

McDonald's case study-Market segmentation

January 29, 2024

```
[2]: #IMPORTING THE IMPORTANT LIBRARIES  
import os  
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
!pip install yellowbrick
```

Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: yellowbrick in ./local/lib/python3.10/site-packages (1.5)
Requirement already satisfied: matplotlib!=3.0.0,>=2.0.2 in /usr/local/lib/python3.10/site-packages (from yellowbrick) (3.6.3)
Requirement already satisfied: scipy>=1.0.0 in /usr/local/lib/python3.10/site-packages (from yellowbrick) (1.9.3)
Requirement already satisfied: scikit-learn>=1.0.0 in /usr/local/lib/python3.10/site-packages (from yellowbrick) (1.3.1)
Requirement already satisfied: numpy>=1.16.0 in /usr/local/lib/python3.10/site-packages (from yellowbrick) (1.23.5)
Requirement already satisfied: cycycler>=0.10.0 in /usr/local/lib/python3.10/site-packages (from yellowbrick) (0.11.0)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (1.0.7)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (4.33.3)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (1.4.3)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (22.0)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (9.1.1)
Requirement already satisfied: pyparsing>=2.2.1 in /usr/local/lib/python3.10/site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (3.0.9)

Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/site-packages (from matplotlib!=3.0.0,>=2.0.2->yellowbrick) (2.8.2)
 Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/site-packages (from scikit-learn>=1.0.0->yellowbrick) (1.2.0)
 Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/site-packages (from scikit-learn>=1.0.0->yellowbrick) (3.1.0)
 Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/site-packages (from python-dateutil>=2.7->matplotlib!=3.0.0,>=2.0.2->yellowbrick) (1.16.0)

[notice] A new release of pip is available: 23.3 -> 23.3.2
 [notice] To update, run:
 pip install --upgrade pip

```
[3]: mcdonalds=pd.read_csv("mcdonalds.csv")
```

```
[4]: mcdonalds.columns
```

```
[4]: Index(['yummy', 'convenient', 'spicy', 'fattening', 'greasy', 'fast', 'cheap',
          'tasty', 'expensive', 'healthy', 'disgusting', 'Like', 'Age',
          'VisitFrequency', 'Gender'],
          dtype='object')
```

```
[5]: mcdonalds.shape
```

```
[5]: (1453, 15)
```

```
[27]: mcdonalds.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1453 entries, 0 to 1452
Data columns (total 15 columns):
#   Column                Non-Null Count  Dtype
---  -
0   yummy                 1453 non-null   object
1   convenient            1453 non-null   object
2   spicy                 1453 non-null   object
3   fattening             1453 non-null   object
4   greasy                1453 non-null   object
5   fast                  1453 non-null   object
6   cheap                 1453 non-null   object
7   tasty                 1453 non-null   object
8   expensive             1453 non-null   object
9   healthy               1453 non-null   object
10  disgusting            1453 non-null   object
```

```

11 Like          1453 non-null  object
12 Age           1453 non-null  int64
13 VisitFrequency 1453 non-null  object
14 Gender        1453 non-null  object
dtypes: int64(1), object(14)
memory usage: 170.4+ KB

```

```
[6]: mcdonalds.describe()
```

```

[6]:
      count  1453.000000
      mean    44.604955
      std     14.221178
      min     18.000000
      25%     33.000000
      50%     45.000000
      75%     57.000000
      max     71.000000

```

```
[28]: print(pd.isnull(mcdonalds).sum())
```

```

yummy          0
convenient      0
spicy           0
fattening       0
greasy          0
fast            0
cheap           0
tasty           0
expensive       0
healthy         0
disgusting      0
Like            0
Age             0
VisitFrequency  0
Gender          0
dtype: int64

