

Homework Spectral Clustering

Preparation:

First, we import the required libraries:

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import sklearn
5 import matplotlib.pyplot as plt
6 import scipy.sparse as sp
7 from scipy.sparse.csgraph import connected_components
8 import scipy.sparse.linalg as splinalg
9 from scipy.sparse.linalg import eigsh
10 from scipy.sparse import csr_matrix
11 from scipy.sparse.csgraph import laplacian
12 from scipy.linalg import eig
13 from sklearn.cluster import KMeans, AgglomerativeClustering, DBSCAN, SpectralClustering, Birch
14 from sklearn.datasets import make_blobs
15 from sklearn.metrics import adjusted_rand_score
16 from sklearn.metrics import silhouette_score, calinski_harabasz_score, davies_bouldin_score
17
```

The next step is loading the datasets.

First, we load the circle dataset:

```
1 X = np.loadtxt("/content/drive/MyDrive/CLA/Clustering/Circle.csv", delimiter=',')
2 plt.scatter(X[:, 0], X[:, 1], s=10)
3 plt.title("Dataset: Circle")
4 plt.show()
5
```

Then we want to see the structure of this dataset. To do this we see the first 5 elements of the dataset:

```
1 circle_df = pd.read_csv('/content/drive/MyDrive/CLA/Clustering/Circle.csv', header=None)
2 circle_df.columns = ['x', 'y']
3 circle_df.head()
```

	x	y
0	2.7147	0.81472
1	2.8054	0.94572
2	2.0253	0.20682
3	2.8096	1.03310
4	2.5257	0.79188

Then we repeat these two steps to see the structure of spiral dataset as well:

```
1 X = np.loadtxt("/content/drive/MyDrive/CLA/Clustering/Spiral.csv", delimiter=',')
2 plt.scatter(X[:, 0], X[:, 1], s=10)
3 plt.title("Dataset: Spiral")
4 plt.show()
5
```

```
1 spiral_df = pd.read_csv('/content/drive/MyDrive/CLA/Clustering/Spiral.csv', header=None)
2 spiral_df.columns = ['x', 'y', 'label']
3 spiral_df.head()
```

	x	y	label
0	31.95	7.95	3
1	31.15	7.30	3
2	30.45	6.65	3
3	29.70	6.00	3
4	28.90	5.55	3

Task No.1

Here we want to implement task 1. In order to do so, first we need to define a function to compute the similarity between each pair of points. function 'compute_similarity_matrix' does this task. actually, it takes a set of points and calculates the distance between each pair of the dataset using the metric presented in the problem:

$$s_{i,j} = \exp\left(-\frac{(\|x_i - x_j\|^2)}{2\sigma^2}\right)$$

```
1 def compute_similarity_matrix(X, sigma):
2     N = X.shape[0]
3     similarity_matrix = np.exp(-np.sum((X[:, np.newaxis] - X)**2, axis=2) / (2 * sigma**2))
4     np.fill_diagonal(similarity_matrix, 0)
5     return similarity_matrix
```

```
1 spiral_x = spiral_df[['x', 'y']].values
2 circle_x = circle_df[['x', 'y']].values
3 sigma = 1
4 Spiral_similarity_matrix = compute_similarity_matrix(spiral_x, sigma)
5 circle_similarity_matrix = compute_similarity_matrix(circle_x, sigma)
6 print('Spiral Shape: ', Spiral_similarity_matrix.shape)
7 print('Circle Shape: ', circle_similarity_matrix.shape)
```

Using the Similarity matrix defined above, here we try to find k nearest neighbor of each point and construct a matrix, such that each row (and consequently each column) has k non-zero values showing which points are the corresponding points's close neighbours.

Note that the (spiral_knn_adj_matrix, circle_knn_adj_matrix) are adjacency matrices for each dataset.

```

1 def compute_knn_similarity_matrix(similarity_matrix, k, sigma):
2     knn_similarity_matrix = np.zeros_like(similarity_matrix)
3     for i in range(similarity_matrix.shape[0]):
4         k_nearest_neighbors = np.argsort(similarity_matrix[i])[-k:]
5         knn_similarity_matrix[i, k_nearest_neighbors] = similarity_matrix[i, k_nearest_neighbors]
6         knn_similarity_matrix[k_nearest_neighbors, i] = similarity_matrix[k_nearest_neighbors, i]
7     return knn_similarity_matrix
8
9
10 k = 10
11 spiral_knn_adj_matrix = compute_knn_similarity_matrix(Spiral_similarity_matrix, k, sigma)
12
13 circle_knn_adj_matrix = compute_knn_similarity_matrix(circle_similarity_matrix, k, sigma)
14
15 #check if the diagonal elements are zero
16 print('Spiral dataset: ', np.diag(spiral_knn_adj_matrix).sum())
17 print('Circle dataset: ', np.diag(circle_knn_adj_matrix).sum())

```

Task No2.

Here we construct degree and Laplacian matrices, which are simple in definition however I describe them here briefly:

- Degree Matrix: degree matrix of an undirected graph is a diagonal matrix which contains information about the degree of each node—that is, the number of edges attached to each node. We use degree matrix to construct Laplacian matrix.
- Laplacian Matrix: Is defined as the differences of degree matrix and adjacency matrix. It captures the relationship between nodes and edges in a graph and is widely used

```

1 def compute_degree_matrix(adjacency_matrix):
2     adjacency_matrix = csr_matrix(adjacency_matrix)
3     degree_vector = np.asarray(adjacency_matrix.sum(axis=1)).flatten()
4     return csr_matrix(np.diag(degree_vector))
5
6 spiral_laplaican_matrix = compute_degree_matrix(spiral_knn_adj_matrix) - spiral_knn_adj_matrix
7 circle_laplaican_matrix = compute_degree_matrix(circle_knn_adj_matrix) - circle_knn_adj_matrix
8 print('Spiral dataset: \n', spiral_laplaican_matrix)
9 print("=====")
10 print('\nCircle dataset: \n', circle_laplaican_matrix)

```

The result:

```

Spiral dataset:
[[ 0.7400033 -0.58786967 -0.13945686 ... 0. 0.
  0. ]
 [-0.58786967 1.38989147 -0.6336554 ... 0. 0.
  0. ]
 [-0.13945686 -0.6336554 1.56387582 ... 0. 0.
  0. ]
 ...
 [ 0. 0. 0. ... 7.50068164 -0.96923323
 -0.94058806]
 [ 0. 0. 0. ... -0.96923323 6.54768684
 -0.99501248]
 [ 0. 0. 0. ... -0.94058806 -0.99501248
 6.10367257]]

