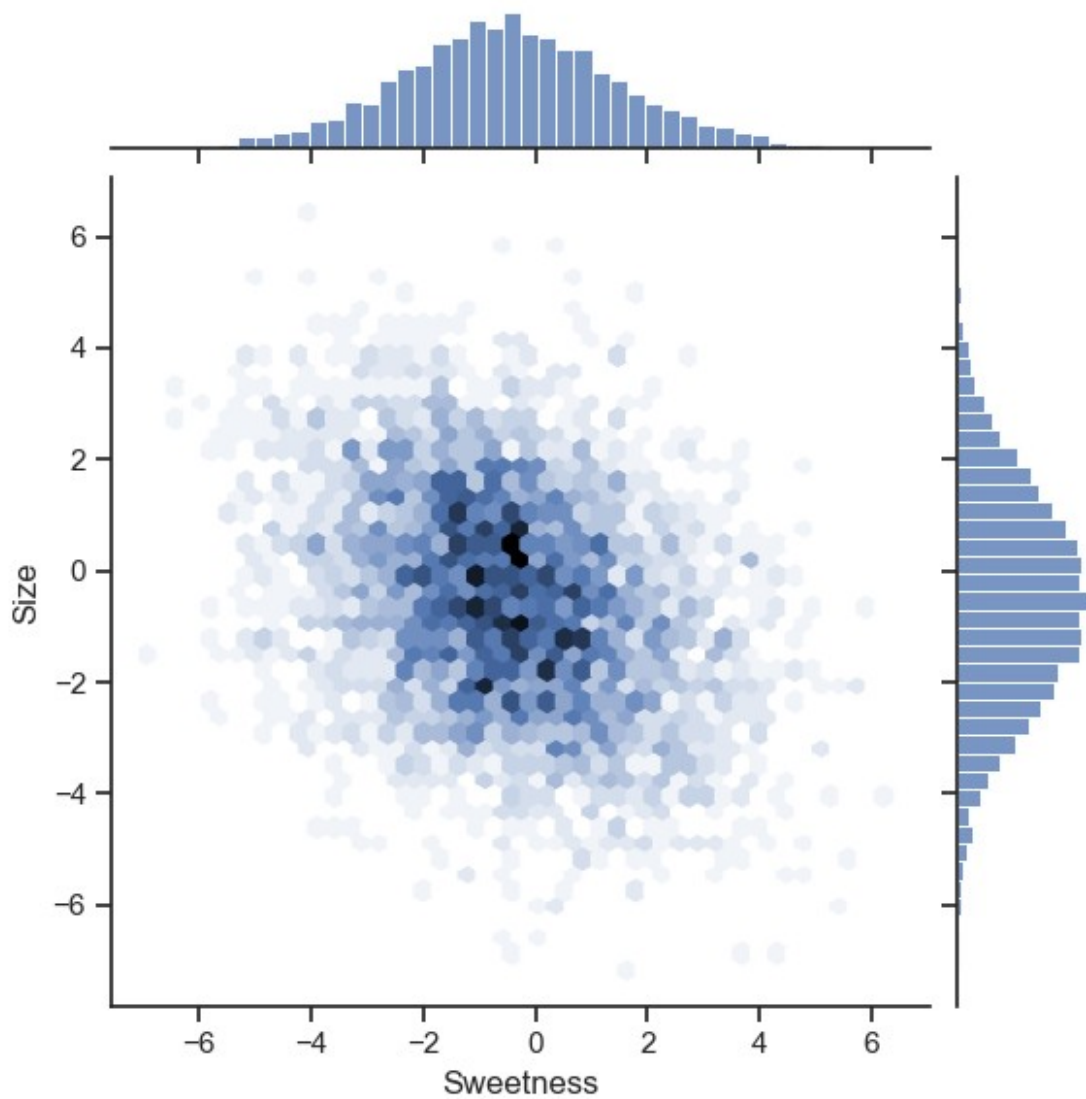<Axes: xlabel='Sweetness', ylabel='Density'>



```
sns.jointplot(x='Sweetness', y='Size', data=df_data)
```
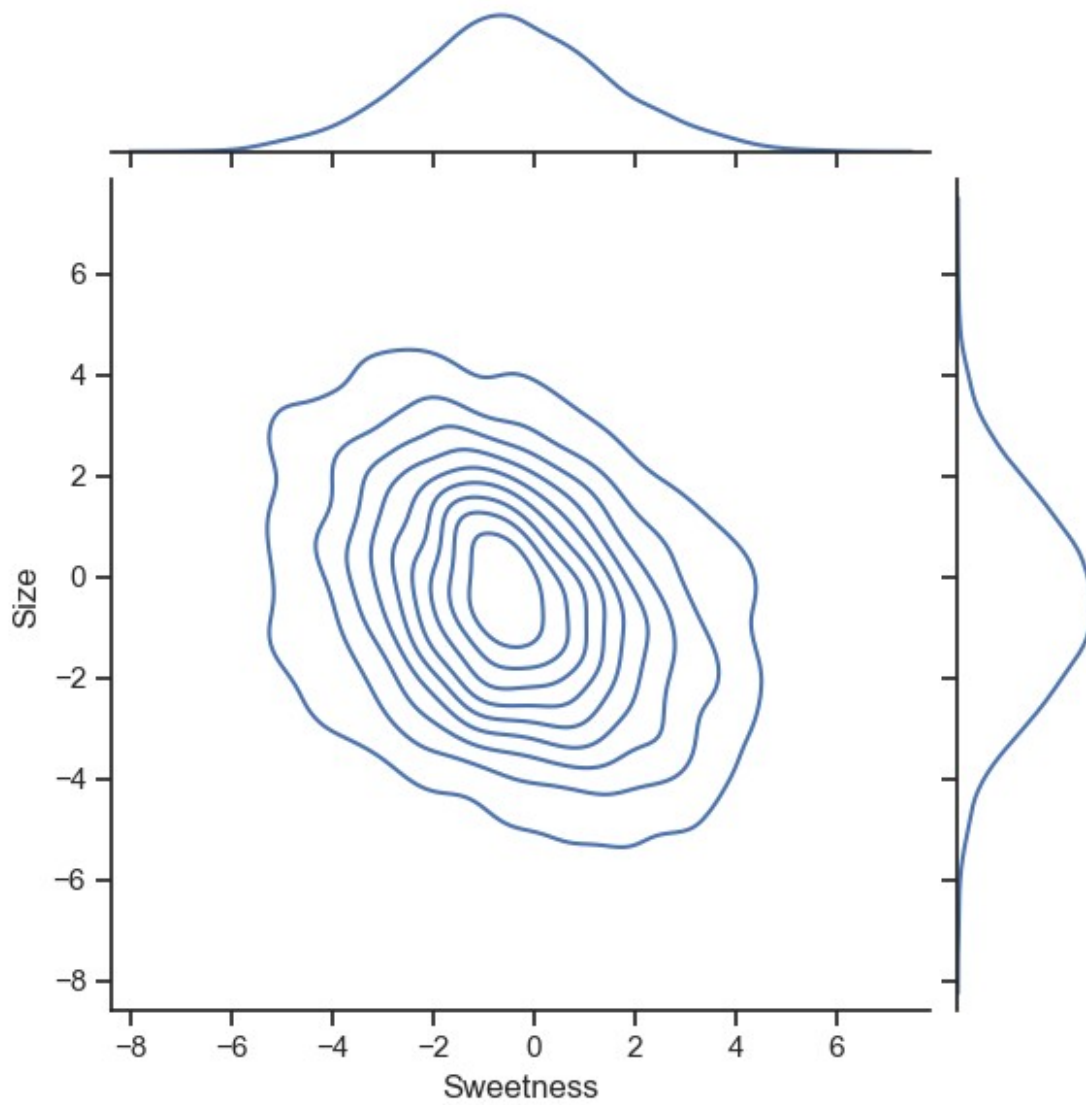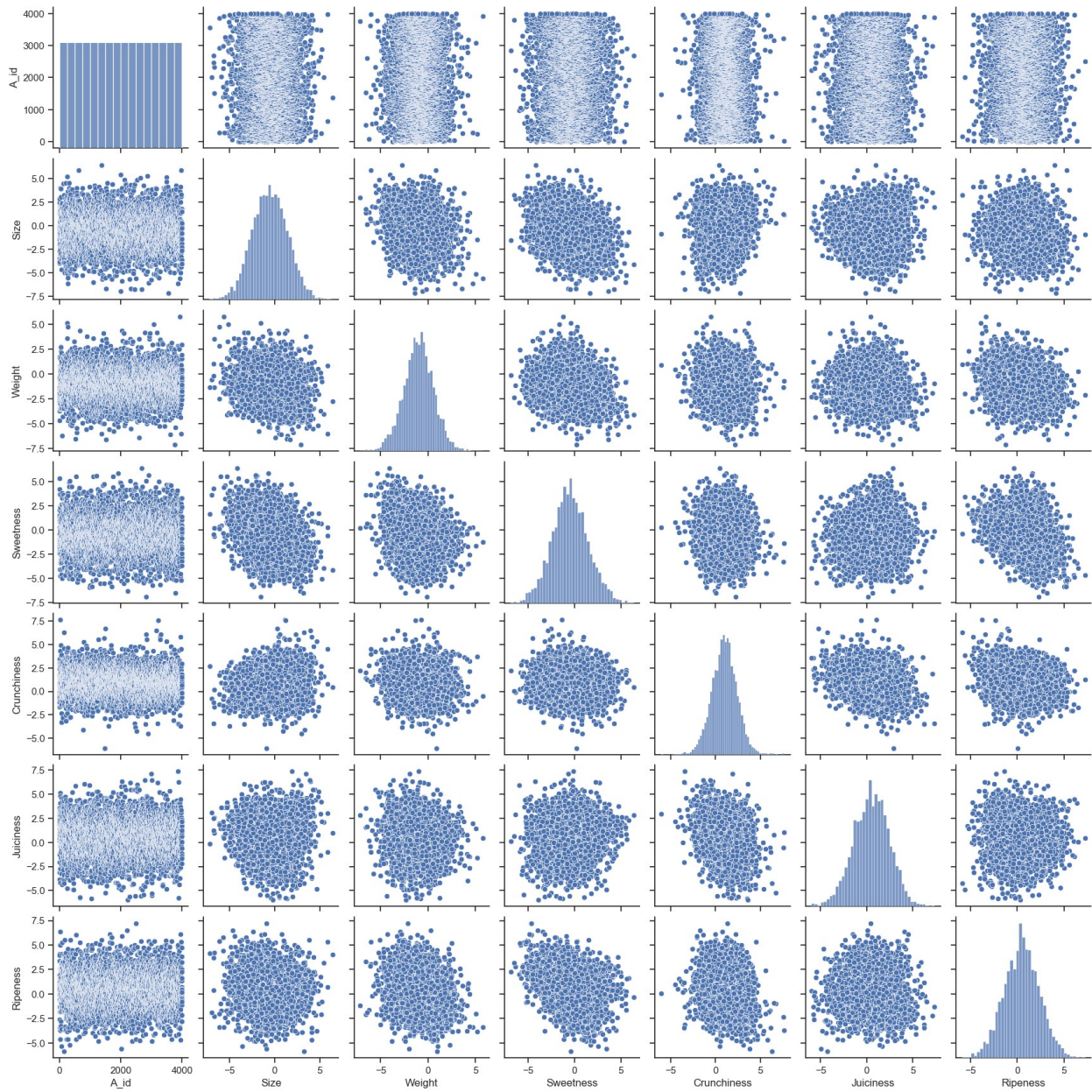
<seaborn.axisgrid.JointGrid at 0x13ca259a0>

```
sns.jointplot(x='Sweetness', y='Size', data=df_data, kind="hex")
<seaborn.axisgrid.JointGrid at 0x13cc3e4f0>
```

```
sns.jointplot(x='Sweetness', y='Size', data=df_data, kind="kde")
<seaborn.axisgrid.JointGrid at 0x13cc3e610>
```
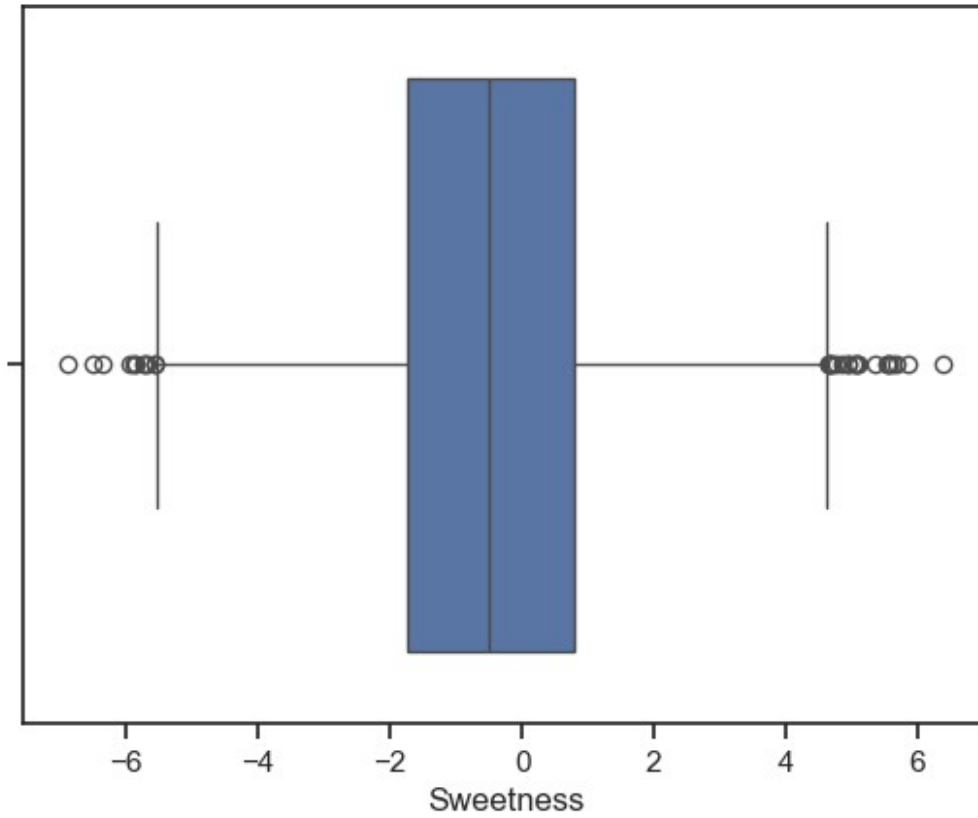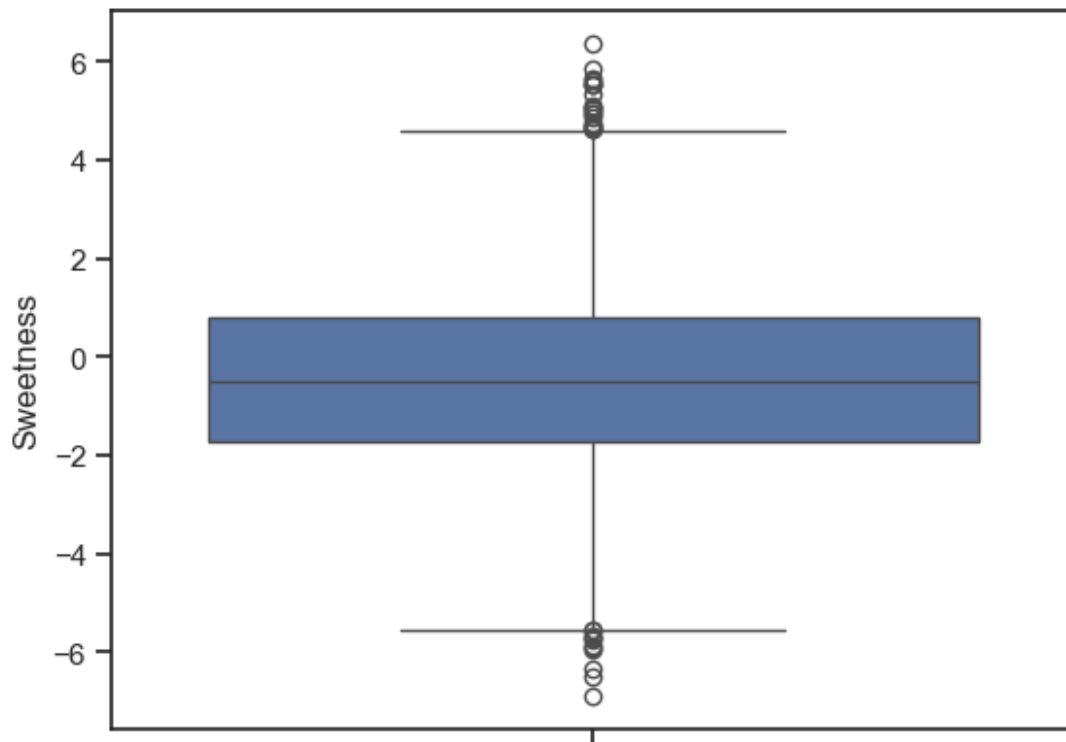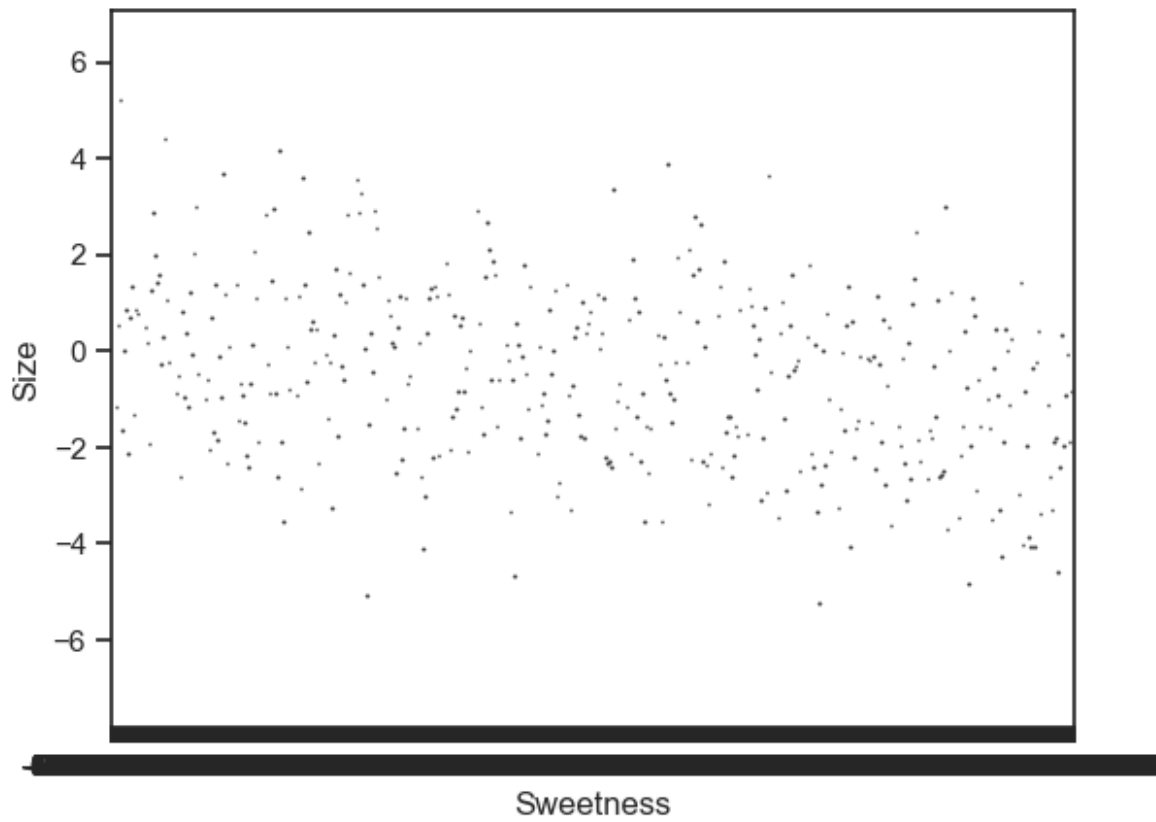
