

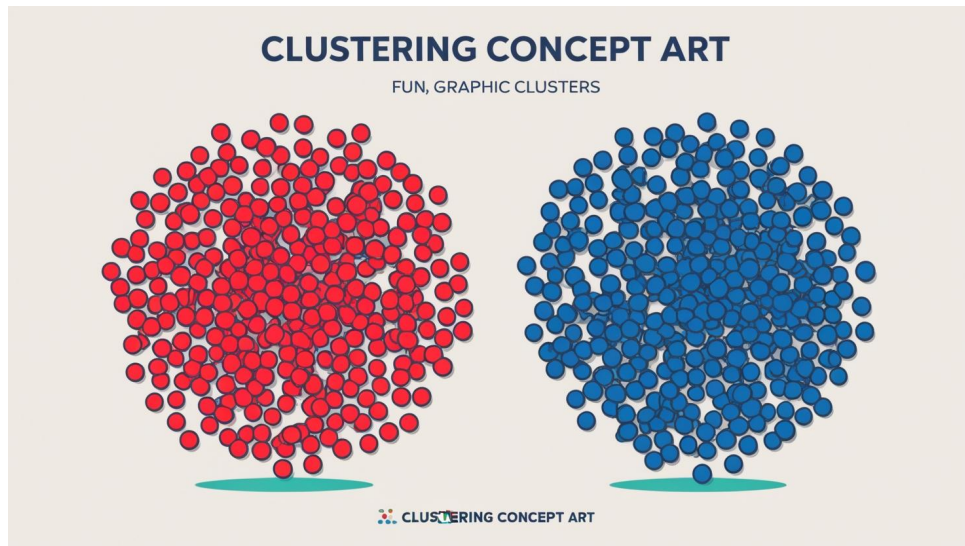
# Segmenting for Success: A Deep Dive into K-Means Clustering

Minor in AI, IIT Ropar

9th May, 2025

## Chapter 2: Beyond the Spreadsheet: K-Means in Action

Welcome back, aspiring data wranglers! In the previous chapter, we laid the groundwork for understanding clustering. We explored how to group similar data points together, focusing on intuition and manual calculations. Now, it's time to level up! We'll trade our spreadsheets for code and unleash the power of the K-Means algorithm.



**Figure 1:** [Placeholder for an illustrative image about data clustering or "leveling up" in data science.]

### 2.1. The Shopkeeper's Inventory: A Recap

Remember our friendly neighborhood shopkeeper from Chapter 1? He's still trying to make sense of his customer data. He wants to understand his clientele better so he can stock the right items and tailor his marketing efforts.

Let's quickly revisit the data:

**Table 1:** Shopkeeper's Customer Data Recap

Customer	Frequency (Visits)	Spend (Amount)
C1	2	500
C2	10	800
C3	4	300
C4	11	1200
C5	3	350
C6	9	1000

We manually grouped these customers based on their frequency of visits and spending habits. But, as you can imagine, manually processing large datasets is a tedious and error-prone task. This is where the K-Means algorithm comes to the rescue!