```

```
[ ]: # There are no null values.
```

```
[7]: mcdonalds.head(5)
```

```

[7]:  yummy convenient spicy fattening greasy fast cheap tasty expensive healthy \
0     No           Yes   No         Yes   No  Yes   Yes   No         Yes    No
1     Yes          Yes   No         Yes   Yes  Yes   Yes   Yes         Yes    No
2     No           Yes   Yes         Yes   Yes  Yes   No   Yes         Yes    Yes
3     Yes          Yes   No         Yes   Yes  Yes   Yes   Yes         No     No

```

4	No	Yes	No	Yes	Yes	Yes	Yes	No	No	Yes
---	----	-----	----	-----	-----	-----	-----	----	----	-----

	disgusting	Like	Age	VisitFrequency	Gender
0	No	-3	61	Every three months	Female
1	No	+2	51	Every three months	Female
2	No	+1	62	Every three months	Female
3	Yes	+4	69	Once a week	Female
4	No	+2	49	Once a month	Male

```
[ ]: # # From this output we can see that the first respondent is a 61 year old,
      ↪female who eats at McDonald's every 3 months and she believes that is not
      ↪yummy, convenient, not spicy, fattening, not greasy, fast, cheap, not tasty,
      ↪expensive,not healthy and not disgusting
```

```
[8]: mcdonalds.dtypes
```

```
[8]: yummy          object
      convenient     object
      spicy          object
      fattening      object
      greasy         object
      fast           object
      cheap          object
      tasty          object
      expensive      object
      healthy        object
      disgusting     object
      Like           object
      Age            int64
      VisitFrequency object
      Gender         object
      dtype: object
```

```
[25]: mcdonalds['yummy'].value_counts()
```

```
[25]: Yes      803
      No       650
      Name: yummy, dtype: int64
```

```
[26]: mcdonalds['Age'].value_counts()
```

```
[26]: 55      53
      60      38
      37      37
      59      36
      57      36
      52      36
```

58	35
36	35
49	34
62	34
50	34
32	33
44	32
56	32
64	32
53	31
26	31
24	30
35	30
51	30
47	30
42	30
23	30
39	29
29	28
34	28
30	28
38	27
40	27
31	27
25	26
33	26
61	26
67	26
48	26
43	25
27	25
63	25
54	24
41	23
22	23
65	23
45	22
20	21
46	19
28	18
66	17
21	16
18	16
70	15
69	14
68	13
19	10

```
71      1
Name: Age, dtype: int64
```

```
[29]: mcdonalds.corr()
```

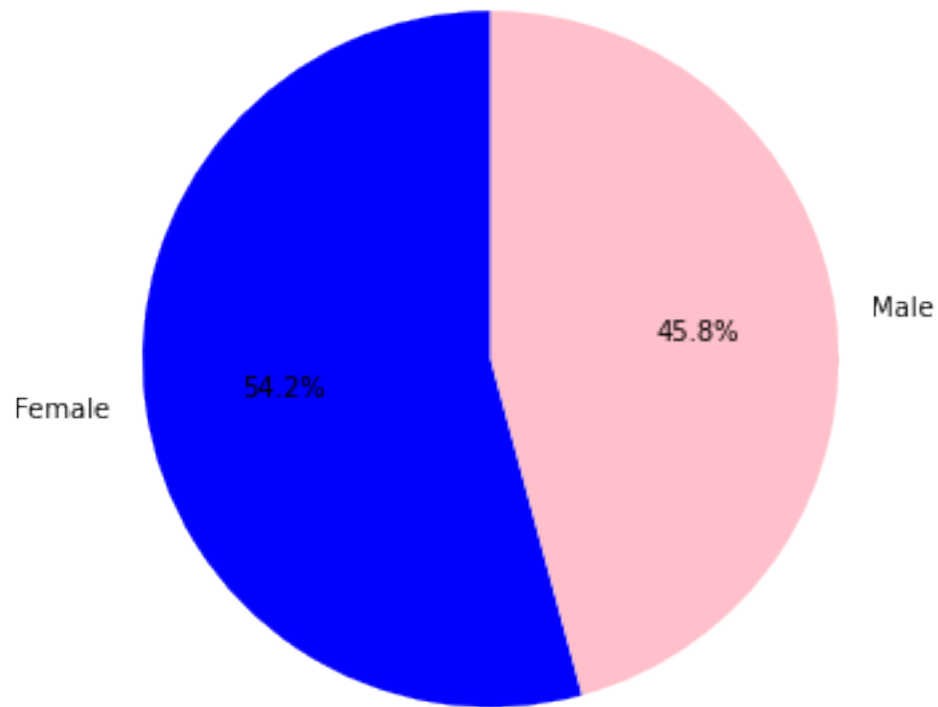
```
/tmp/ipykernel_758/2370541813.py:1: FutureWarning: The default value of
numeric_only in DataFrame.corr is deprecated. In a future version, it will
default to False. Select only valid columns or specify the value of numeric_only
to silence this warning.
    mcdonalds.corr()
```

```
[29]:      Age
Age  1.0
```

```
[30]: # Basic market segmentation with gender as a Variable.
      # Count the number of instances for each gender
      gender_counts = mcdonalds['Gender'].value_counts()

      # Plotting the pie chart
      plt.figure(figsize=(6, 6))
      plt.pie(gender_counts, labels=gender_counts.index, autopct='%1.1f%%',
      ↪startangle=90, colors=['blue', 'pink'])
      plt.title('Distribution of Gender in McDonald\'s Dataset')
      plt.show()
```