```
sns.pairplot(df_data)
```

```
<seaborn.axisgrid.PairGrid at 0x13cc3e550>
```

```
sns.boxplot(x=df_data['Sweetness'])
```

```
<Axes: xlabel='Sweetness'>
```

```
sns.boxplot(y=df_data['Sweetness'])
```
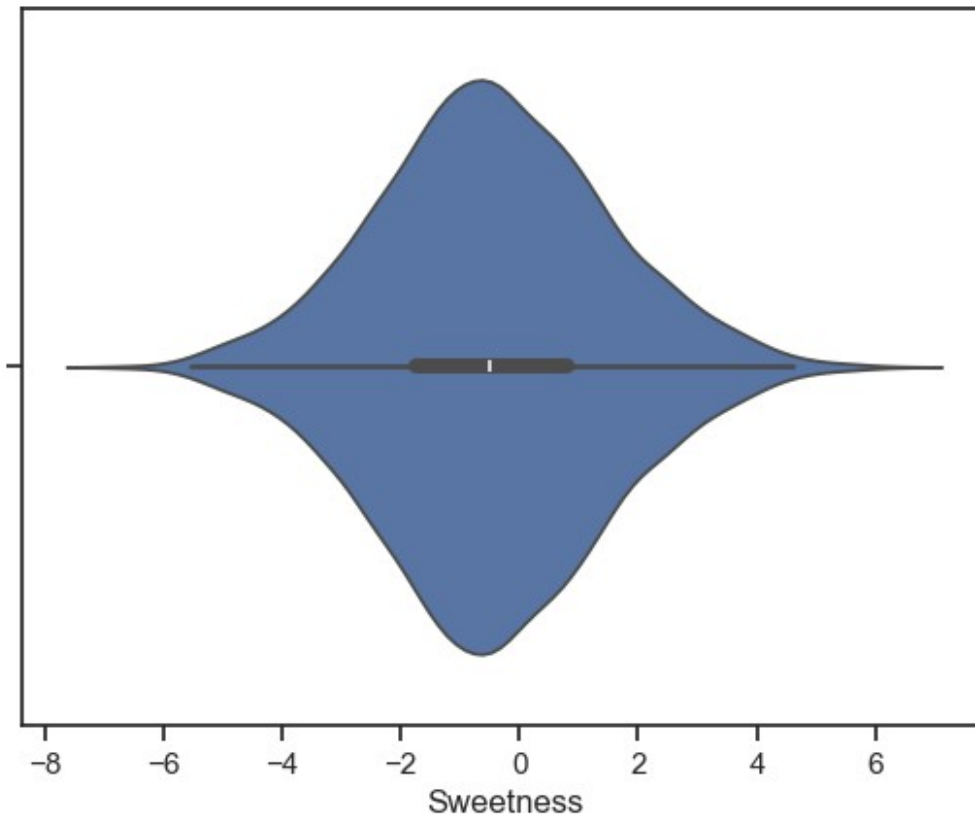
```
<Axes: ylabel='Sweetness'>
```

```
sns.boxplot(x='Sweetness', y='Size', data=df_data)
<Axes: xlabel='Sweetness', ylabel='Size'>
```

```
sns.violinplot(x=df_data['Sweetness'])

<Axes: xlabel='Sweetness'>
```

```
fig, ax = plt.subplots(2, 1, figsize=(10,10))
sns.violinplot(ax=ax[0], x=df_data['Sweetness'])