## 2.2. Translating Intuition into Code

Before diving into libraries, let's peek under the hood and see what a raw K-Means implementation would look like. This helps us appreciate the power and convenience that libraries offer.

```
1 import math
2
3 # Data
4 customers = {
5     "C1": (2, 500),
6     "C2": (10, 800),
7     "C3": (4, 300),
8     "C4": (11, 1200),
9     "C5": (3, 350),
10    "C6": (9, 1000)
11 }
12
13 # Initialize centroids manually
14 centroid_A = customers["C1"]
15 centroid_B = customers["C4"]
16
17 def euclidean(p1, p2):
18     return math.sqrt((p1[0] - p2[0]) ** 2 + (p1[1] - p2[1]) ** 2)
19
20 def compute_new_centroid(cluster):
21     freq = sum(c[0] for c in cluster) / len(cluster)
22     spend = sum(c[1] for c in cluster) / len(cluster)
23     return (freq, spend)
24
25 # Run until convergence
26 iteration = 1
27 prev_assignment = {}
28 while True:
29     print(f"\n--- Iteration {iteration} ---")
30
31     cluster_A = []
32     cluster_B = []
33     assignment = {}
34
35     for cid, (f, s) in customers.items():
36         dA = euclidean((f, s), centroid_A)
37         dB = euclidean((f, s), centroid_B)
38
39         if dA < dB:
40             cluster_A.append((f, s))
41             assignment[cid] = 'A'
42         else:
43             cluster_B.append((f, s))
44             assignment[cid] = 'B'
45
46         print(f"{cid}: Dist to A={dA:.2f}, B={dB:.2f} => Cluster {assignment[cid]}")
47
48     if assignment == prev_assignment:
49         print("\nConverged.")
50         break
51
52     prev_assignment = assignment
53     centroid_A = compute_new_centroid(cluster_A)
54     centroid_B = compute_new_centroid(cluster_B)
55
56     print(f"\nNew Centroid A: {centroid_A}")
57     print(f"New Centroid B: {centroid_B}")
58
59     iteration += 1
60
61 # Final Output
```

```

62 print("\nFinal Clusters:")
63 print("Cluster A:", [cid for cid, grp in assignment.items() if grp == 'A'])
64 print("Cluster B:", [cid for cid, grp in assignment.items() if grp == 'B'])

```

**Listing 1:** Raw K-Means Implementation from Scratch

This code does the same thing we did manually in the spreadsheet:

1. **Data:** Defines the customer data as a dictionary.
2. **euclidean function:** Calculates the Euclidean distance between two points.
3. **compute\_new\_centroid function:** Calculates the average frequency and spending for a given cluster.
4. **The while loop:** This is the heart of the K-Means algorithm. It iterates until the clusters stop changing (convergence).
  - It assigns each customer to the nearest centroid.
  - It recalculates the centroids based on the new cluster assignments.

This code, while functional, is a bit lengthy. Imagine writing this for hundreds or thousands of data points! Thankfully, libraries like `scikit-learn` offer a much cleaner and efficient way to perform K-Means clustering.

### 2.3. K-Means with scikit-learn: A Breath of Fresh Air

Let's see how we can achieve the same results using the `scikit-learn` library:

```

1 import numpy as np
2 from sklearn.cluster import KMeans
3
4 # Data
5 customer_ids = ["C1", "C2", "C3", "C4", "C5", "C6"]
6 data = np.array([
7     [2, 500],
8     [10, 800],
9     [4, 300],
10    [11, 1200],
11    [3, 350],
12    [9, 1000]
13 ])
14
15 # Initial centroids from C1 and C4
16 initial_centroids = np.array([
17     [2, 500], # C1
18     [11, 1200] # C4
19 ])
20
21 # KMeans model with manually specified initial centroids
22 kmeans = KMeans(n_clusters = 2, init = initial_centroids, n_init = 1, max_iter = 25)
23 kmeans.fit(data)
24
25 # Results
26 labels = kmeans.labels_
27
28 # Display results
29 for cid, label in zip(customer_ids, labels):
30     print(f"{cid} - Cluster {label}")
31
32 print("\nCluster Centers:")
33 print(kmeans.cluster_centers_)

```

**Listing 2:** K-Means with scikit-learn (Manual Centroids)

Let's break down the key parts:

- **import KMeans from sklearn.cluster:** Imports the `KMeans` class.
- **data = np.array(...):** Converts our data into a NumPy array, which is the preferred format for `scikit-learn`.