Distribution of Gender in McDonald's Dataset



```
[ ]: # The number of female costumers is more than that of male customers.
```

```
[10]: # Extracting the columns containing segmentation variables(1-11) and converting  
      ↪ the data to a matrix containing 0 and 1 instead of YES and NO.
```

```
[11]: # Extracting columns 1 to 11 and converting to a binary matrix  
MD_x = mcdonalds.iloc[:, 0:11]  
# Converting values to 0 and 1  
MD_x = (MD_x == "Yes").astype(int)  
  
# Calculating and rounding column means to two decimal places  
print(np.round(MD_x.mean(), 2))
```

yummy	0.55
convenient	0.91
spicy	0.09
fattening	0.87
greasy	0.53

```
fast          0.90
cheap         0.60
tasty         0.64
expensive     0.36
healthy       0.20
disgusting    0.24
dtype: float64
```

```
[12]: # The output indicates that about 55% of the respondents believe that
      ↪ McDonald's food is YUMMY, 53 believe that the food at McDonald's is GREASY
      ↪ and 36% believes that it is expensive.
```

```
[13]: # PRINCIPAL COMPONENT ANALYSIS
```

```
[14]: from sklearn.decomposition import PCA
      pca = PCA()
      MD_pca = pca.fit_transform(MD_x)

      # Extracting relevant information
      result_summary = pd.DataFrame({
          'PC': [f'PC{i+1}' for i in range(pca.n_components_)], 'Standard deviation':
          ↪ pca.explained_variance_**0.5, 'Proportion of Variance': pca.
          ↪ explained_variance_ratio_, 'Cumulative Proportion': pca.
          ↪ explained_variance_ratio_.cumsum()
      })

      # Set the precision for better formatting
      pd.set_option("display.precision", 4)

      # Display the result summary
      print(result_summary)
```