sns.distplot(df_data['Sweetness'], ax=ax[1])
```

```
/var/folders/bg/1zs8qp8d26v62zscyhsgmff40000gq/T/
ipykernel_2543/2581262117.py:3: UserWarning:

`distplot` is a deprecated function and will be removed in seaborn
v0.14.0.

Please adapt your code to use either `displot` (a figure-level
function with
similar flexibility) or `histplot` (an axes-level function for
histograms).

For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751

  sns.distplot(df_data['Sweetness'], ax=ax[1])
```
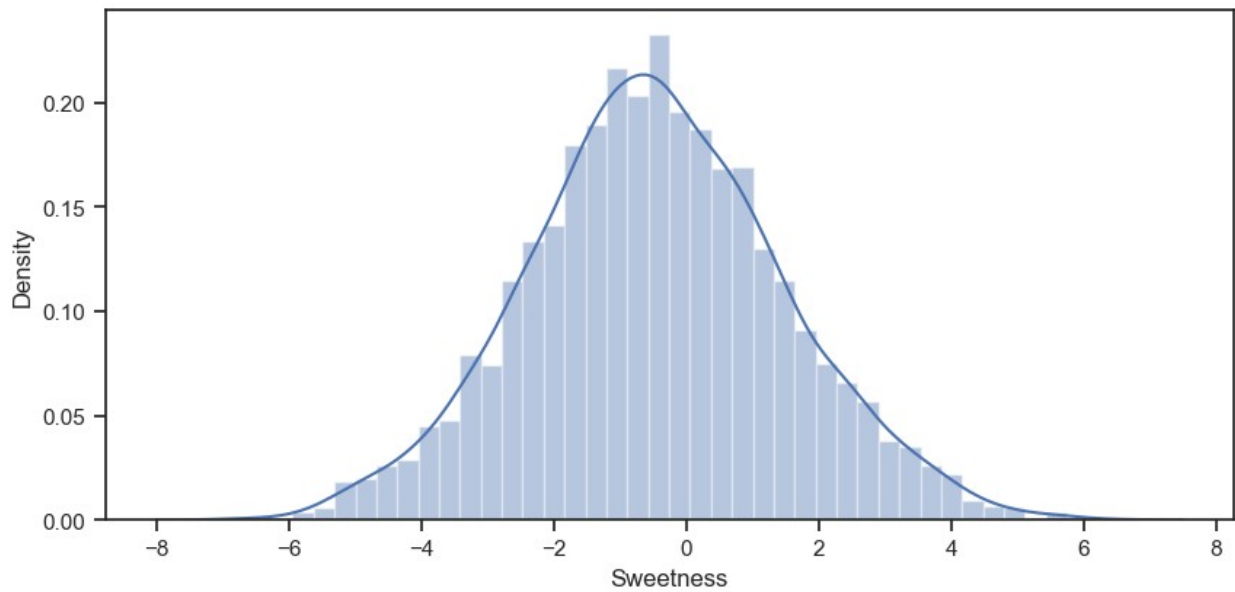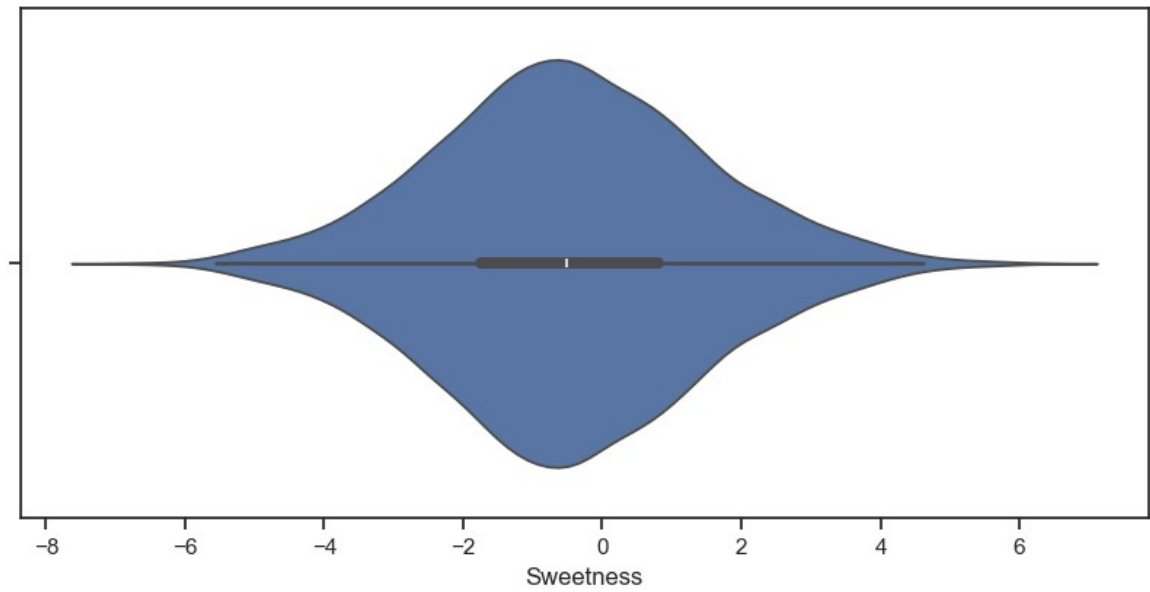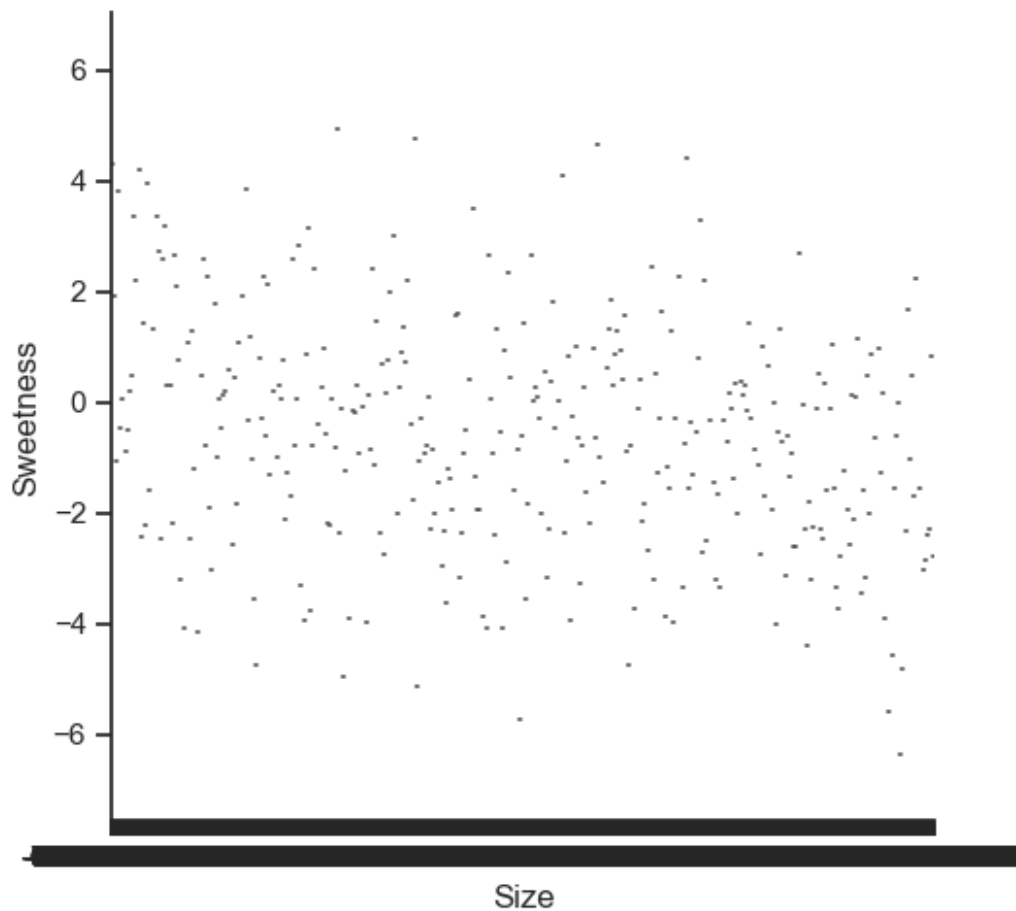
```
<Axes: xlabel='Sweetness', ylabel='Density'>
```

```
sns.catplot(y='Sweetness', x='Size', data=df_data, kind="violin",
split=True)
```

```
<seaborn.axisgrid.FacetGrid at 0x2908ebcd0>
```

```
df_data.drop(columns=["Quality", "Acidity"]).corr()

                A_id       Size     Weight   Sweetness   Crunchiness
Juiciness  \
A_id        1.000000  -0.028911  -0.005730   -0.002378     -0.013111
0.006179
Size       -0.028911   1.000000  -0.170702   -0.324680      0.169868   -
0.018892
Weight     -0.005730  -0.170702   1.000000   -0.154246     -0.095882   -
0.092263
Sweetness  -0.002378  -0.324680  -0.154246    1.000000     -0.037552
0.095882
Crunchiness -0.013111   0.169868  -0.095882   -0.037552      1.000000   -
0.259607
Juiciness   0.006179  -0.018892  -0.092263    0.095882     -0.259607
1.000000
Ripeness    0.000742  -0.134773  -0.243824   -0.273800     -0.201982   -
0.097144


            Ripeness
A_id        0.000742
```

```
Size        -0.134773
Weight      -0.243824
Sweetness   -0.273800
Crunchiness -0.201982
Juiciness   -0.097144
Ripeness     1.000000

df_data.drop(columns=["Quality", "Acidity"]).corr(method='pearson')