- `initial_centroids = np.array(...)`: Sets the initial centroid values.
- `kmeans = KMeans(n_clusters=2, init = initial_centroids, n_init = 1, max_iter = 25)`: Creates a `KMeans` object.
  - `n_clusters`: Specifies the number of clusters we want (in this case, 2). This is the "K" in K-Means.
  - `init`: Defines the method for initialization of the centroids. Here we are initializing the centroids manually by passing the `initial_centroids` array.
  - `n_init`: Sets the number of times the K-means algorithm will be run with different centroid seeds. If `init` is 'k-means++' or 'random', the argument is used to declare the number of runs.
  - `max_iter`: Sets the maximum number of iterations for the algorithm to converge.
- `kmeans.fit(data)`: Trains the K-Means model on our data. This is where the magic happens! The algorithm finds the best cluster assignments based on the data and the specified parameters.
- `labels = kmeans.labels_`: Retrieves the cluster labels for each data point. `labels_` is an attribute of the trained `kmeans` object. It's a NumPy array where each element corresponds to a data point and indicates which cluster that data point belongs to.
- `kmeans.cluster_centers_`: Retrieves the coordinates of the final cluster centroids.

See how much simpler this is compared to the manual implementation? `scikit-learn` handles all the complex calculations and iterations for us.

## 2.4. Random Centroid Initialization: Letting the Algorithm Choose

In the previous example, we manually chose the initial centroids (C1 and C4). But what if we want the algorithm to pick them randomly?

```

1 import numpy as np
2 from sklearn.cluster import KMeans
3
4 # Data
5 customer_ids = ["C1", "C2", "C3", "C4", "C5", "C6"]
6 data = np.array([
7     [2, 500],
8     [10, 800],
9     [4, 300],
10    [11, 1200],
11    [3, 350],
12    [9, 1000]
13 ])
14
15 # KMeans model without manually specified initial centroids
16 kmeans = KMeans(n_clusters = 2, n_init = 10, max_iter = 25) # n_init='auto' is default in
17 # newer versions
18 kmeans.fit(data)
19
20 # Results
21 labels = kmeans.labels_
22
23 # Display results
24 for cid, label in zip(customer_ids, labels):
25     print(f"{cid} - Cluster {label}")
26
27 print("\nCluster Centers:")
28 print(kmeans.cluster_centers_)

```

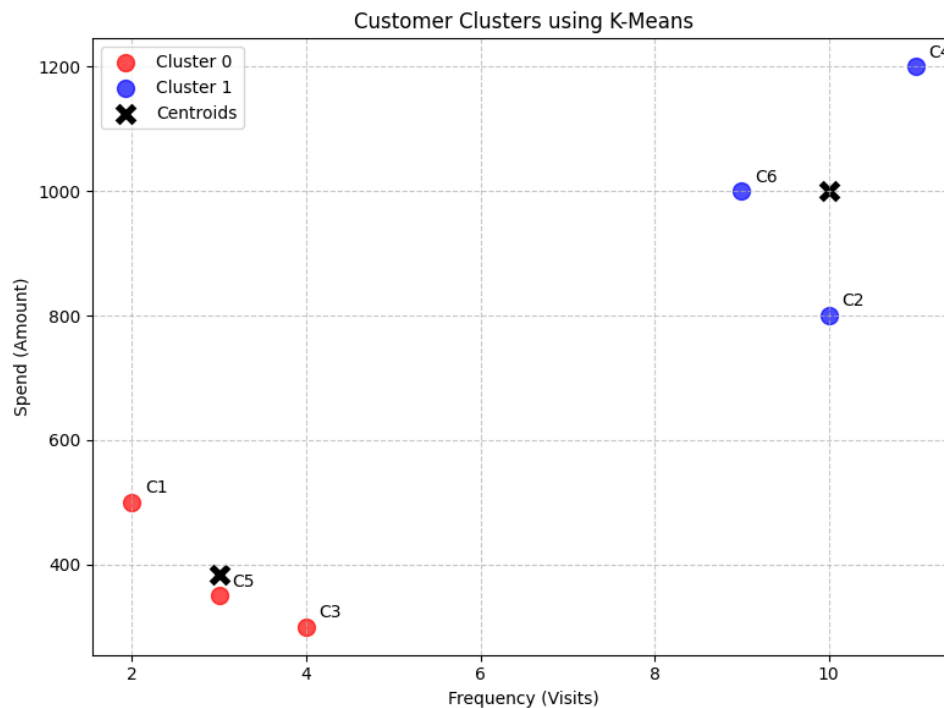
**Listing 3:** K-Means with `scikit-learn` (Random Centroids)

The key difference here is that we **removed the `init` parameter**. By default, `scikit-learn` uses the "k-means++" initialization method, which intelligently chooses initial centroids to speed up convergence. Also, we can tweak the `n_init` parameter, which decides how many

times the K-means algorithm is run with different centroid seeds. (Note: newer versions of scikit-learn might default `n_init` to 'auto' when `init='k-means++'` is used).