	PC	Standard deviation	Proportion of Variance	Cumulative Proportion
0	PC1	0.7570	0.2994	0.2994
1	PC2	0.6075	0.1928	0.4922
2	PC3	0.5046	0.1330	0.6253
3	PC4	0.3988	0.0831	0.7084
4	PC5	0.3374	0.0595	0.7679
5	PC6	0.3103	0.0503	0.8182
6	PC7	0.2897	0.0438	0.8620
7	PC8	0.2751	0.0395	0.9016
8	PC9	0.2653	0.0368	0.9383
9	PC10	0.2488	0.0324	0.9707
10	PC11	0.2369	0.0293	1.0000

```
[15]: standard_deviations = pca.explained_variance_**0.5
      print("Standard deviations (1, ..., p=11):")
```



```
print(np.round(standard_deviations, 1))
```

```
Standard deviations (1, .., p=11):  
[0.8 0.6 0.5 0.4 0.3 0.3 0.3 0.3 0.3 0.2 0.2]
```

```
[16]: # printing the rotation matrix:  
column_names = MD_x.columns  
rotation_matrix = pd.DataFrame(np.round(pca.components_.T, 2),  
                                columns=column_names)  
print(f"Rotation (n x k) = ({MD_x.shape[1]} x {MD_x.shape[1]}):")  
print(rotation_matrix)
```

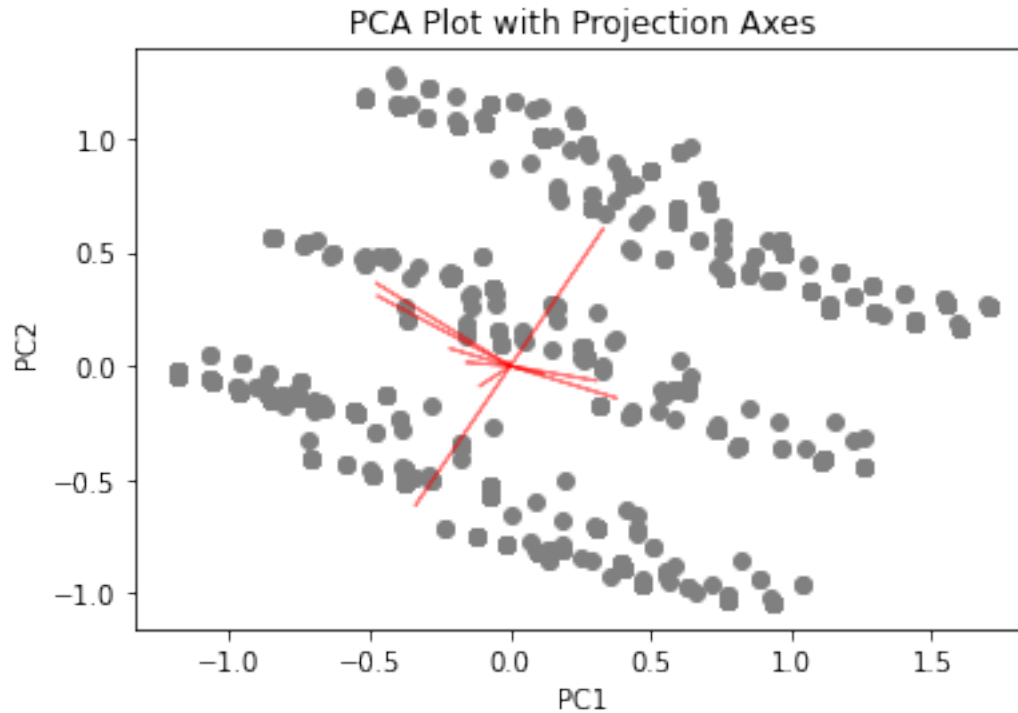
```
Rotation (n x k) = (11 x 11):
```

	yummy	convenient	spicy	fattening	greasy	fast	cheap	tasty	\
0	-0.48	0.36	-0.30	0.06	-0.31	0.17	-0.28	0.01	
1	-0.16	0.02	-0.06	-0.14	0.28	-0.35	-0.06	-0.11	
2	-0.01	0.02	-0.04	0.20	0.07	-0.36	0.71	0.38	
3	0.12	-0.03	-0.32	-0.35	-0.07	-0.41	-0.39	0.59	
4	0.30	-0.06	-0.80	0.25	0.36	0.21	0.04	-0.14	
5	-0.11	-0.09	-0.06	-0.10	0.11	-0.59	-0.09	-0.63	
6	-0.34	-0.61	-0.15	0.12	-0.13	-0.10	-0.04	0.14	
7	-0.47	0.31	-0.29	-0.00	-0.21	-0.08	0.36	-0.07	
8	0.33	0.60	0.02	0.07	-0.00	-0.26	-0.07	0.03	
9	-0.21	0.08	0.19	0.76	0.29	-0.18	-0.35	0.18	
10	0.37	-0.14	-0.09	0.37	-0.73	-0.21	-0.03	-0.17	