                 A_id       Size     Weight   Sweetness   Crunchiness
Juiciness  \
A_id         1.000000 -0.028911 -0.005730   -0.002378     -0.013111
0.006179
Size        -0.028911  1.000000 -0.170702   -0.324680      0.169868   -
0.018892
Weight      -0.005730 -0.170702  1.000000   -0.154246     -0.095882   -
0.092263
Sweetness   -0.002378 -0.324680 -0.154246    1.000000     -0.037552
0.095882
Crunchiness -0.013111  0.169868 -0.095882   -0.037552      1.000000   -
0.259607
Juiciness    0.006179 -0.018892 -0.092263    0.095882     -0.259607
1.000000
Ripeness     0.000742 -0.134773 -0.243824   -0.273800     -0.201982   -
0.097144

            Ripeness
A_id         0.000742
Size        -0.134773
Weight      -0.243824
Sweetness   -0.273800
Crunchiness -0.201982
Juiciness   -0.097144
Ripeness     1.000000

df_data.drop(columns=["Quality", "Acidity"]).corr(method='kendall')

                 A_id       Size     Weight   Sweetness   Crunchiness
Juiciness  \
A_id         1.000000 -0.022124 -0.004756    0.001090     -0.010822
0.002903
Size        -0.022124  1.000000 -0.097221   -0.211004      0.118658   -
0.023001
Weight      -0.004756 -0.097221  1.000000   -0.080836     -0.058782   -
0.060676
Sweetness    0.001090 -0.211004 -0.080836    1.000000     -0.011565
0.065046
Crunchiness -0.010822  0.118658 -0.058782   -0.011565      1.000000   -
0.161359
Juiciness    0.002903 -0.023001 -0.060676    0.065046     -0.161359
```

```
1.000000
Ripeness      -0.003643 -0.101724 -0.166940  -0.171992      -0.125027  -
0.085860

            Ripeness
A_id          -0.003643
Size          -0.101724
Weight        -0.166940
Sweetness     -0.171992
Crunchiness   -0.125027
Juiciness     -0.085860
Ripeness       1.000000
```

```python
sns.heatmap(df_data.drop(columns=["Quality", "Acidity"]).corr())
```

```
<Axes: >
```