## 2.5. Visualizing the Clusters

To gain a better understanding of the results, let's visualize the clusters using a scatter plot:



**Figure 2:** [Placeholder for a Scatter Plot Image showing the Customer Clusters. This plot would visually separate customers based on their frequency and spend, with centroids marked.]

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.cluster import KMeans
4
5 # Data
6 customer_ids = ["C1", "C2", "C3", "C4", "C5", "C6"]
7 data = np.array([
8     [2, 500], [10, 800], [4, 300], [11, 1200], [3, 350], [9, 1000]
9 ])
10
11 # Initial centroids from C1 and C4 (for reproducibility with example)
12 initial_centroids = np.array([ [2, 500], [11, 1200] ])
13
14 # Fit KMeans
15 kmeans = KMeans(n_clusters=2, init=initial_centroids, n_init=1, max_iter=100, random_state=42)
16 kmeans.fit(data)
17 labels = kmeans.labels_
18 centroids = kmeans.cluster_centers_
19
20 # Plot
21 colors = ['red', 'blue']
22 plt.figure(figsize=(8,6)) # Added figure size for better layout
23 for i in range(2):
24     cluster_points = data[labels == i]
25     plt.scatter(cluster_points[:, 0], cluster_points[:, 1], c=colors[i], label=f'Cluster {i}',
26                 s=100, alpha=0.7)
27
28 # Plot centroids
29 plt.scatter(centroids[:, 0], centroids[:, 1], c='black', marker='X', s=200, label='Centroids',
30             edgecolors='w')

```

```

30 # Add customer IDs as labels
31 for i, txt in enumerate(customer_ids):
32     plt.annotate(txt, (data[i, 0] + 0.15, data[i, 1] + 15)) # Adjusted offset
33
34 plt.xlabel('Frequency (Visits)')
35 plt.ylabel('Spend (Amount)')
36 plt.title('Customer Clusters using K-Means')
37 plt.legend()
38 plt.grid(True, linestyle='--', alpha=0.7)
39 plt.tight_layout() # Adjust plot to ensure everything fits
40 plt.show()

```

**Listing 4:** Plotting Customer Clusters with Matplotlib

This code generates a scatter plot where:

- Each customer is represented by a point.
- The color of the point indicates the cluster it belongs to.
- The "X" marks represent the centroids of each cluster.

Visualizing the clusters helps us assess the quality of the clustering and identify potential patterns.

## 2.6. The Elbow Method: Finding the Optimal "K"

Choosing the right number of clusters ("K") is crucial for K-Means. How do we know if we should use 2 clusters, 3 clusters, or more? The elbow method provides a visual approach to determine the optimal "K" value.