	expensive	healthy	disgusting
0	0.57	-0.11	0.05
1	-0.02	-0.67	-0.54
2	0.40	-0.08	0.14
3	-0.16	-0.01	0.25
4	-0.00	0.01	0.00
5	0.17	0.24	0.34
6	0.08	0.43	-0.49
7	-0.64	0.08	0.02
8	0.07	0.45	-0.49
9	-0.19	-0.04	0.16
10	-0.07	-0.29	-0.04

```
[17]: # Plotting the data points in the reduced dimensional space  
plt.scatter(MD_pca[:, 0], MD_pca[:, 1], color='grey')  
  
# Adding projection axes  
for i in range(MD_pca.shape[1]):  
    plt.arrow(0, 0, pca.components_[0, i], pca.components_[1, i], color='red',  
              alpha=0.5)
```

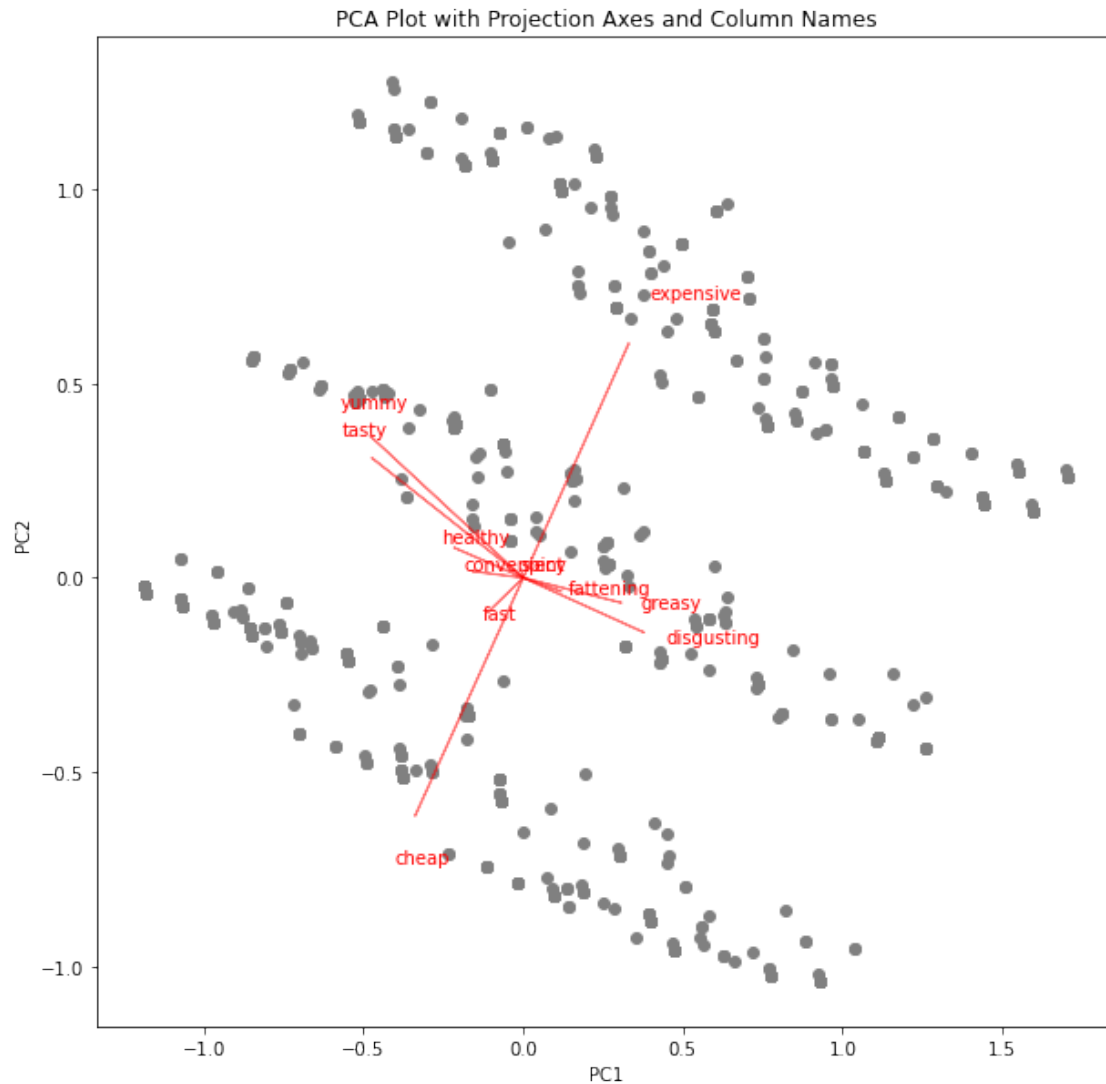
```
plt.xlabel('PC1')
plt.ylabel('PC2')
plt.title('PCA Plot with Projection Axes')
plt.show()
```



```
[18]: plt.figure(figsize=(10, 10))
# Plotting the data points in the reduced dimensional space
plt.scatter(MD_pca[:, 0], MD_pca[:, 1], color='grey')

# Adding projection axes and annotating with column names
for i in range(MD_pca.shape[1]):
    plt.arrow(0, 0, pca.components_[0, i], pca.components_[1, i], color='red',
    ↪alpha=0.5)
    plt.text(pca.components_[0, i] * 1.2, pca.components_[1, i] * 1.2, MD_x.
    ↪columns[i], color='red')

plt.xlabel('PC1')
plt.ylabel('PC2')
plt.title('PCA Plot with Projection Axes and Column Names')
plt.show()
```



```
[19]: pip install scikit-learn-extra
```

```
Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: scikit-learn-extra in
./local/lib/python3.10/site-packages (0.3.0)
Requirement already satisfied: numpy>=1.13.3 in /usr/local/lib/python3.10/site-
packages (from scikit-learn-extra) (1.23.5)
Requirement already satisfied: scipy>=0.19.1 in /usr/local/lib/python3.10/site-
packages (from scikit-learn-extra) (1.9.3)
Requirement already satisfied: scikit-learn>=0.23.0 in
/usr/local/lib/python3.10/site-packages (from scikit-learn-extra) (1.3.1)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/site-
packages (from scikit-learn>=0.23.0->scikit-learn-extra) (1.2.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in
```

/usr/local/lib/python3.10/site-packages (from scikit-learn>=0.23.0->scikit-learn-extra) (3.1.0)

[notice] A new release of pip is available: 23.3 -> 23.3.2

[notice] To update, run:

`pip install --upgrade pip`

Note: you may need to restart the kernel to use updated packages.

```
[20]: # Extracting segments using k-means
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import make_pipeline
from sklearn_extra.cluster import KMedoids

# Set a seed for reproducibility
np.random.seed(1234)
cluster_range = [2, 8]

# Performing k-means clustering with ten random restarts
results = {}
for num_clusters in cluster_range:
    kmeans = KMeans(n_clusters=num_clusters, n_init=10, random_state=1234)
    pipeline = make_pipeline(StandardScaler(), kmeans)
    labels = pipeline.fit_predict(MD_x)
    results[num_clusters] = labels

# Relabel segments to be consistent across segmentations
for num_clusters in cluster_range:
    unique_labels = np.unique(results[num_clusters])
    relabeling_dict = {old_label: new_label for new_label, old_label in
    ↪ enumerate(unique_labels)}
    results[num_clusters] = np.vectorize(relabeling_dict.
    ↪ get)(results[num_clusters])
```

```
[21]: # Access the results for 2 clusters
labels_2_clusters = results[2]
print("Labels for 2 clusters:", labels_2_clusters)

# Access the results for 8 clusters
labels_8_clusters = results[8]
print("Labels for 8 clusters:", labels_8_clusters)
```

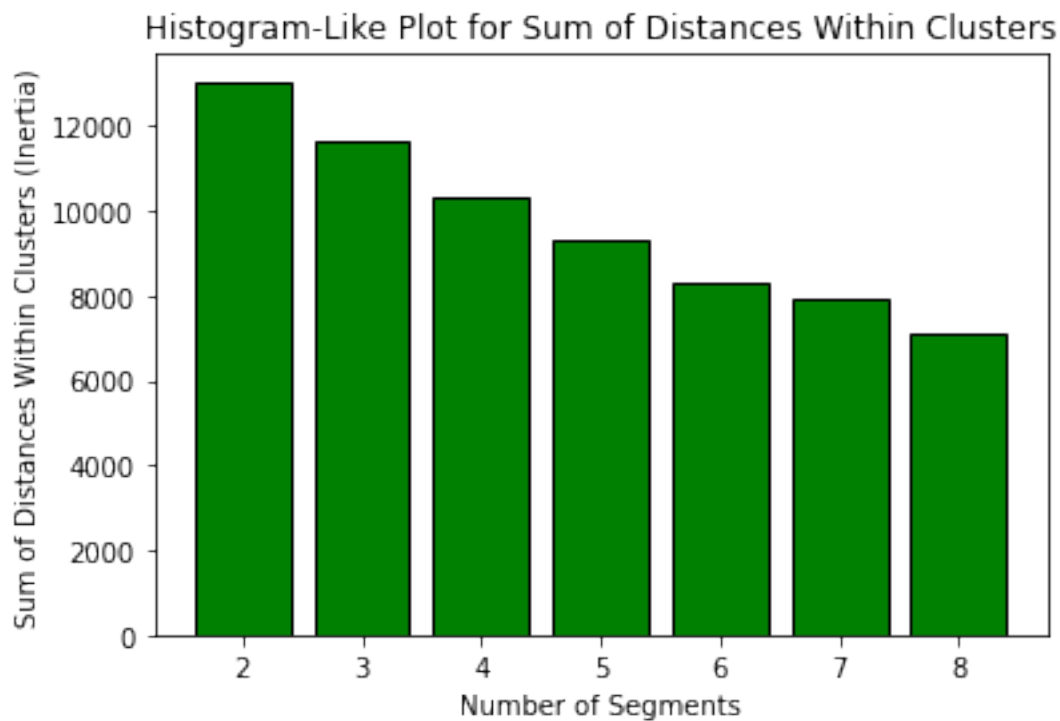
Labels for 2 clusters: [1 0 0 ... 0 0 1]

Labels for 8 clusters: [7 0 6 ... 3 5 2]

```
[22]: # Specify the range of cluster numbers
cluster_range = range(2, 9)

# Compute the inertia (sum of distances within clusters) for each number of
clusters
inertia_values = []
for num_clusters in cluster_range:
    kmeans = KMeans(n_clusters=num_clusters, n_init=10, random_state=1234)
    pipeline = make_pipeline(StandardScaler(), kmeans)
    pipeline.fit(MD_x)
    inertia_values.append(pipeline.named_steps['kmeans'].inertia_)

# Plotting the histogram-like plot
plt.bar(cluster_range, inertia_values, color='green', edgecolor='black')
plt.xlabel('Number of Segments')
plt.ylabel('Sum of Distances Within Clusters (Inertia)')
plt.title('Histogram-Like Plot for Sum of Distances Within Clusters')
plt.show()
```



```
[37]: #Extracting segments-k means

from yellowbrick.cluster import KElbowVisualizer
k_means_model = KMeans()
```