=====

Circle dataset:
[[ 9.56687889 -0.98738648 0. ... 0. 0.
  0. ]
 [-0.98738648 10.39869095 0. ... 0. 0.
  0. ]
 [ 0. 0. 10.29397786 ... 0. 0.
  0. ]
 ...
 [ 0. 0. 0. ... 8.77786359 0.
  0. ]
 [ 0. 0. 0. ... 0. 11.61410051
  0. ]
 [ 0. 0. 0. ... 0. 0.
  7.533866 ]]

```

Task No3.

Here we try to compute the number of connected components of the similarity graph. In other words, we want to derive how many classes are proper for classifying the points.

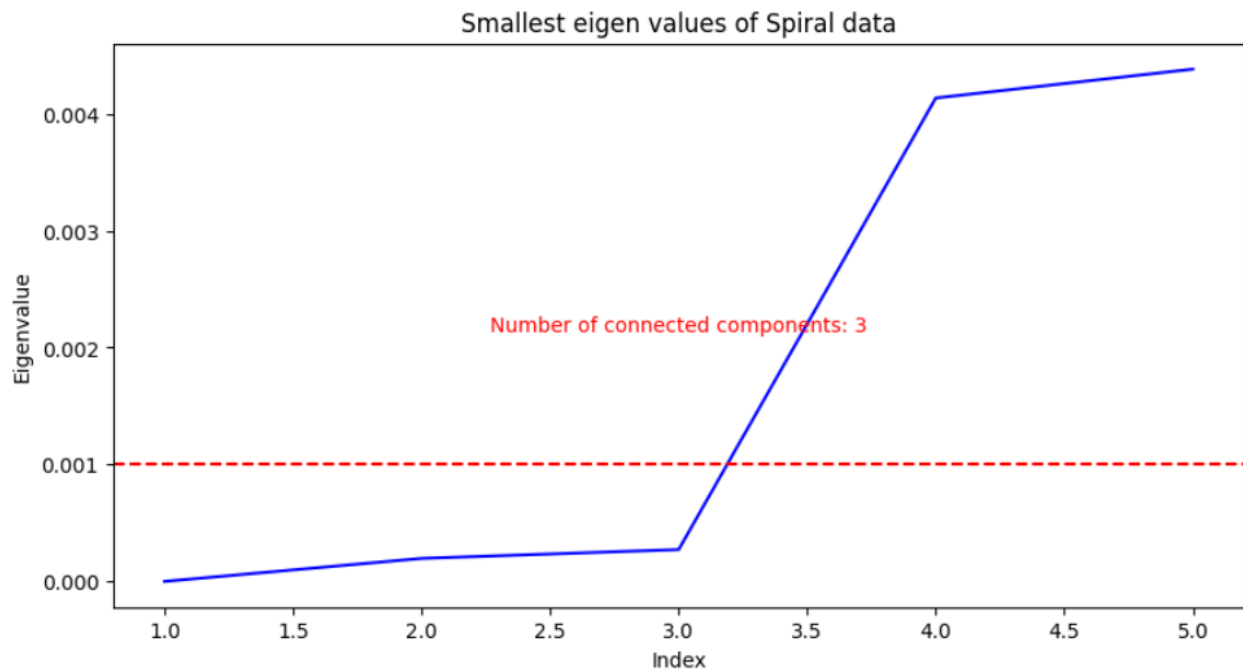
Our way to reach this purpose is to find the factor of zero eigenvalues of the laplacian matrix for each dataset. Keep in mind that in the concept of similarity metric (points are not connected explicitly), by zero eigenvalue, we mean eigenvalues below a small threshold.

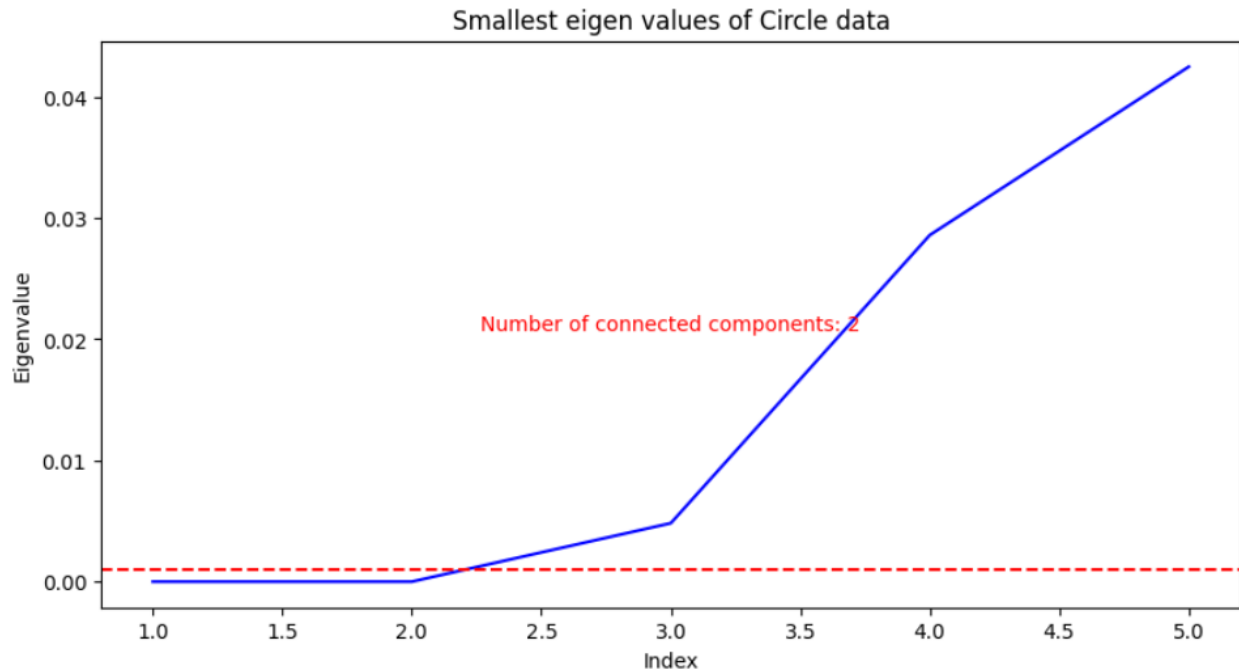
Function 'compute_sorted_eigenvalues_eigenvectors()', calculates the first m eigenvectors and eigenvalues of the input matrix. by first, I mean first eigenvectors based on the sorting of the eigenvalues(eigenvectors corresponding to the smallest eigenvectors).

Then, function 'compute_connected_components()' compares each eigenvalue with a relatively small threshold and outputs number of the eigenvalues which are bigger than the threshold value.

```
1 def compute_sorted_eigenvalues_eigenvectors(laplacian_matrix, m = None):
2     eigenvalues, eigenvectors = eigsh(laplacian_matrix, k=5, which='SM')
3     sorted_indices = np.argsort(eigenvalues)
4     if m is None:
5         return eigenvalues[sorted_indices], eigenvectors[:, sorted_indices]
6     return eigenvalues[sorted_indices][:m], eigenvectors[:, sorted_indices][:, :m]
7
8
9 def compute_connected_components(laplacian_matrix, threshold=1.0e-4, m=5):
10     eigenvalues, _ = compute_sorted_eigenvalues_eigenvectors(laplacian_matrix, m)
11     return sum(eigen_val < threshold for eigen_val in eigenvalues)
12
13
14 def plot_eigenvalues(sorted_eigenvalues, title, threshold=1.0e-3):
15     plt.figure(figsize=(10, 5))
16     plt.plot(range(1, len(sorted_eigenvalues) + 1), sorted_eigenvalues, 'b-')
17     plt.xlabel('Index')
18     plt.ylabel('Eigenvalue')
19     num_connected_components = sum(eigen_val < threshold for eigen_val in sorted_eigenvalues)
20     plt.text(0.5, 0.5, f'Number of connected components: {num_connected_components}', horizontalalignment='center', verticalalignment='center', transform=plt.gca().transData)
21     plt.axhline(y=threshold, color='r', linestyle='--')
22     plt.title(title)
23     plt.show()
```

```
1 spiral_num_connected_components = compute_connected_components(spiral_laplaican_matrix)
2 circle_num_connected_components = compute_connected_components(circle_laplaican_matrix)
3
4 for matrix, title, th in [(spiral_laplaican_matrix, "Smallest eigen values of Spiral data", 1e-3),
5                             (circle_laplaican_matrix, "Smallest eigen values of Circle data", 1e-2)]:
6     eigenvalues, _ = compute_sorted_eigenvalues_eigenvectors(matrix)
7     plot_eigenvalues(eigenvalues, title, threshold=th)
8
```





Task No4.

As completely computed above, we assume the number of clusters is the same as the number of the connected components. Also the function to compute both eigen vectors and eigenvalues was defined in the previous section, therefore, here we just run that function:

```
1 sorted_eigen_val_circle, sorted_eigen_vec_circle = compute_sorted_eigenvalues_eigenvectors(circle_laplaican_matrix, 3)
2 sorted_eigen_val_spiral, sorted_eigen_vec_spiral = compute_sorted_eigenvalues_eigenvectors(spiral_laplaican_matrix, 3)
```

Task No5.

In the code blocks of Task 3, we defined a function that computes both eigenvalues and eigenvectors of a matrix. Also we have to again note that the eigenvalues were sorted (and their corresponding eigenvectors), therefore, here we don't need to do anything more. We just normalize the previous results.

```
1 spiral_u = sorted_eigen_vec_spiral / np.linalg.norm(sorted_eigen_vec_spiral, axis=1, keepdims=True)
2 circle_u = sorted_eigen_vec_circle / np.linalg.norm(sorted_eigen_vec_circle, axis=1, keepdims=True)
```

Task No6 and 7.

Here we are utilizing the KMeans clustering algorithm from the sklearn.cluster module in Python to group data into clusters.

```

1 from sklearn.cluster import KMeans
2
3 def perform_clustering(original_data, u_normalized, n_clusters):
4     # Perform k-means clustering and return the clustered data
5     kmeans = KMeans(n_clusters=n_clusters).fit(u_normalized)
6     clusters = [[] for _ in range(n_clusters)]
7     for i, label in enumerate(kmeans.labels_):
8         clusters[label].append(original_data[i])
9     return clusters, kmeans.labels_
10
11 # Spiral and Circle clustering
12 spiral_clusters, spiral_clusters_labels = perform_clustering(spiral_x, spiral_u, n_clusters=3)
13 circle_clusters, circle_clusters_labels = perform_clustering(circle_x, circle_u, n_clusters=3)
14

```

- Note: This loop iterates through the labels assigned by the KMeans algorithm (kmeans.labels_). For each data point, it appends the corresponding original data point to its respective cluster based on the assigned label.

```

7     for i, label in enumerate(kmeans.labels_):
8         clusters[label].append(original_data[i])

```

The function(perform_clustering)returns two values:

clusters: A list of clusters, where each cluster contains the original data points.

kmeans.labels_ : The labels assigned to each data point, indicating which cluster they belong to.

Task No8.

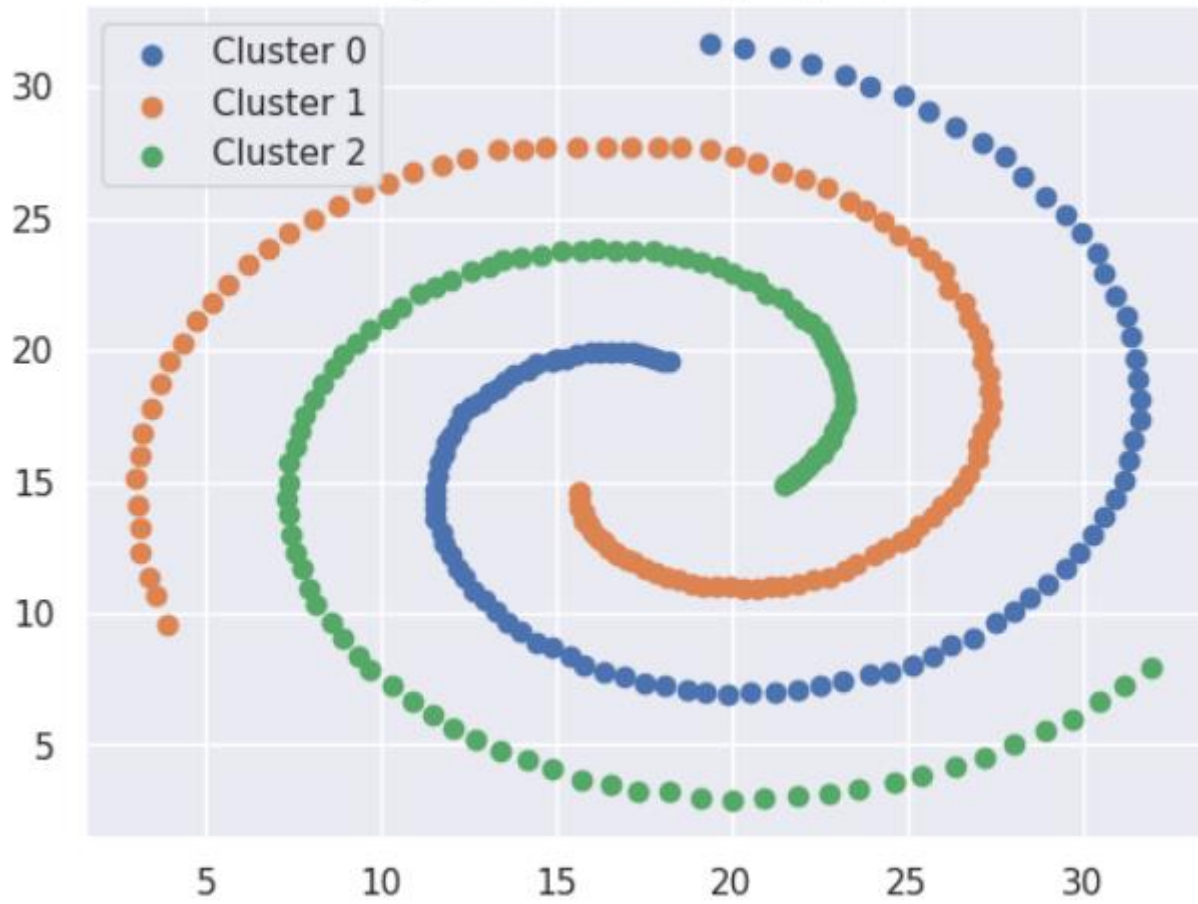
Here, we visualize different clusters using scatter plots by defining a reusable function. It uses Seaborn for styling and Matplotlib for plotting, presenting a clear visual representation of how the clusters are distributed in the dataset.

```

1 def plot_clusters(data, clusters, title):
2     sns.set(style="darkgrid")
3     fig, ax = plt.subplots()
4     for i, cluster in enumerate(clusters):
5         ax.scatter(np.array(cluster)[: , 0], np.array(cluster)[: , 1], label=f'Cluster {i}')
6     ax.set_title(title)
7     ax.legend()
8     plt.show()
9
10 # Plot the clusters for both datasets
11 for clusters, title, data in [(spiral_clusters, 'Spectral Clustering - Spiral', spiral_x),
12                               (circle_clusters, 'Spectral Clustering - Circle', circle_x)]:
13     plot_clusters(data, clusters, title)
14

```

Spectral Clustering - Spiral





In the next step we want to evaluate the performance of clustering results by using three metrics: **Silhouette Score**, **Calinski-Harabasz Index**, and **Davies-Bouldin Index**. These metrics show how well the data points are grouped into clusters.

```
1 from sklearn.metrics import silhouette_score, calinski_harabasz_score, davies_bouldin_score
2
3 def evaluate_clustering_performance(data, clusters_labels):
4     # Compute and return the clustering performance metrics
5     return (silhouette_score(data, clusters_labels, metric='euclidean'),
6             calinski_harabasz_score(data, clusters_labels),
7             davies_bouldin_score(data, clusters_labels))
8
9 # Evaluate clustering performance for Spiral and Circle datasets
10 spiral_metrics = evaluate_clustering_performance(spiral_x, spiral_clusters_labels)
11 circle_metrics = evaluate_clustering_performance(circle_x, circle_clusters_labels)
12
13 # Unpack the results
14 spiral_silhouette_coefficient, spiral_calinski_harabasz_score_coefficient, spiral_davies_bouldin_score_coefficient = spiral_metrics
15 circle_silhouette_coefficient, circle_calinski_harabasz_score_coefficient, circle_davies_bouldin_score_coefficient = circle_metrics
16
```

Silhouette Score: Checks how similar points within the same cluster are compared to points in the nearest cluster. Helps assess whether the clusters are compact and well-separated.

Calinski-Harabasz Index: Measures cluster dispersion, balancing intra-cluster and inter-cluster variance. Higher values mean better-defined clusters.

Davies-Bouldin Index: Captures the tradeoff between cluster compactness and separation. Lower values indicate more compact and better-separated clusters.

In the next cell, we want to compare our clustering accuracy by these three metrics.

```
1 # Define the metrics for Spiral and Circle datasets
2 metrics = {
3     'Silhouette Coefficient': [spiral_silhouette_coefficient, circle_silhouette_coefficient],
4     'Calinski-Harabasz Score': [spiral_calinski_harabasz_score_coefficient, circle_calinski_harabasz_score_coefficient],
5     'Davies-Bouldin Score': [spiral_davies_bouldin_score_coefficient, circle_davies_bouldin_score_coefficient]
6 }
7
8 # Create the results DataFrame
9 results = pd.DataFrame(metrics, index=['Spiral', 'Circle'])
10 results.index.name = 'Metrics'
11 results
12
```

Here is the result:

	Silhouette Coefficient	Calinski-Harabasz Score	Davies-Bouldin Score
Metrics			
Spiral	0.001344	5.797852	5.882023
Circle	0.323268	902.993795	82.265848

The clustering metrics show some challenges in clustering spiral-shaped data:

Overlap between clusters is evident due to the low Silhouette Coefficient, and clusters are separated not in the best way and dispersed due to low Calinski-Harabasz and high Davies-Bouldinrate. However, the Circle dataset achieves successful clustering. Due to the Silhouette Coefficient which indicates some separation, and the high Calinski-Harabasz score reflects reasonably compact clusters.

Task No9.

The first step is converting the dataset into a C-contiguous array for optimized memory access and compatibility with some clustering algorithms. Then we use four methods (K-Means, Agglomerative, Birch, Spectral) to cluster our circle and spiral dataset.

```

1 # Apply the clustering methods mentioned above to the datasets
2
3 def perform_clustering(data, n_clusters, clustering_method):
4
5     # Ensure data is in C-contiguous order.
6     data = np.ascontiguousarray(data)
7
8     # Perform clustering.
9     if clustering_method == AgglomerativeClustering:
10         clustering = clustering_method(n_clusters=n_clusters, linkage='ward').fit(data) # Removed affinity
11     elif clustering_method == Birch:
12         clustering = clustering_method(n_clusters=n_clusters).fit(data)
13     elif clustering_method == SpectralClustering:
14         clustering = clustering_method(n_clusters=n_clusters, affinity='nearest_neighbors').fit(data)
15     else:
16         clustering = clustering_method(n_clusters=n_clusters).fit(data)
17
18     # Group the data based on the cluster labels.
19     clusters = [[] for _ in range(n_clusters)]
20     for i in range(data.shape[0]):
21
22         # Append the data point to the corresponding cluster.
23         clusters[clustering.labels_[i]].append(data[i])
24
25     return clusters, clustering.labels_
26
27
28 def plot_clusters(data, clusters, title):
29     sns.set(style="darkgrid")
30     fig, ax = plt.subplots()
31
32     for i in range(len(clusters)):
33         cluster = np.array(clusters[i])
34         # if array is one-dimensional, reshape it to two-dimensional
35         if len(cluster.shape) == 1:
36             cluster = cluster.reshape(-1, 1)
37
38         ax.scatter(cluster[:, 0], cluster[:, 1], label='Cluster {}'.format(i))
39     ax.set_title(title)
40     ax.legend()
41     plt.show()
42
43 # Spiral
44 spiral_kmeans_clusters, spiral_kmeans_clusters_labels = perform_clustering(spiral_x, n_clusters=3, clustering_method=KMeans)
45 spiral_agglomerative_clusters, spiral_agglomerative_clusters_labels = perform_clustering(spiral_x, n_clusters=3, clustering_method=AgglomerativeClustering)
46 spiral_birch_clusters, spiral_birch_clusters_labels = perform_clustering(spiral_x, n_clusters=3, clustering_method=Birch)
47 spiral_spec_clusters, spiral_spec_clusters_labels = perform_clustering(spiral_x, n_clusters=3, clustering_method=SpectralClustering)
48
49 # Circle
50 circle_kmeans_clusters, circle_kmeans_clusters_labels = perform_clustering(circle_x, n_clusters=3, clustering_method=KMeans)
51 circle_agglomerative_clusters, circle_agglomerative_clusters_labels = perform_clustering(circle_x, n_clusters=3, clustering_method=AgglomerativeClustering)
52 circle_birch_clusters, circle_birch_clusters_labels = perform_clustering(circle_x, n_clusters=3, clustering_method=Birch)
53 circle_spec_clusters, circle_spec_clusters_labels = perform_clustering(circle_x, n_clusters=3, clustering_method=SpectralClustering)
54
55 # Plot the clusters for Spiral dataset
56 plot_clusters(spiral_x, spiral_kmeans_clusters, 'K-Means Clustering - Spiral')
57 plot_clusters(spiral_x, spiral_agglomerative_clusters, 'Agglomerative Clustering - Spiral')
58 plot_clusters(spiral_x, spiral_birch_clusters, 'Birch Clustering - Spiral')
59 plot_clusters(spiral_x, spiral_spec_clusters, 'Spectral Clustering - Spiral')
60
61 # Plot the clusters for Circle dataset
62 plot_clusters(circle_x, circle_kmeans_clusters, 'K-Means Clustering - Circle')
63 plot_clusters(circle_x, circle_agglomerative_clusters, 'Agglomerative Clustering - Circle')
64 plot_clusters(circle_x, circle_birch_clusters, 'Birch Clustering - Circle')
65 plot_clusters(circle_x, circle_spec_clusters, 'Spectral Clustering - Circle')
66

```

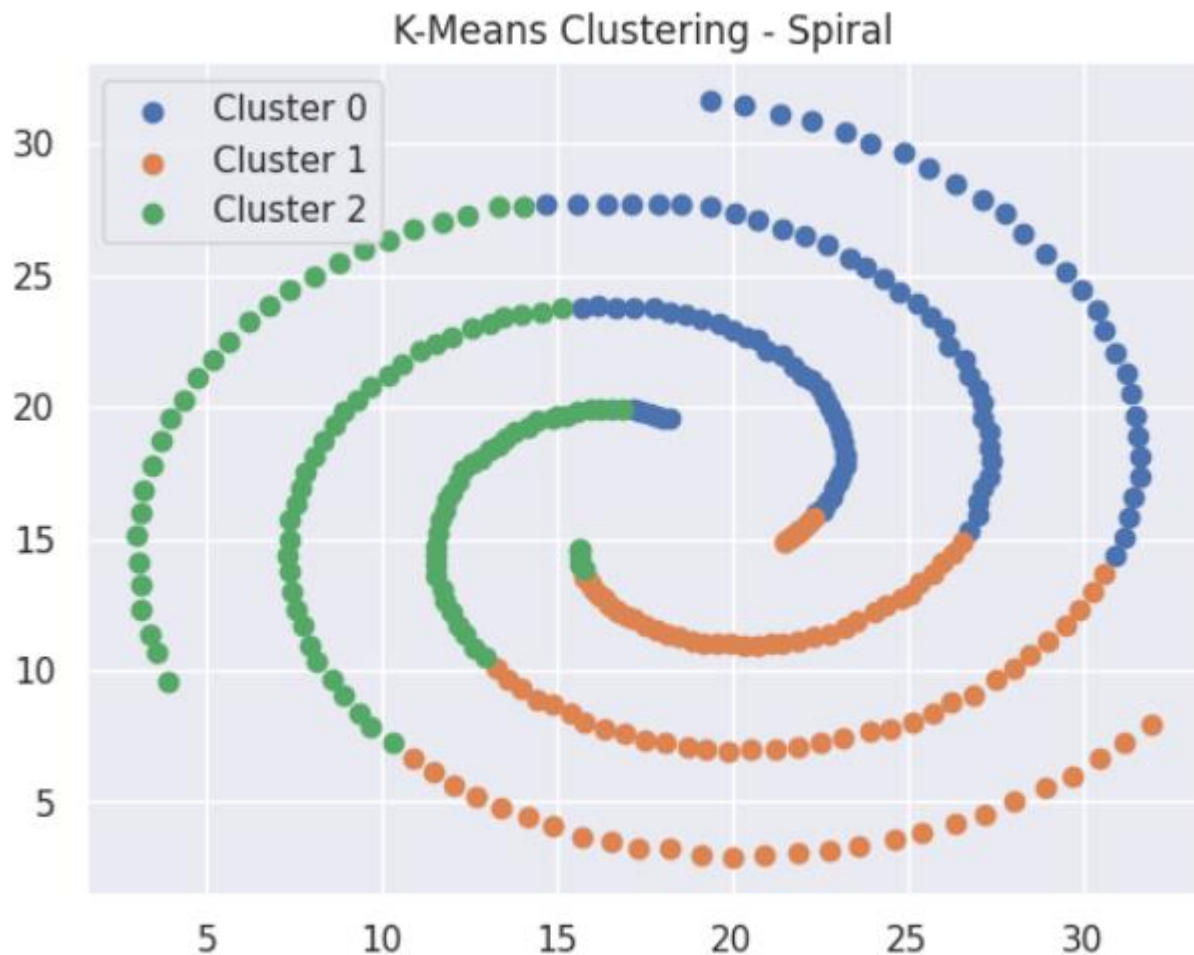
K-Means Clustering: Partitions data into K clusters by minimizing intra-cluster variance but may not be satisfying for non-linear clusters such as Spiral and Circle.

Agglomerative Clustering: Merges clusters iteratively based on similarity and uses Ward's linkage, which minimizes the sum of squared differences. Actually, it works better on some non-linear clusters in comparison to K-Means.

Birch Clustering: Builds a clustering tree incrementally for large datasets.

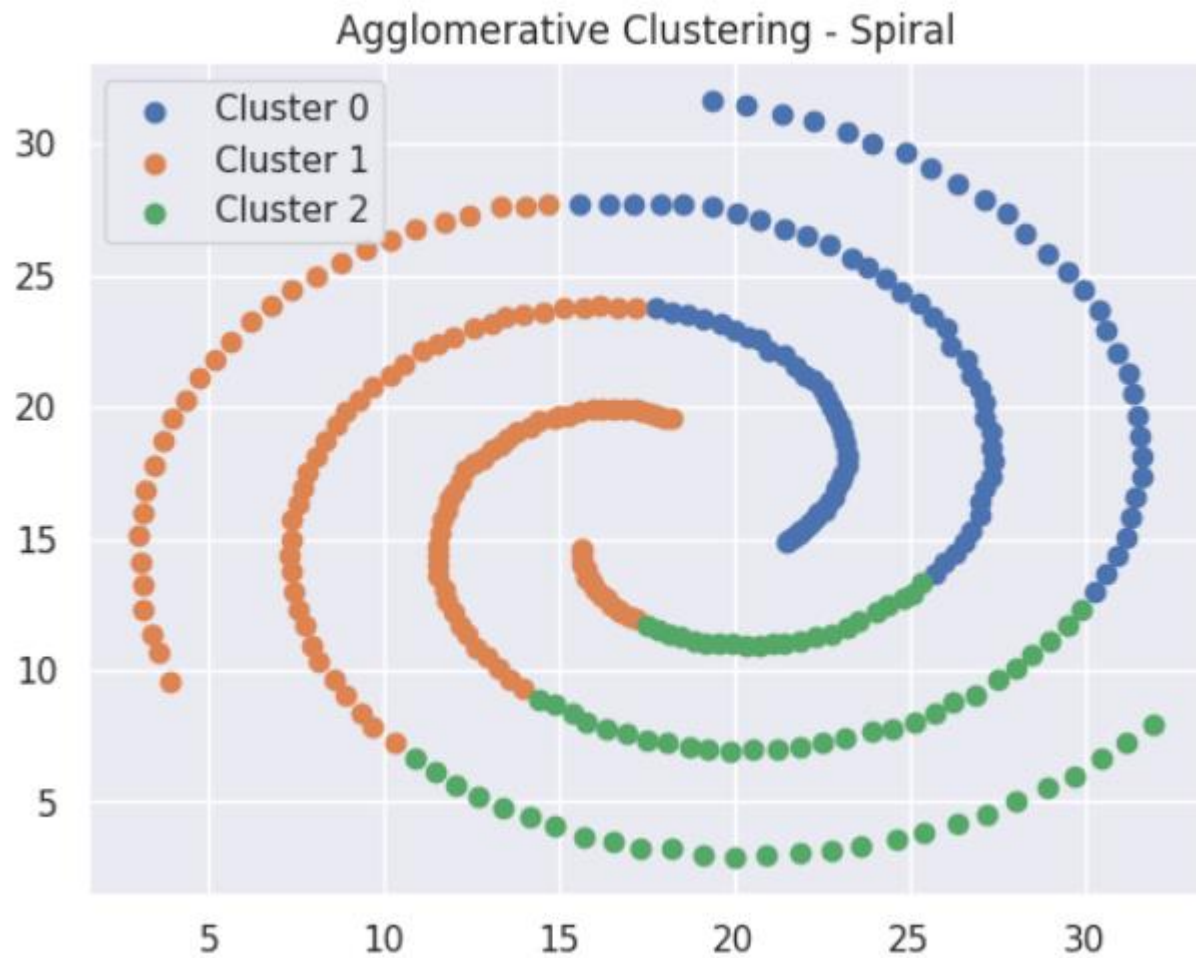
Spectral Clustering: Graph-based clustering using eigenvalues of the Laplacian matrix and captures non-linear cluster structures. It is the best option for non-linear clusters like Spiral and Circle.

Results for spiral dataset:

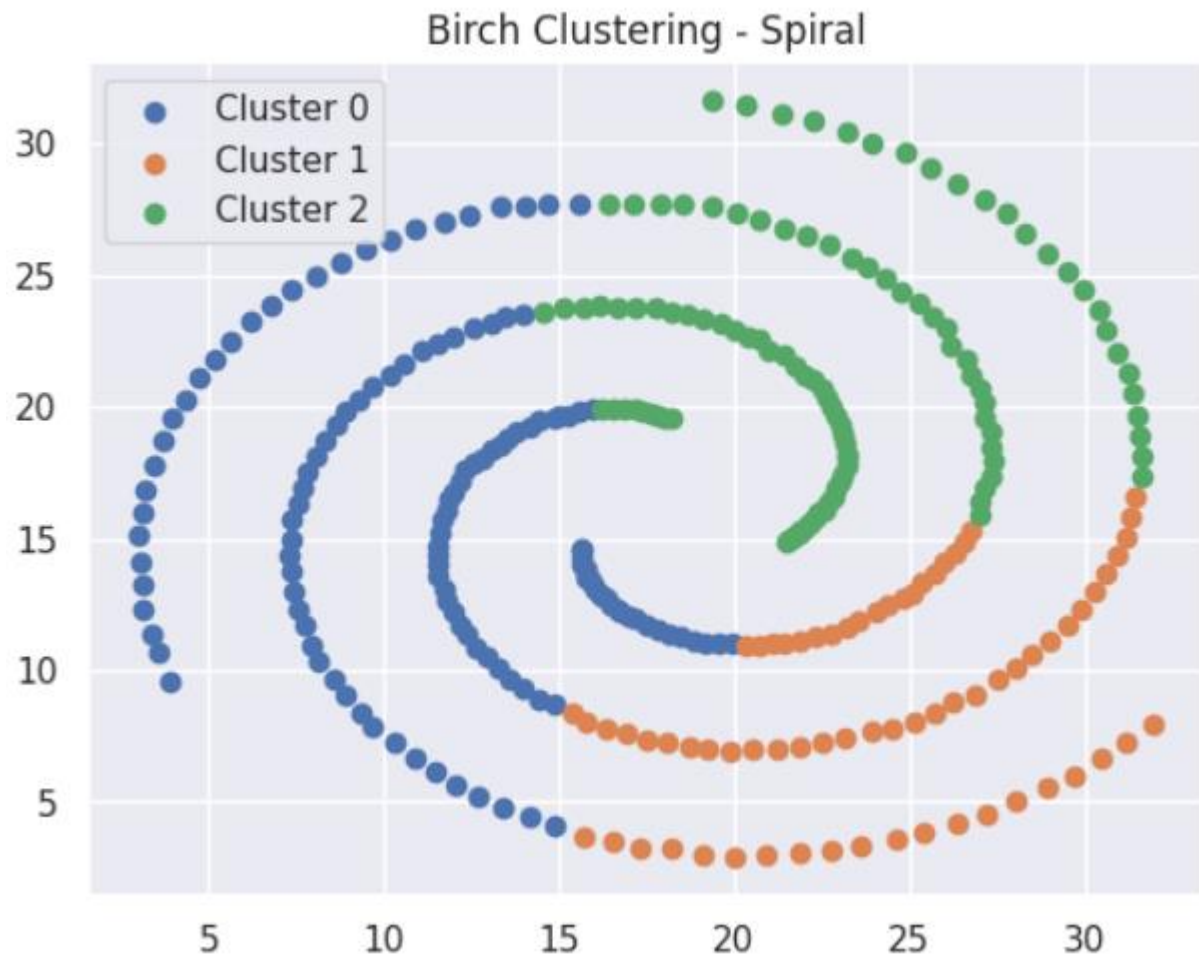


The algorithm divides the data into linear segments, which is a limitation of K-Means due to its reliance on Euclidean distances and spherical cluster assumptions.

Clusters are artificially divided along radial segments rather than aligning with the spiral structure. Therefore, K-Means is unsuitable for the spiral dataset because the clusters are non-linearly separable.

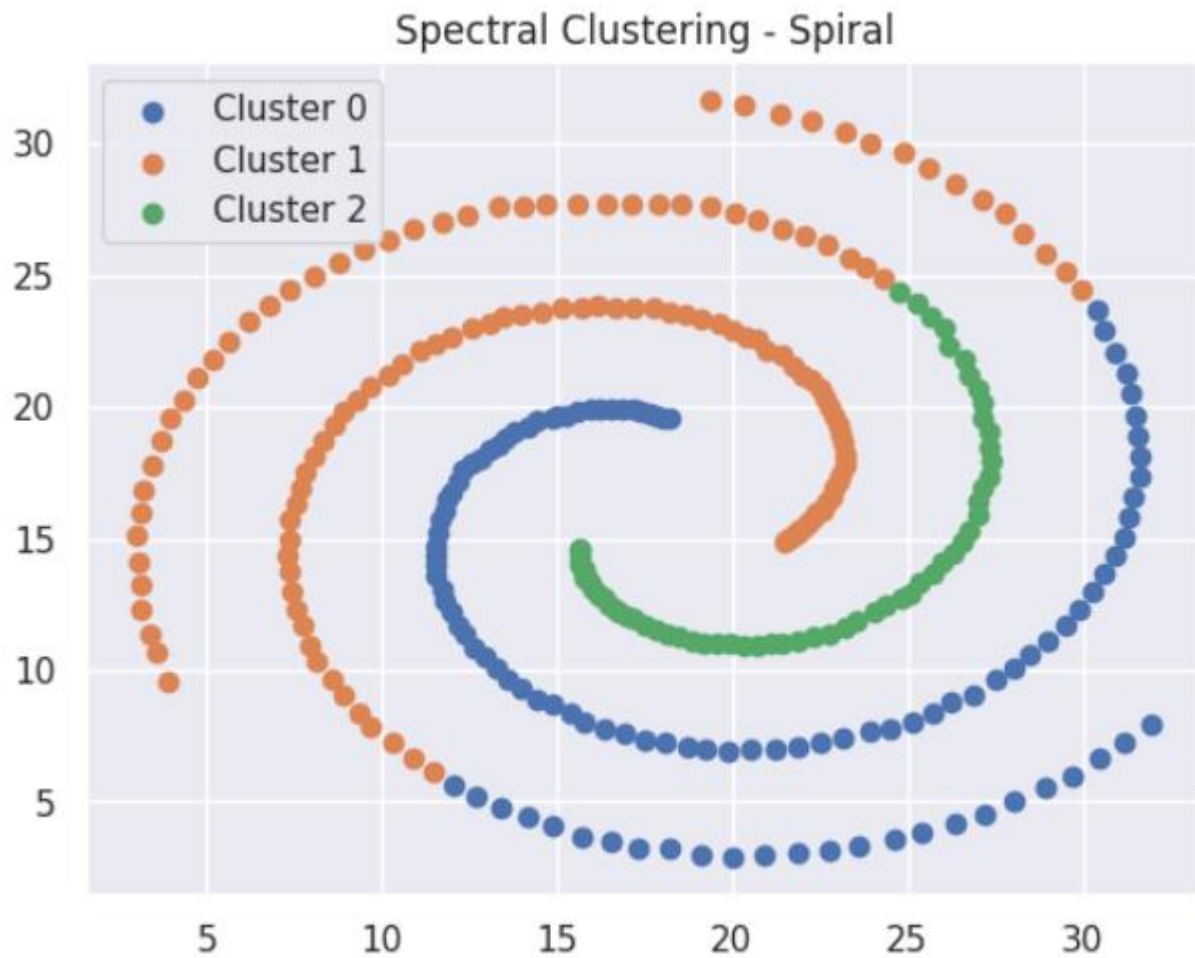


Agglomerative clustering performs significantly better than K-Means. The clusters align more closely with the spiral structure. The Ward linkage minimizes intra-cluster variance, capturing the non-linear separability to some extent. Hence, Agglomerative clustering is better suited for the spiral dataset than K-Means, but it is not perfect due to some overlapping areas.



Birch clustering exhibits similar performance to agglomerative clustering. It captures the spiral structure but has some misclassifications at the intersections of clusters. Birch is efficient for larger datasets, but its tree-based approach may struggle with complex non-linear structures like the spiral.

Birch clustering performs well but may require fine-tuning for complex geometries.



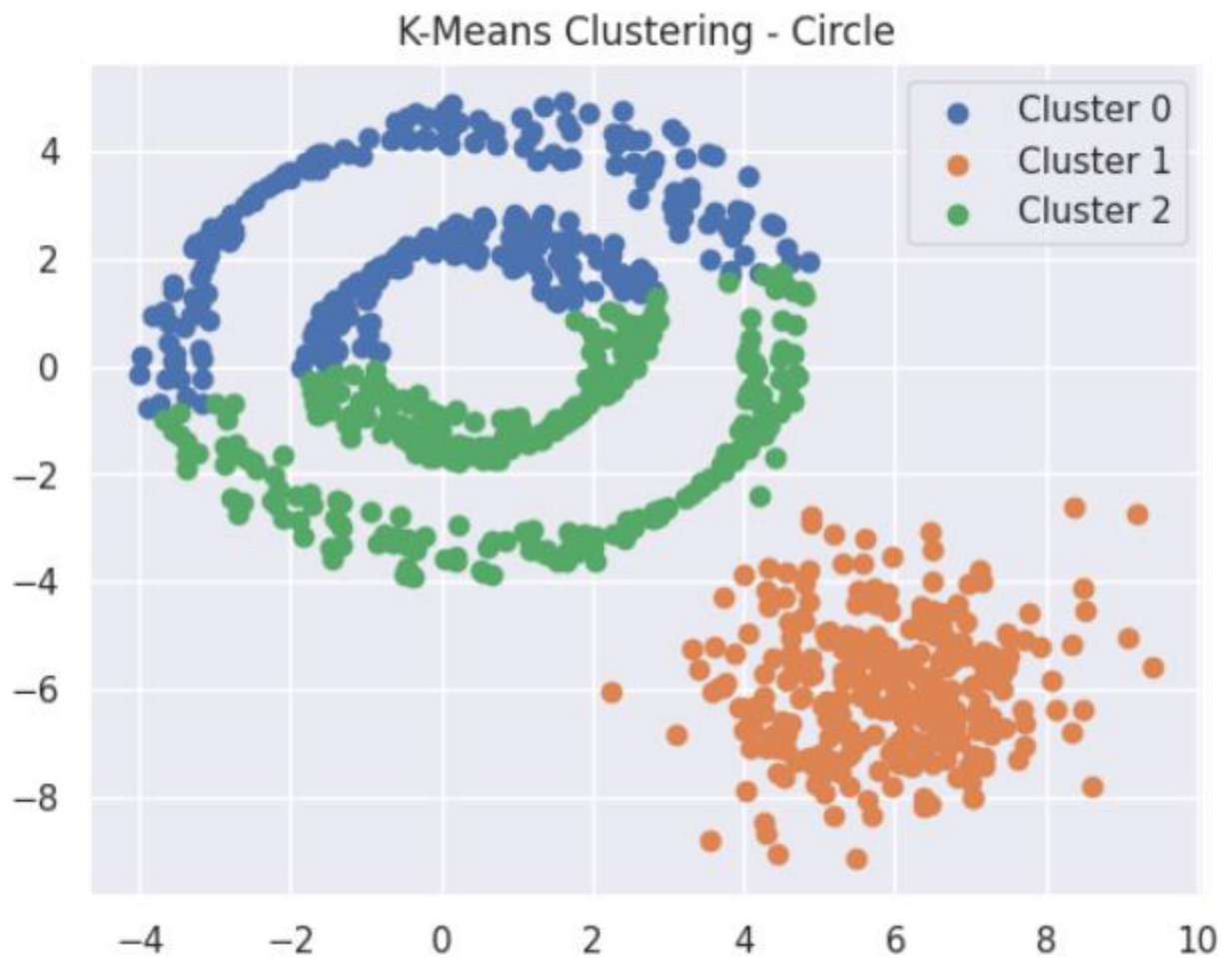
Spectral clustering performs the best for the spiral dataset.

The clusters are cleanly separated along the spiral arms, accurately reflecting the underlying structure.

By constructing a similarity graph and using eigenvalues, spectral clustering effectively handles non-linear separability.

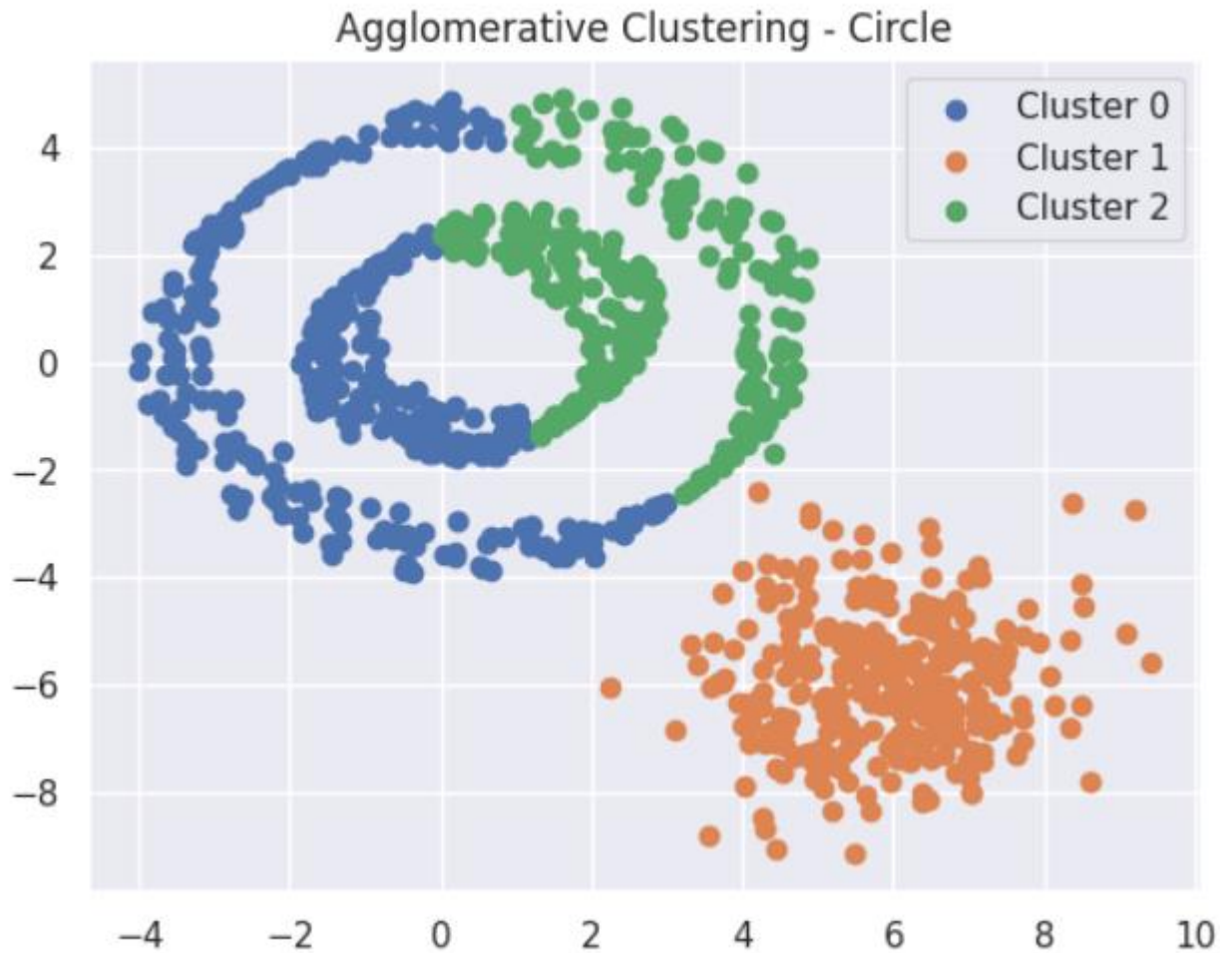
Spectral clustering is the most appropriate method for the spiral dataset, as it accurately identifies the non-linear structure.

Results for circle dataset:



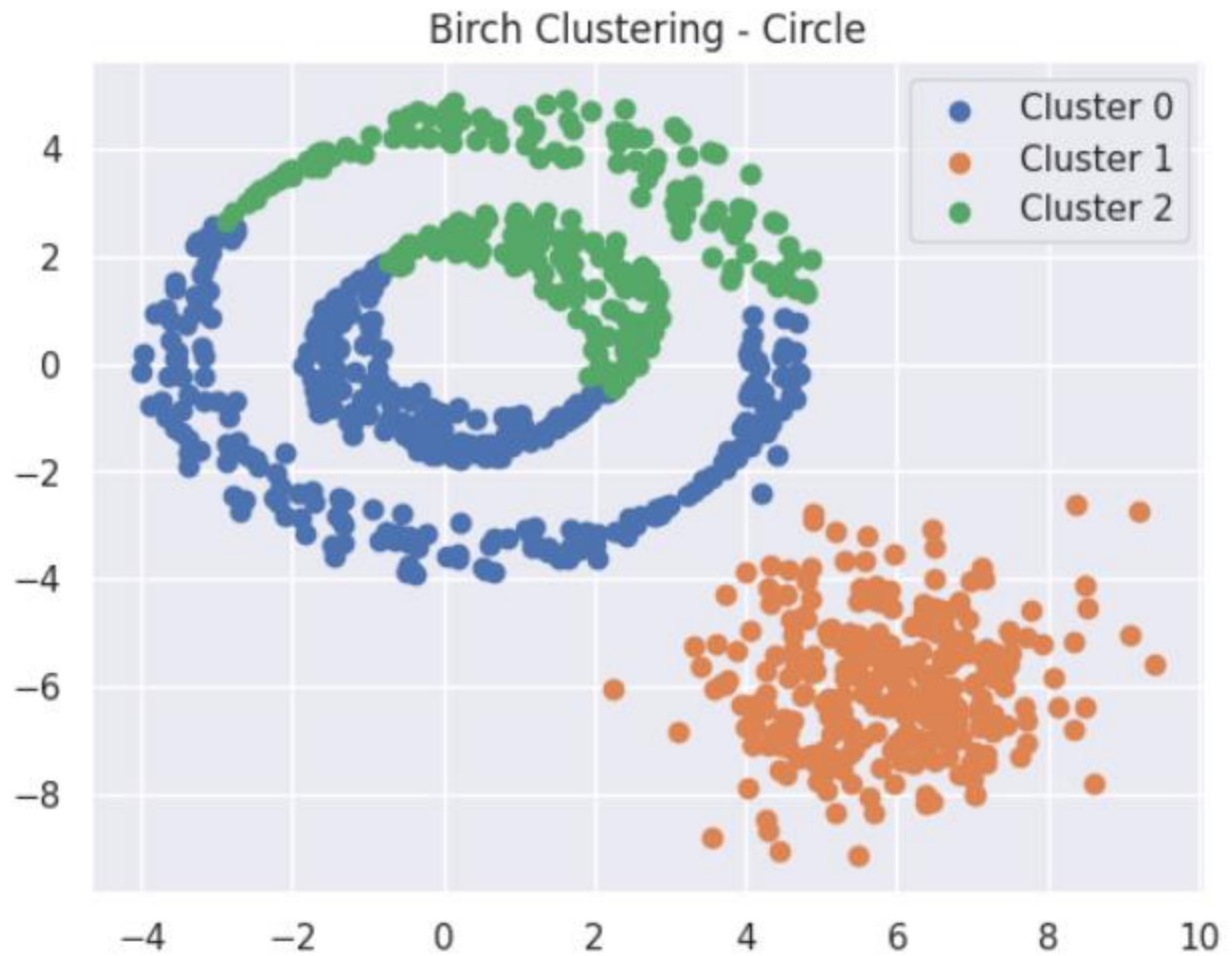
K-Means fails to capture the concentric circular structure. It divides the data into spherical clusters, cutting across the circular rings and assigning incorrect labels. This behavior arises from K-Means' assumption of isotropic cluster shapes.

K-Means is unsuitable for the circle dataset, as it cannot handle non-linear boundaries.



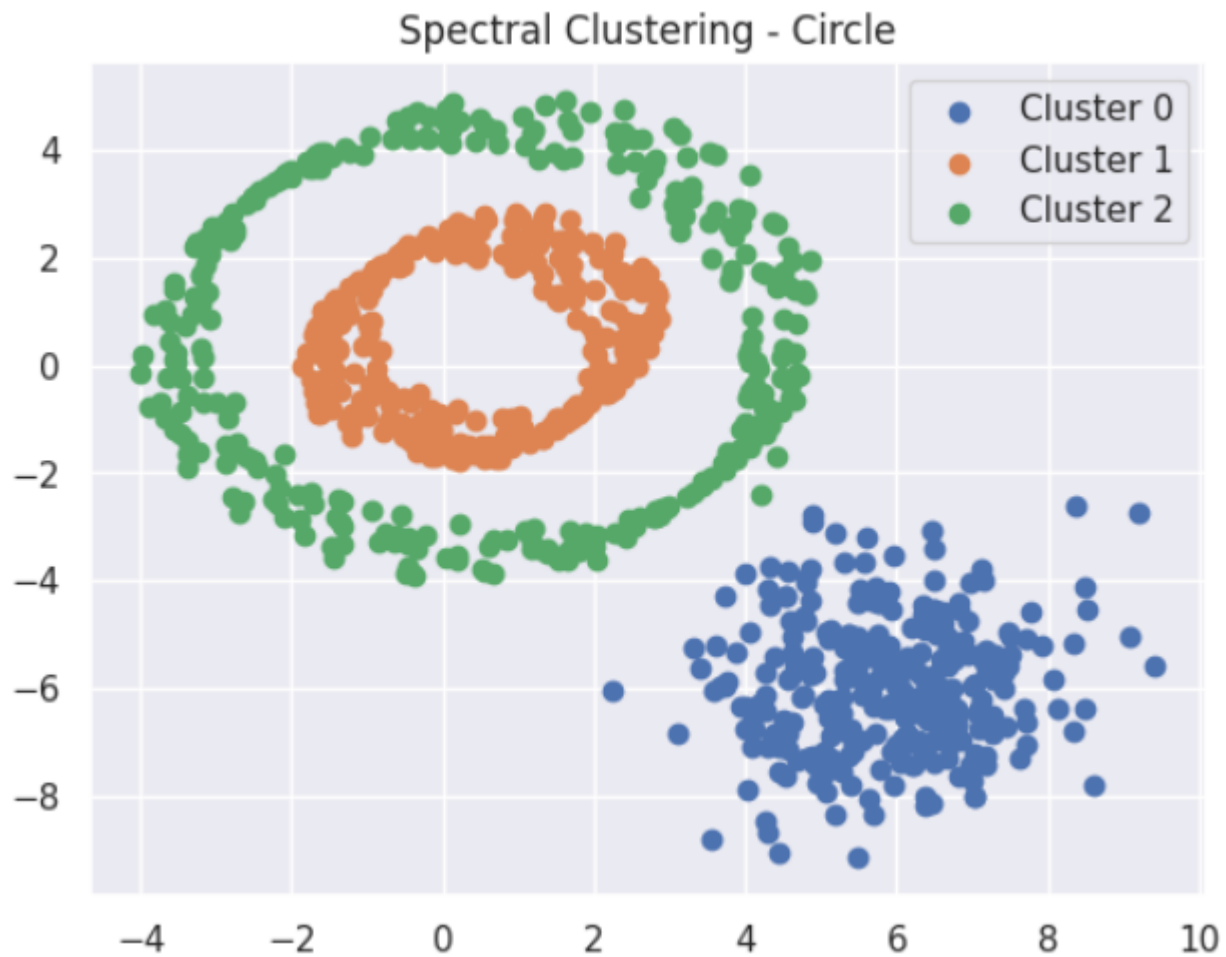
Agglomerative clustering performs better than K-Means but still struggles to separate the inner and outer rings cleanly. It captures the general structure of the dataset but shows some overlap between clusters.

Agglomerative clustering is better than K-Means but still falls short of accurately capturing the concentric structure.



Birch clustering produces results similar to agglomerative clustering. The clusters are reasonably well-separated but overlap in certain areas, especially between the inner and outer rings.

Birch clustering is efficient but less effective for datasets with concentric circular structures.



Spectral clustering performs the best for the circle dataset.

The clusters are cleanly separated along the inner and outer rings, accurately reflecting the circular structure. By leveraging graph-based techniques, spectral clustering excels at capturing non-linear separability.

Spectral clustering is the most suitable method for the circle dataset, accurately identifying the underlying structure.