```

1 import matplotlib.pyplot as plt
2 from sklearn.cluster import KMeans
3 from sklearn.preprocessing import StandardScaler
4 import pandas as pd
5 import numpy as np
6
7 # Example data (expanded for a more illustrative elbow curve)
8 data_dict = { # Renamed from 'data' to avoid conflict
9     'Frequency': [2, 10, 4, 11, 3, 9, 3, 7, 9, 4, 12, 5, 8, 5, 1, 13, 6, 8, 2, 10],
10    'Spend': [500, 800, 300, 1200, 350, 1000, 750, 800, 1000, 1200, 450, 300, 200, 600, 400,
11             1300, 320, 900, 550, 1100]
12 }
13 df = pd.DataFrame(data_dict)
14
15 # Scale the features
16 scaler = StandardScaler()
17 scaled_data = scaler.fit_transform(df)
18
19 # Try KMeans for different values of k
20 inertia = []
21 K_range = range(1, 11) # Renamed from K
22 for k_val in K_range: # Renamed from k
23     kmeans_model = KMeans(n_clusters=k_val, random_state=42, n_init=10) # Added n_init
24     kmeans_model.fit(scaled_data)
25     inertia.append(kmeans_model.inertia_)
26
27 # Plot the Elbow
28 plt.figure(figsize=(8, 4))
29 plt.plot(K_range, inertia, 'bo-')
30 plt.xlabel('Number of clusters (k)')
31 plt.ylabel('Inertia (Sum of squared distances)')
32 plt.title('Elbow Method For Optimal k')
33 plt.xticks(K_range) # Ensure all k values are shown as ticks
34 plt.grid(True, linestyle='--', alpha=0.7)
35 plt.show()

```

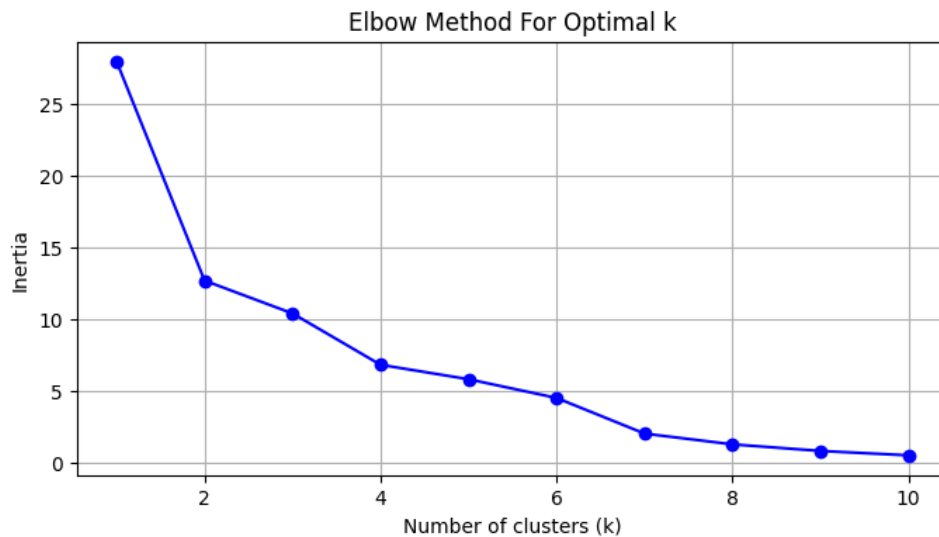
**Listing 5:** Elbow Method for Optimal K

This code does the following:

1. **Calculates Inertia:** For each value of "K" (from 1 to 10), it runs K-Means and calculates the inertia. Inertia is the sum of squared distances of samples to their closest cluster center.

It's a measure of how internally coherent the clusters are. Lower inertia is generally better.

2. **Plots the Elbow Curve:** It plots the inertia values against the corresponding "K" values.



**Figure 3:** [Placeholder for the Elbow Method Curve Image. This graph would show inertia decreasing as K increases, with an "elbow" indicating the optimal K.]

The elbow method suggests choosing the "K" value at the "elbow" of the curve – the point where the rate of decrease in inertia starts to slow down significantly. This point is considered a good balance between minimizing inertia and avoiding overfitting (having too many clusters).

## 2.7. Scaling Up: Clustering a Real-World Dataset

Now, let's tackle a more challenging scenario with a real-world dataset: online retail data. This dataset contains transaction information for an online store, with over 500,000 rows!

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 # import seaborn as sns # Not used in this snippet, can be removed if not needed later
5
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.cluster import KMeans
8 # from sklearn.metrics import silhouette_score # Not used in this snippet
9
10 # Step 1: Load the dataset
11 # Assuming 'Online Retail.xlsx' is in the same directory or path is specified
12 # df = pd.read_excel('Online Retail.xlsx') # This line would be active
13
14 # Create a dummy dataframe for demonstration as the file is not available
15 data_sample = {
16     'InvoiceNo': ['536365', '536365', 'C536379', '536381', '536381', '536382', '536383', '
17     536384', '536385', '536385'],
18     'StockCode': ['85123A', '71053', '22485', '22728', '22727', '21755', '84879', '22423', '
19     85099B', '22633'],
20     'Description': ['CREAM CUPID HEARTS COAT HANGER', 'WHITE METAL LANTERN', 'MANUAL', 'ALARM
21     CLOCK BAKELIKE PINK', 'ALARM CLOCK BAKELIKE RED', 'STRAWBERRY CERAMIC TRINKET BOX', '
22     ASSORTED COLOUR BIRD ORNAMENT', 'REGENCY CAKESTAND 3 TIER', 'JUMBO BAG RED RETROSPOT', '
23     HAND WARMER UNION JACK'],
24     'Quantity': [6, 6, -1, 2, 2, 12, 32, 1, 10, 6],
25     'InvoiceDate': ['12/1/2010 8:26', '12/1/2010 8:26', '12/1/2010 9:41', '12/1/2010 9:41', '
26     12/1/2010 9:41', '12/1/2010 9:45', '12/1/2010 9:49', '12/1/2010 9:53', '12/1/2010 9:56', '
27     12/1/2010 9:56'],
28     'UnitPrice': [2.55, 3.39, 2.1, 3.75, 3.75, 1.65, 1.69, 12.75, 1.95, 1.85],
29     'CustomerID': [17850.0, 17850.0, 14527.0, 15311.0, 15311.0, 16133.0, 13047.0, 18074.0,
30     17850.0, 17850.0]
31 }

```

```

23 }
24 df = pd.DataFrame(data_sample)
25 df['InvoiceDate'] = pd.to_datetime(df['InvoiceDate'])
26
27
28 # Step 2: Data Preprocessing
29 # Remove the rows with missing customer ids
30 df = df[pd.notnull(df['CustomerID'])]
31
32 # Remove canceled orders (InvoiceNo starting with 'C')
33 df = df[~df['InvoiceNo'].astype(str).str.startswith('C')]
34
35 # Remove negative quantities (and zero, if applicable)
36 df = df[df['Quantity'] > 0]
37
38 # Step 3: Feature Engineering
39 # Calculate TotalPrice
40 df['TotalPrice'] = df['Quantity'] * df['UnitPrice']
41
42 # Aggregate data by CustomerID
43 customer_df = df.groupby('CustomerID').agg(
44     Frequency=('InvoiceNo', 'nunique'), # Frequency
45     TotalQuantity=('Quantity', 'sum'),
46     Monetary=('TotalPrice', 'sum')
47 ).reset_index()
48
49 # Step 4: Feature Scaling
50 features = ['Frequency', 'TotalQuantity', 'Monetary']
51 scaler = StandardScaler()
52 scaled_features = scaler.fit_transform(customer_df[features])
53
54 # Step 5: Determine the optimal number of clusters using the Elbow Method
55 sse = [] # Sum of Squared Errors (Inertia)
56 k_range = range(1, 11)
57 for k in k_range:
58     kmeans = KMeans(n_clusters=k, random_state=42, n_init=10) # Added n_init
59     kmeans.fit(scaled_features)
60     sse.append(kmeans.inertia_)
61
62 print(sse)
63 print("\n") # For spacing in output if run
64
65 # Plot the Elbow Curve
66 plt.figure(figsize=(8, 5))
67 plt.plot(k_range, sse, marker='o')
68 plt.title('Elbow Method for Online Retail Data')
69 plt.xlabel('Number of clusters (k)')
70 plt.ylabel('Sum of squared distances (Inertia)')
71 plt.xticks(k_range)
72 plt.grid(True, linestyle='--', alpha=0.7)
73 plt.show()

```

**Listing 6:** Processing Online Retail Data and Elbow Method

This code performs the following steps:

1. **Data Loading:** Loads the online retail dataset (here, a sample is used for demonstration).
2. **Data Preprocessing:**
  - Removes rows with missing customer IDs.
  - Removes cancelled orders (identified by "C" in the InvoiceNo).
  - Removes rows with negative quantities (indicating returns).
3. **Feature Engineering:**
  - Calculates the total price for each transaction (Quantity \* UnitPrice).
  - Aggregates the data by customer ID to create features like frequency (number of unique invoices), total quantity purchased, and total monetary value.
4. **Feature Scaling:** Scales the features using StandardScaler. This is important because K-Means is sensitive to the scale of the features.



5. **Elbow Method:** Calculates and plots the inertia for different values of "K" to help determine the optimal number of clusters.

## 2.8. Silhouette Score: A Complementary Metric

While the elbow method provides a visual guide, we can also use a more quantitative measure to evaluate the clustering performance: the silhouette score.

```
1 # Continuing from the previous section (scaled_features is available)
2 from sklearn.metrics import silhouette_score
3 # Ensure scaled_features is defined from the previous step.
4 # For standalone running, you'd need to regenerate customer_df and scaled_features.
5 # Assuming scaled_features from section 2.7 is in memory.
6 # If not, re-run the data prep steps from 2.7 or load a pre-scaled dataset.
7 # For this example, let's re-use the 'scaled_data' from section 2.6 for a quick run,
8 # as 'scaled_features' from 2.7 depends on the larger (dummy) retail dataset.
9 # In a real notebook, 'scaled_features' would persist.
10 # Using 'scaled_data' (from the smaller example in 2.6) for silhouette plot:
11 # Example: scaled_features_for_silhouette = scaled_data
12 # If using the Online Retail data, 'scaled_features' should be used.
13 # Here, we'll assume 'scaled_features' is ready from the previous block.
14 # If it's not (e.g., running this section independently with a new kernel),
15 # the following will error or use a different 'scaled_features'.
16 # For demonstration, let's assume scaled_features (from retail data) is ready.
17 # If 'scaled_features' is small (e.g. from dummy data), the plot might not be very informative
18 .
19 # Step 6: Evaluate clustering with Silhouette Score
20 silhouette_scores = []
21 k_range_silhouette = range(2, 11) # Silhouette score is not defined for k=1
22 for k in k_range_silhouette:
23     kmeans = KMeans(n_clusters=k, random_state=42, n_init=10) # Added n_init
24     labels = kmeans.fit_predict(scaled_features) # Use scaled_features from retail data
25     score = silhouette_score(scaled_features, labels)
26     silhouette_scores.append(score)
27     print(f'For n_clusters = {k}, the Silhouette Score is {score:.4f}')
28
29 # Plot Silhouette Scores
30 plt.figure(figsize=(8, 5))
31 plt.plot(k_range_silhouette, silhouette_scores, marker='o')
32 plt.title('Silhouette Scores for Various Clusters (Online Retail Data)')
33 plt.xlabel('Number of clusters (k)')
34 plt.ylabel('Silhouette Score')
35 plt.xticks(k_range_silhouette)
36 plt.grid(True, linestyle='--', alpha=0.7)
37 plt.show()
```

Listing 7: Silhouette Score Evaluation

- **Silhouette Score:** The silhouette score measures how well each data point fits within its cluster. It ranges from -1 to +1, where:
  - +1: Indicates that the data point is well-clustered.
  - 0: Indicates that the data point is close to the decision boundary between two clusters.
  - -1: Indicates that the data point might be assigned to the wrong cluster.

We calculate the silhouette score for different values of "K" and choose the "K" that maximizes the silhouette score.

## 2.9 Next Steps

Finally, we would choose the optimal number of clusters based on both the elbow method and the silhouette score. Then, we would apply K-Means with that optimal "K" value, assign cluster labels to the customers, and analyze the characteristics of each cluster.

## 2.10. Conclusion

This concludes our journey into the world of K-Means clustering with code. We've explored how to implement K-Means using `scikit-learn`, how to choose the optimal number of clusters using the elbow method and the silhouette score, and how to apply K-Means to a real-world dataset.

Remember, K-Means is a powerful tool for segmentation and can be applied to a wide range of problems. So, go out there and start clustering!