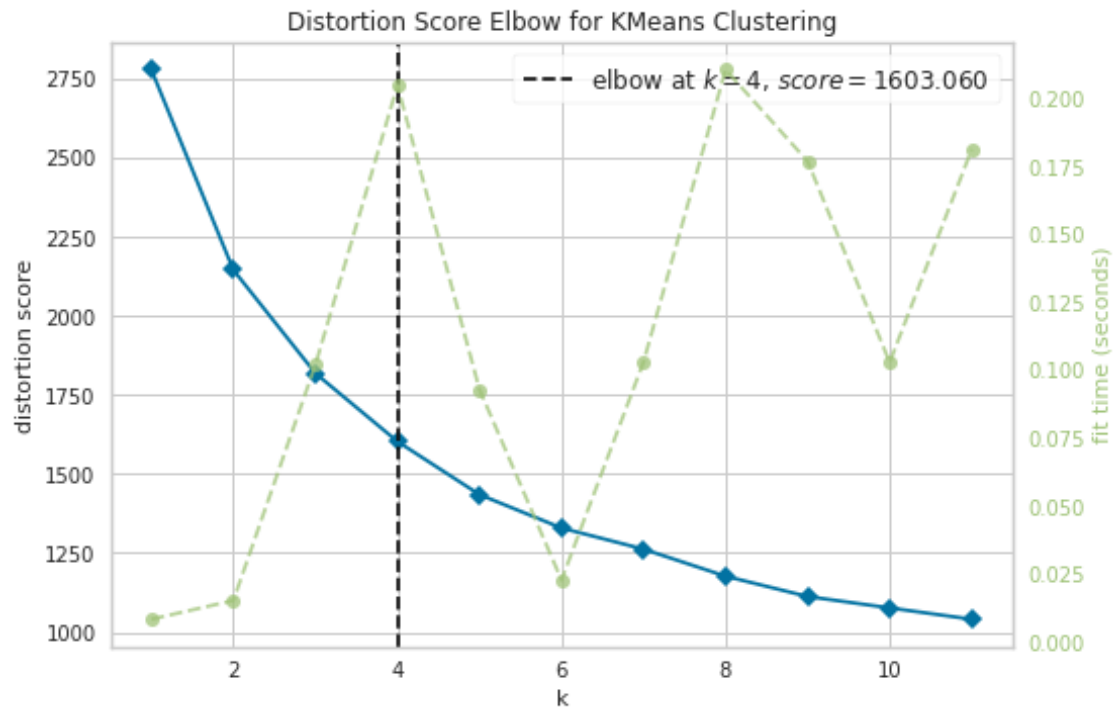
```
visualizer = KElbowVisualizer(k_means_model, k=(1,12)).fit(MD_x)
visualizer.show()
```

```
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
/usr/local/lib/python3.10/site-packages/sklearn/cluster/_kmeans.py:1416:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
    super()._check_params_vs_input(X, default_n_init=10)
```

[illegible]

[illegible]

families were found: Arial, Liberation Sans, Bitstream Vera Sans, sans-serif
 findfont: Generic family 'sans-serif' not found because none of the following
 families were found: Arial, Liberation Sans, Bitstream Vera Sans, sans-serif
 findfont: Generic family 'sans-serif' not found because none of the following
 families were found: Arial, Liberation Sans, Bitstream Vera Sans, sans-serif



[37]: <AxesSubplot: title={'center': 'Distortion Score Elbow for KMeans Clustering'},
 xlabel='k', ylabel='distortion score'>

[]:

[]:

[]:

[]: