

Fraud Analytics (CS6890)

Assignment:

1

Title : Identify clusters using (Node2Vec Embedding, Spectral, and GCN) embeddings

	Name	Roll Number
	<i>Shreesh Gupta</i>	CS23MTECH12009
	<i>Hrishikesh Hemke</i>	CS23MTECH14003
Team Details :		
	<i>Manan Patel</i>	CS23MTECH14006
	<i>Yug Patel</i>	CS23MTECH14019
	<i>Bhargav Patel</i>	CS23MTECH11026

In []: 1 # Required libraries
2 !pip install pandas numpy matplotlib scikit-learn torch torch-geometric ne

In [1]:

```

1 import pandas as pd
2 import numpy as np
3 from matplotlib import cm as cm
4 import matplotlib.pyplot as plt
5 from sklearn.cluster import DBSCAN
6 from sklearn.cluster import k_means
7 from sklearn.decomposition import PCA
8 import torch
9 import torch.nn.functional as F
10 from torch_geometric.datasets import Planetoid
11 from torch_geometric.nn import GCNConv
12 from torch_geometric.transforms import NormalizeFeatures
13 from torch_geometric.data import Data
14 import networkx as nx
15 import matplotlib.pyplot as plt
16 from sklearn.manifold import TSNE
17 import random
18 from gensim.models import Word2Vec
19 from joblib import Parallel, delayed
20 from node2vec import Node2Vec
21 from collections import Counter
22 from sklearn.preprocessing import StandardScaler
23 from sklearn.metrics import silhouette_score, calinski_harabasz_score, davies_bouldin_score
24 from sklearn.cluster import KMeans
25 from sklearn import datasets
26 from scipy.cluster.hierarchy import dendrogram, linkage
27 from scipy.cluster.hierarchy import fcluster
28 from torch_geometric.utils import (negative_sampling, remove_self_loops, add_self_loops)
29 from torch.optim import Adam
30 from numpy.linalg import eigh # Used for symmetric matrices, like the Laplacian matrix

```

Data Analysis

- This code snippet loads a dataset from a CSV file named 'Payments - Payments.csv' into a DataFrame and then constructs a directed graph (DiGraph) using NetworkX, where Sender and Receiver columns represent the nodes and the Amount column

In [2]:

```

1 df = pd.read_csv('Payments - Payments.csv')
2 # Display the first few rows of the dataset to understand its structure
3 df.head()

```

Out[2]:

	Sender	Receiver	Amount
0	1309	1011	123051
1	1309	1011	118406
2	1309	1011	112456
3	1309	1011	120593
4	1309	1011	166396

```
In [3]: 1 G = nx.from_pandas_edgelist(df, 'Sender', 'Receiver', ['Amount'], create_
2 node_to_index = {node: idx for idx, node in enumerate(G.nodes())}
3
4 # Print basic information about the graph to ensure it's constructed corre-
5 print('Number of nodes', len(G.nodes))
6 print('Number of edges', len(G.edges))
7 print('Average degree', sum(dict(G.degree).values()) / len(G.nodes))
8
```

Number of nodes 799
 Number of edges 5358
 Average degree 13.411764705882353

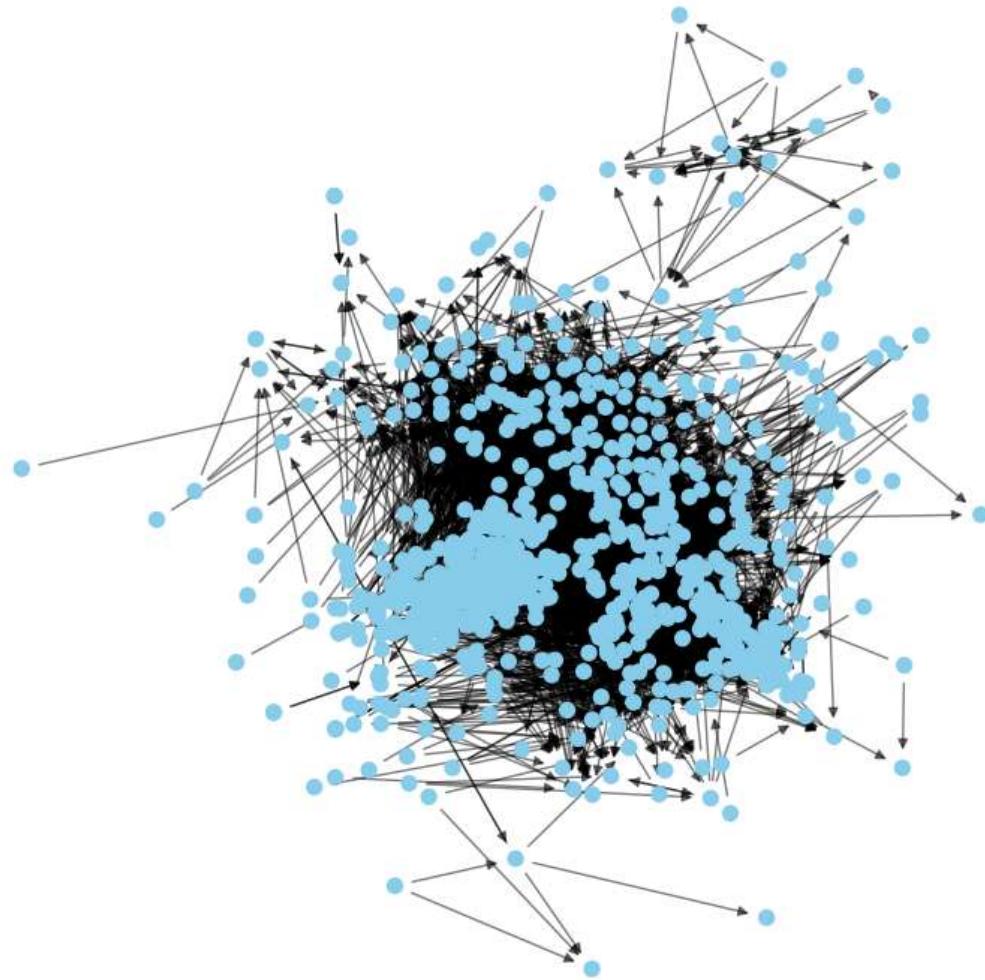
Plots

- This code snippet defines a color in RGB format and adjusts it for use in matplotlib by scaling to a 0-1 range. It includes a function `plot_graph` that plots a network graph using NetworkX, where all nodes are uniformly colored using the defined RGB color. The graph's layout is determined using the Kamada-Kawai layout algorithm

```
In [4]: 1 # Define a color, this will be used for the nodes in the graph.
2 # Here we're using a light blue color in RGB format, divided by 255 to use
3 CLR = (135, 206, 235)
4 CLR = [x/255 for x in CLR]
5
6 # A simple function to plot the graph using NetworkX
7 def plot_graph(graph):
8     """A simple function to plot the graph using NetworkX."""
9     # Set the color for each node. Here, all nodes will have the same color
10    node_colors = [CLR] * len(graph.nodes())
11
12    # Set the figure size
13    plt.figure(figsize=(10, 10))
14
15    # Define the layout for the graph. You can experiment with different lay-
16    pos = nx.kamada_kawai_layout(graph)
17
18    # Draw the networkx graph -- nodes, edges, and labels.
19    nx.draw_networkx_edges(graph, pos, alpha=0.5)
20    nx.draw_networkx_nodes(graph, pos, node_color=node_colors, cmap='autum-
21    # nx.draw_networkx_labels(graph, pos, font_size=8, font_color='black')
22
23    # Turn off the axis, as they are not meaningful for this kind of plot.
24    plt.axis('off')
25
26    # Display the plot.
27    plt.show()
```

```
In [5]: 1 plot_graph(G)
```

```
c:\Users\gupta\anaconda3\Lib\site-packages\networkx\drawing\nx_pylab.py:433:  
UserWarning: No data for colormapping provided via 'c'. Parameters 'cmap'  
will be ignored  
    node_collection = ax.scatter(
```



```
In [6]: 1 def plot_clusters(embeddings, labels):
2     """
3         Function to plot the clusters using embeddings and labels from DBSCAN.
4
5         :param embeddings: The PCA-reduced embeddings of the data.
6         :param labels: The cluster labels from DBSCAN.
7         """
8
9         # Create a scatter plot of the embeddings colored by Labels
10        unique_labels = set(labels)
11        colors = plt.cm.jet(np.linspace(0, 1, len(unique_labels)))
12
13        plt.figure(figsize=(12, 8))
14        for label, col in zip(unique_labels, colors):
15            if label == -1:
16                # Use black color for noise (label == -1)
17                col = 'k'
18                marker = 'x' # Use a different marker for noise points
19                label_class = 'Noise'
20            else:
21                marker = 'o'
22                label_class = f'Cluster {label}'
23
24            class_member_mask = (labels == label)
25            xy = embeddings[class_member_mask]
26            plt.plot(xy[:, 0], xy[:, 1], marker, markerfacecolor=col,
27                      markeredgecolor='k', markersize=10, label=label_class)
28
29        plt.title('Clustering')
30        plt.xlabel('PCA Component 1')
31        plt.ylabel('PCA Component 2')
32        plt.legend()
33        plt.grid(True)
34        plt.show()
```

```
In [7]: 1 def plot_tsne(data, labels, title='t-SNE plot of Clusters'):
2     """
3         This function applies t-SNE to reduce data dimensions and plots the resulting clusters.
4
5         Parameters:
6             - data: ndarray, the high-dimensional data to be reduced and plotted.
7             - labels: array, cluster labels for each point in `data`.
8             - title: str, the title of the plot.
9     """
10    # Applying t-SNE to reduce dimension to 2
11    tsne = TSNE(n_components=2, random_state=42)
12    tsne_result = tsne.fit_transform(data)
13
14    # Setting up the plot
15    plt.figure(figsize=(10, 8))
16    scatter = plt.scatter(tsne_result[:, 0], tsne_result[:, 1], c=labels,
17
18    # Adding color bar and labels
19    plt.colorbar(scatter)
20    plt.title(title)
21    plt.xlabel('t-SNE Axis 1')
22    plt.ylabel('t-SNE Axis 2')
23    plt.grid(True)
24
25    # Show plot
26    plt.show()
27
```

```
In [8]: 1 def plot_all_pca_components(X_, labels, n_components=5):
2     """
3         Plot all the principal component pairs against each other tightly.
4
5         Parameters:
6             - X_pca: The transformed data from PCA.
7             - labels: The labels for the data points.
8             - n_components: The number of components used in PCA.
9     """
10
11    pca = PCA(n_components)
12    X_pca = pca.fit_transform(X_)
13
14    # Set up the matplotlib figure and axes, based on the number of components
15    fig, axes = plt.subplots(n_components, n_components, figsize=(15, 15))
16
17    # Color map
18    cmap = plt.get_cmap('viridis', np.unique(labels).max() + 1)
19
20    # Plot the scatter plots for each combination of components
21    for i in range(n_components):
22        for j in range(n_components):
23            ax = axes[i, j]
24
25            scatter = ax.scatter(X_pca[:, i], X_pca[:, j], c=labels, cmap=cmap)
26            ax.set_xlabel(f'Component {i+1}', fontsize=9)
27            ax.set_ylabel(f'Component {j+1}', fontsize=9)
28
29
30    # Color bar for labels
31    # cbar = plt.colorbar(scatter, ax=axes[:, -1], location='right', shrink=0.5)
32    # cbar.set_label('Labels', rotation=270, labelpad=15)
33
34
35    # Set a main title for all subplots
36    plt.suptitle('PCA Components vs. Each Other', fontsize=16)
37
38    # Adjust Layout to make it tight
39    plt.tight_layout(rect=[0, 0.03, 1, 0.95])
40
41    # Show the plot
42    plt.show()
43
44
```

Node2Vec

This Python implementation encapsulates the Node2Vec algorithm within a class structure. Node2Vec is designed for generating node embeddings that capture the network's topology while also considering a node's neighborhood's structure. The algorithm leverages the flexibility of random walks to explore diverse neighborhoods and uses these walks to train a Word2Vec model.

Key Parameters:

- **dimensions** : The number of dimensions of the embedding space.
- **walk_length** : The number of steps in each random walk.
- **num_walks** : The number of walks to start at each node.
- **p** and **q** : Parameters which control the random walk process. **p** (return parameter) controls the likelihood of immediately revisiting a node in the walk, while **q** (in-out parameter) differentiates between inward and outward nodes in terms of walk transitions.

Transition Probabilities:

The transition probabilities between nodes during the walks are defined as follows:

$$P(c_i = x | c_{i-1} = v) = \begin{cases} \frac{\pi_{vx}}{Z} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases}$$

Where:

- (π_{vx}) represents the unnormalized transition probability between nodes (v) and (x),
- (Z) is a normalizing constant ensuring probabilities sum to 1 over all choices for (x).

The probabilities (π_{vx}) are computed based on the edge weights and the parameters (p) and (q):

$$\pi_{vx} = \alpha_{pq}(t, x) \cdot w_{vx}$$

Where (w_{vx}) is the weight of the edge from (v) to (x), and ($\alpha_{pq}(t, x)$) is given by:

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } x = t \\ 1 & \text{if } x \neq t \text{ and } x \in N(t) \\ \frac{1}{q} & \text{if } x \notin N(t) \end{cases}$$

Random Walk Simulation:

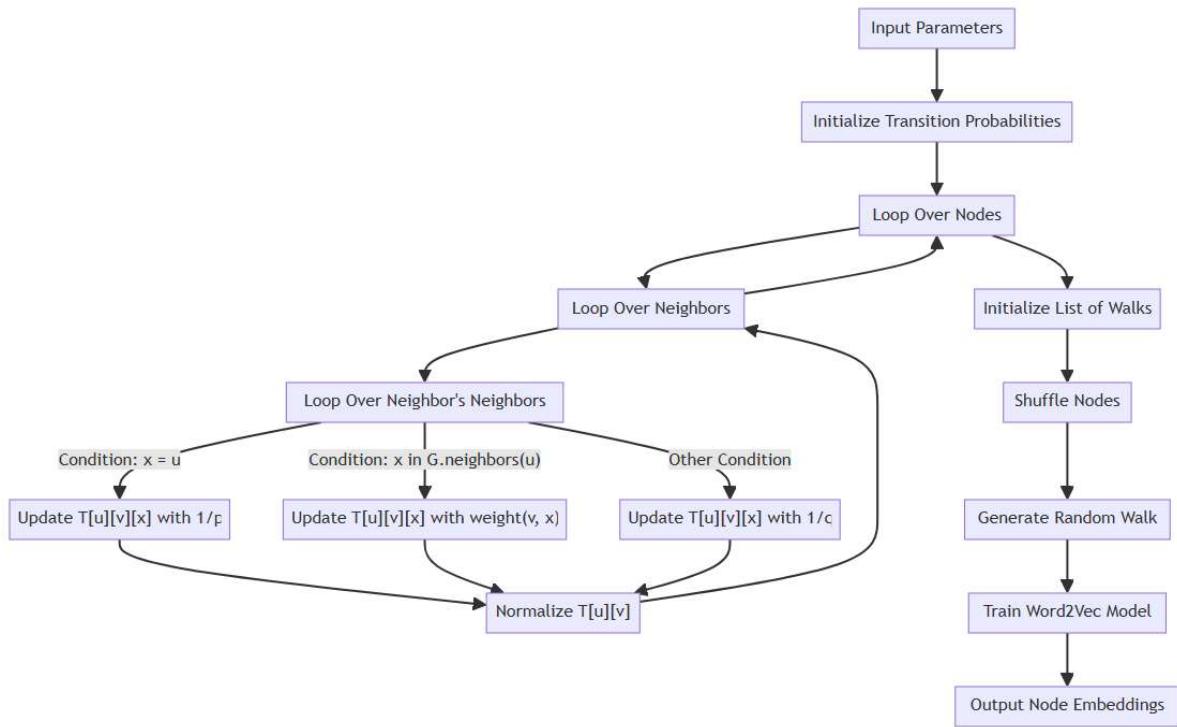
Random walks are simulated based on the precomputed transition probabilities. Each walk begins at a selected node and sequentially transitions to a neighboring node based on the probabilities, creating paths that reflect the local network structure.

Model Training:

After simulating the walks, they are used to train a Word2Vec model. This model learns vector representations for each node that are useful for machine learning tasks on graphs such as node classification or clustering.

Implementation Details:

The class `Node2Vec` includes methods to precompute transition probabilities, simulate random walks, and fit the Word2Vec model. The embeddings produced encapsulate both local and global structural information about the nodes, potentially improving performance on downstream tasks.



In [9]:

```
1  class Node2Vec:
2      def __init__(self, graph, dimensions=64, walk_length=30, num_walks=10,
3          self.graph = graph
4          self.dimensions = dimensions
5          self.walk_length = walk_length
6          self.num_walks = num_walks
7          self.p = p
8          self.q = q
9          self.workers = workers
10         self.transition_probs = {}
11
12         self.precompute_transition_probs()
13
14     def precompute_transition_probs(self):
15         """
16             Precompute transition probabilities for each node based on the graph
17         """
18         for node in self.graph.nodes():
19             self.transition_probs[node] = {}
20             for neighbor in self.graph.neighbors(node):
21                 self.transition_probs[node][neighbor] = {}
22                 for destination in self.graph.neighbors(neighbor):
23                     if destination == node:
24                         prob = self.graph[neighbor][destination].get('weight')
25                     elif destination in self.graph[node]:
26                         prob = self.graph[neighbor][destination].get('weight')
27                     else:
28                         prob = self.graph[neighbor][destination].get('weight')
29                     self.transition_probs[node][neighbor][destination] = prob
30
31         # Normalize the probabilities
32         total_prob = sum(self.transition_probs[node][neighbor].values()
33         for destination in self.transition_probs[node][neighbor])
34         self.transition_probs[node][neighbor][destination] /= total_prob
35
36     def _walk(self, start_node):
37         """
38             Simulate a single random walk of fixed length from a starting node
39         """
40         walk = [start_node]
41         while len(walk) < self.walk_length:
42             cur = walk[-1]
43             if len(self.graph[cur]) > 0: # Continue walk if there are neighbors
44                 if len(walk) == 1: # First step is special
45                     walk.append(random.choices(
46                         list(self.graph[cur]),
47                         weights=[self.graph[cur][n].get('weight', 1) for n in
48                         k=1
49                         )[0]))
50                 else:
51                     prev = walk[-2]
52                     next_node = random.choices(
53                         list(self.transition_probs[prev][cur].keys()),
54                         weights=list(self.transition_probs[prev][cur].values())
55                         k=1
56                         )[0]
57                     walk.append(next_node)
```

```

58         else:
59             break
60     return walk
61
62     def simulate_walks(self):
63         """
64             Repeatedly simulate random walks from each node.
65         """
66         walks = []
67         nodes = list(self.graph.nodes())
68         for _ in range(self.num_walks):
69             random.shuffle(nodes)
70             for node in nodes:
71                 walks.append(self._walk(node))
72     return walks
73
74     def fit(self, window=10, min_count=1, batch_words=4):
75         """
76             Train Word2Vec model.
77         """
78         walks = self.simulate_walks()
79         walks = [[str(node) for node in walk] for walk in walks] # Convert
80         model = Word2Vec(sentences=walks, vector_size=self.dimensions, window=window, min_count=min_count, batch_size=batch_size, workers=workers, epochs=epochs, sample=sample)
81     return model
82
83
84
85     node2vec_m = Node2Vec(G, dimensions=64, walk_length=30, num_walks=500, p=p, q=q)
86     model = node2vec_m.fit()
87
88
89     embeddings_m = {node: model.wv[str(node)] for node in G.nodes}
90     print(len(embeddings_m))

```

799

In [10]:

```

1 # Node2Vec parameters - you might adjust these based on your dataset specification
2 dimensions = 64
3 walk_length = 30
4 num_walks = 500
5 workers = 1
6
7 node2vec = Node2Vec(G, dimensions=dimensions, walk_length=walk_length, num_walks=num_walks, workers=workers)
8 model = node2vec.fit(window=10, min_count=1, batch_words=4)
9

```

In [11]:

```

1 embeddings = {node: model.wv[str(node)] for node in G.nodes}
2 print(len(embeddings))

```

799

DBSCAN Clustering Hyperparameter Tuning

This Python script performs hyperparameter tuning to identify the optimal `eps` parameter for the DBSCAN clustering algorithm. The `eps` parameter determines the maximum distance between two samples for them to be considered as in the same neighborhood. The aim is to find the `eps` value that results in the highest number of clusters without considering outliers.

Preprocessing:

- Standardization:** The feature set `X` (presumed to be a collection of embeddings) is standardized using `StandardScaler`. This normalization ensures that each feature contributes equally to the distance computation, crucial for distance-based algorithms like DBSCAN.

$$X' = \frac{X - \mu}{\sigma}$$

Where (`X'`) is the scaled data, (`\mu`) is the mean, and (`\sigma`) is the standard deviation.

Hyperparameter Tuning:

- Exploring `eps` Values:** The `eps` values are explored over a range from 0.1 to 5, divided into 50 intervals. For each value of `eps`, the DBSCAN algorithm is executed to compute the number of clusters formed.
- DBSCAN Execution:**
 - DBSCAN is run with the current `eps` and a fixed `min_samples` of 3 (minimum points to form a cluster).
 - The number of clusters is determined by counting unique labels, excluding noise points (labeled as `-1`).

Visualization:

- Plotting:** A line plot is generated to visualize the relationship between `eps` values and the number of clusters formed. This helps in visually identifying the `eps` value that maximizes the number of meaningful clusters.

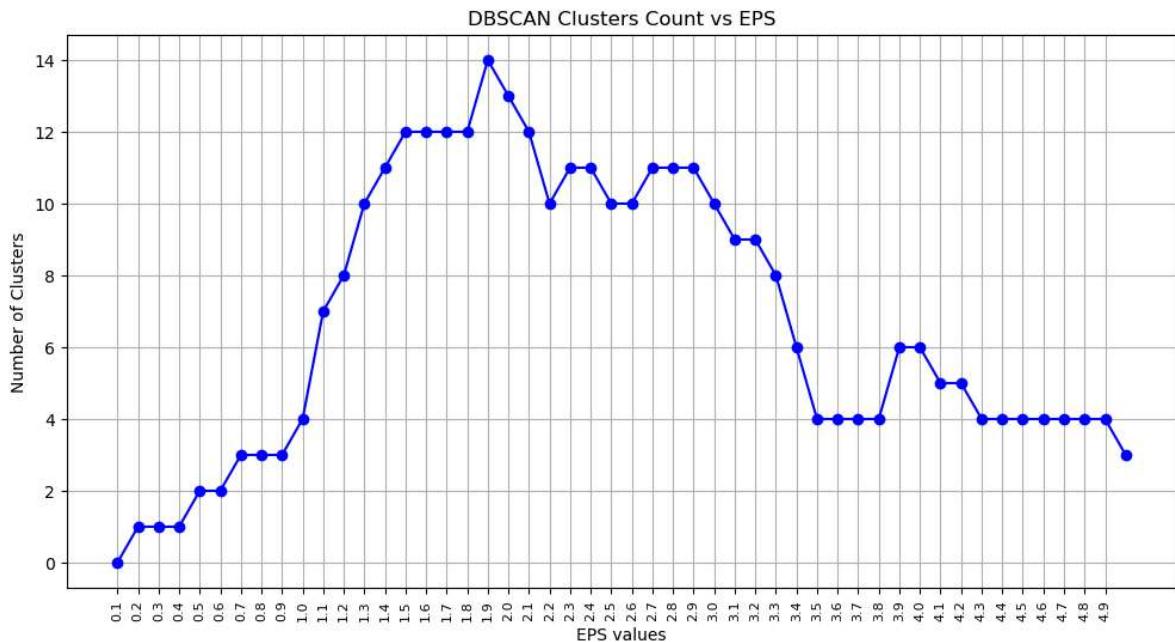
Results:

- Optimal Parameters:** The script outputs the best `eps` value and the corresponding number of clusters. This `eps` value is deemed optimal based on the criterion of maximizing the number of clusters.

By adjusting the `eps` parameter based on this analysis, the DBSCAN algorithm can be finely tuned to capture the natural clustering structure of the dataset more effectively.

In [12]:

```
1 # Hyperparameter tuning for EPS value
2 eps_range = np.linspace(0.1, 5, 50)
3 num_clusters = []
4 best_eps, best_cluster = 0.1, 0
5 X = list(embeddings_m.values())
6
7 scaler = StandardScaler()
8 embedding_scaled = scaler.fit_transform(X)
9 samples = 3
10 X = embedding_scaled
11
12 # Testing for max cluster EPS
13 for eps in eps_range:
14     dbscan = DBSCAN(eps=eps, min_samples=samples)
15     dbscan.fit(X)
16     num_cls = len(set(dbscan.labels_)) - (1 if -1 in dbscan.labels_ else 0)
17     num_clusters.append(num_cls)
18     if best_cluster < num_cls:
19         best_cluster = num_cls
20         best_eps = eps
21
22 # Plot EPS vs Number of Clusters
23 plt.figure(figsize=(12, 6))
24 plt.plot(eps_range, num_clusters, 'o-', color='blue')
25 plt.xticks(np.arange(min(eps_range), max(eps_range), 0.1))
26 plt.xticks(rotation=90)
27 plt.tick_params(axis='x', which='major', labelsize=8)
28 plt.xlabel('EPS values')
29 plt.ylabel('Number of Clusters')
30 plt.title('DBSCAN Clusters Count vs EPS')
31 plt.grid(True)
32 plt.show()
33
34 # Display the best EPS value found and the number of clusters associated
35 print(f"Best EPS based on hyperparameter tuning: {best_eps}")
36 print(f"Number of clusters for the best EPS: {best_cluster}")
37
```



Best EPS based on hyperparameter tuning: 1.9000000000000001

Number of clusters for the best EPS: 14

```
In [13]: 1 clusters = DBSCAN(eps=best_eps, min_samples=samples)
          2 clusters.fit(X)
          3 labels = clusters.labels_
          4
          5 num_cls = len(set(clusters.labels_)) - (1 if -1 in clusters.labels_ else 0)
          6 print(num_cls)
```

14

```
In [14]: 1 set(clusters.labels_)
```

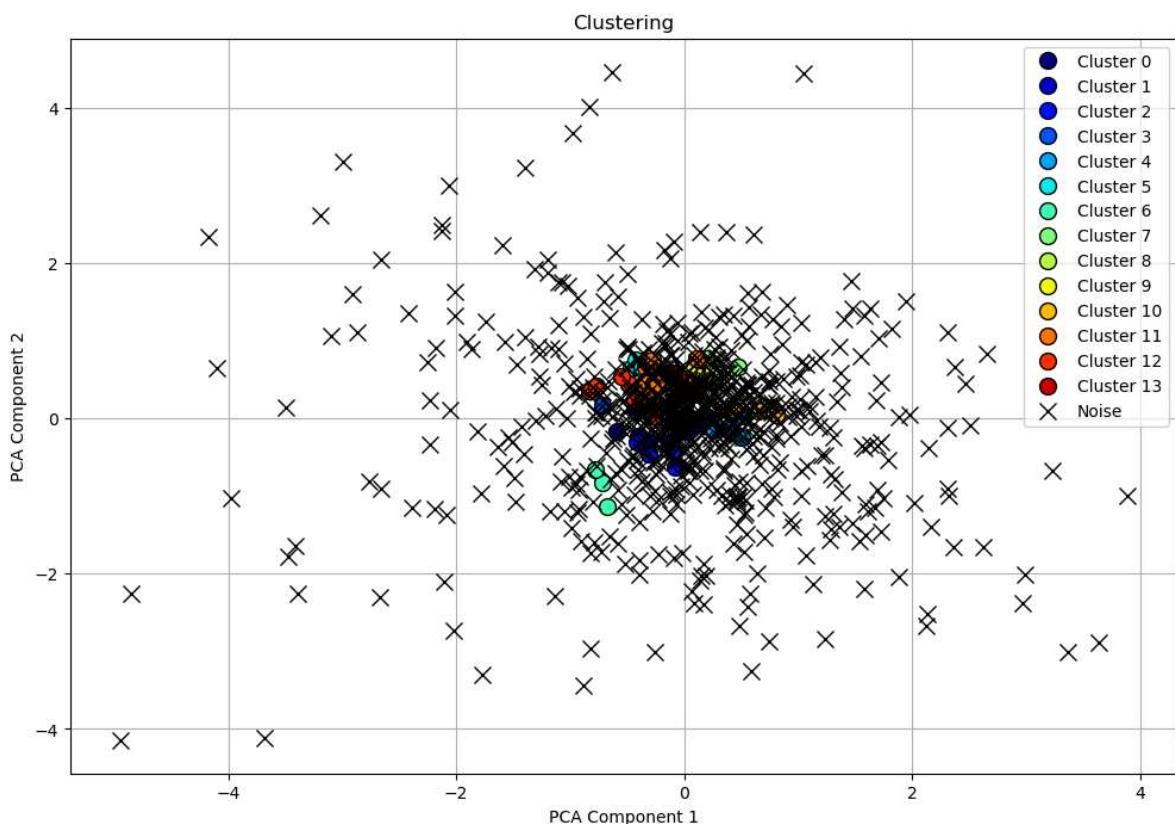
Out[14]: {-1, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13}

DBSCAN Results

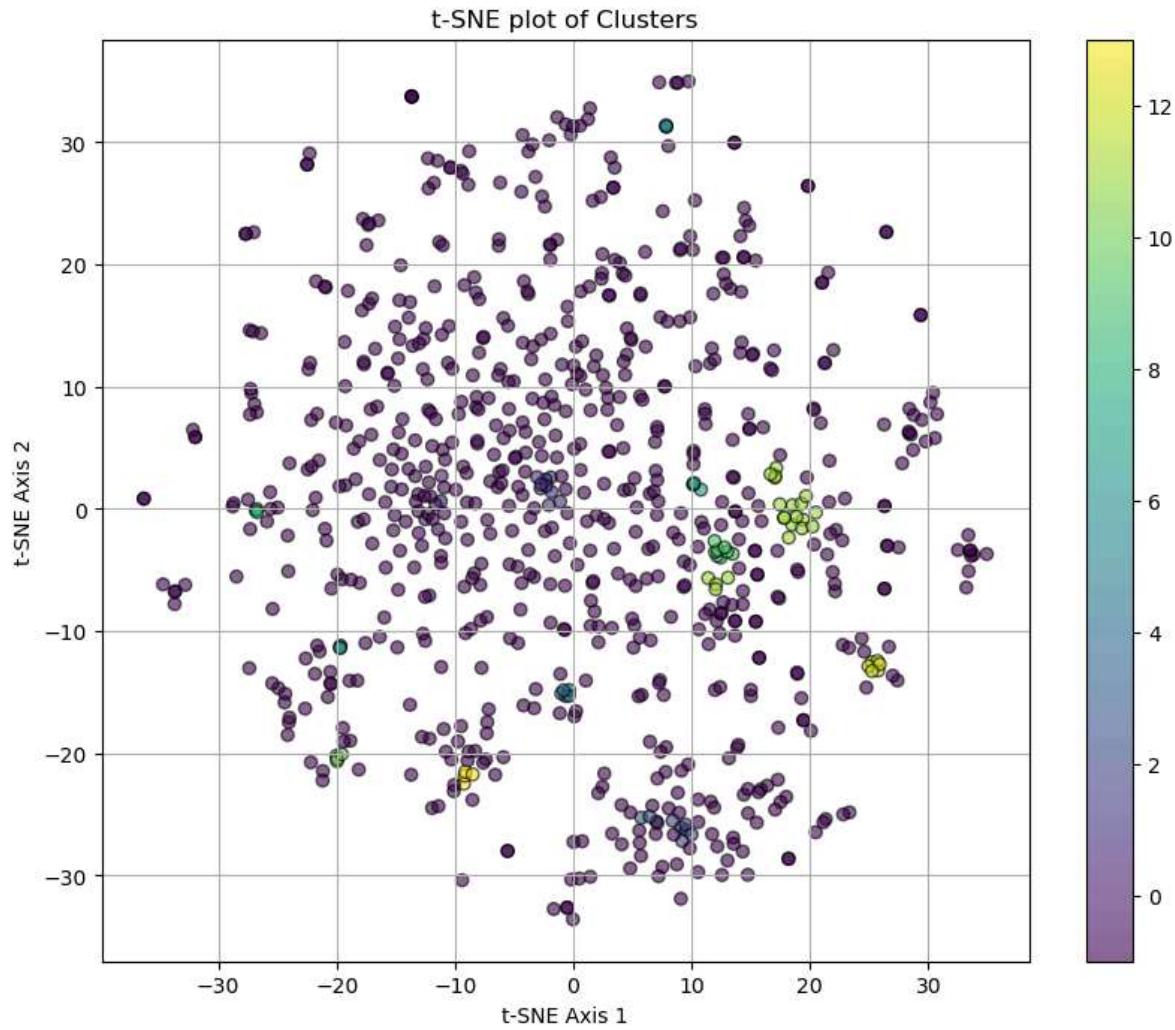
```
In [15]: 1 def count_labels(labels):
2     """
3         Function to count the occurrences of each label.
4
5         :param labels: A list or array of cluster labels.
6         :return: A dictionary with labels as keys and counts as values.
7     """
8     label_counts = Counter(labels)
9     return label_counts
10
11 # Assuming `clusters.labels_` is your list of Labels from DBSCAN
12 label_counts = count_labels(clusters.labels_)
13 print(label_counts)
14
```

Counter({-1: 715, 11: 21, 1: 9, 9: 8, 2: 6, 12: 6, 0: 5, 4: 5, 7: 4, 10: 4, 1
3: 4, 3: 3, 5: 3, 6: 3, 8: 3})

```
In [16]: 1 plot_clusters(X, clusters.labels_)
2
```



```
In [17]: 1 plot_tsne(X, clusters.labels_)
```



Clustering Evaluation Metrics

The effectiveness of the DBSCAN clustering algorithm is evaluated using three different metrics, each providing insights into different aspects of the clustering quality:

Silhouette Score

- **Range:** -1 to 1
- **Interpretation:** A higher score indicates that clusters are compact and well-separated compared to other clusters. A score close to 1 denotes that clusters are dense and well-separated. A score near 0 indicates overlapping clusters, and a negative score suggests that samples might have been assigned to the wrong clusters.

Calinski-Harabasz Index

- **Interpretation:** A higher score is better. It signifies that the clusters are dense and well-separated. This index is particularly useful for datasets with clusters of roughly equal size.

Davies-Bouldin Index

- **Range:** 0 to Infinity
- **Interpretation:** A lower index indicates better clustering. This index evaluates how much the clusters are separated and how compact they are, with lower values showing better clustering performance.

These metrics collectively provide a comprehensive view of the clustering performance, helping to understand the strengths and weaknesses of the applied clustering approach. The scores are computed for the DBSCAN algorithm's results, reflecting how well it has identified dense and distinct clusters within the dataset.

In [18]:

```

1 # Assuming 'X' is your dataset and 'clusters.labels_' are the labels from
2 def matices(X, labels):
3     silhouette_avg = silhouette_score(X, labels)
4     ch_score = calinski_harabasz_score(X, labels)
5     db_index = davies_bouldin_score(X, labels)
6     metrics = pd.DataFrame({
7         'Metric': ['Silhouette Score', 'Calinski-Harabasz Index', 'Davies-
8         'Score': [silhouette_avg, ch_score, db_index]
9     })
10
11     return metrics
12     # print(table.strip()) # Using strip() to remove any Leading/trailing
13
14 matices(X,clusters.labels_)
```

Out[18]:

	Metric	Score
0	Silhouette Score	-0.319318
1	Calinski-Harabasz Index	1.079690
2	Davies-Bouldin Index	3.649824

K Means

This code snippet defines a function `KMeansClassFinder` that is designed to determine the optimal number of clusters for a dataset using the Elbow Method, a popular heuristic in KMeans clustering. The function takes two parameters: `X`, which is the dataset to be clustered, and `n_classes`, which defines the range for the number of clusters to test (from 1 to `n_classes` - 1).

Within the function:

- **WCSS Calculation:** It first initializes an empty list `wcss` to store the within-cluster sum of squares (WCSS). WCSS is a measure of the compactness of the clusters, with lower values generally indicating better clustering.
- **KMeans Iteration:** The function then iterates over the specified range of cluster numbers. For each possible number of clusters, it initializes a KMeans clustering with parameters `init='k-means++'` for efficient centroid placement, `max_iter=300` for a maximum number of iterations, `n_init=10` for multiple initial centroid seeds, and a fixed `random_state=0` for reproducibility. It then fits the KMeans algorithm to the dataset `X` and appends the inertia (WCSS) of the resulting model to the `wcss` list.

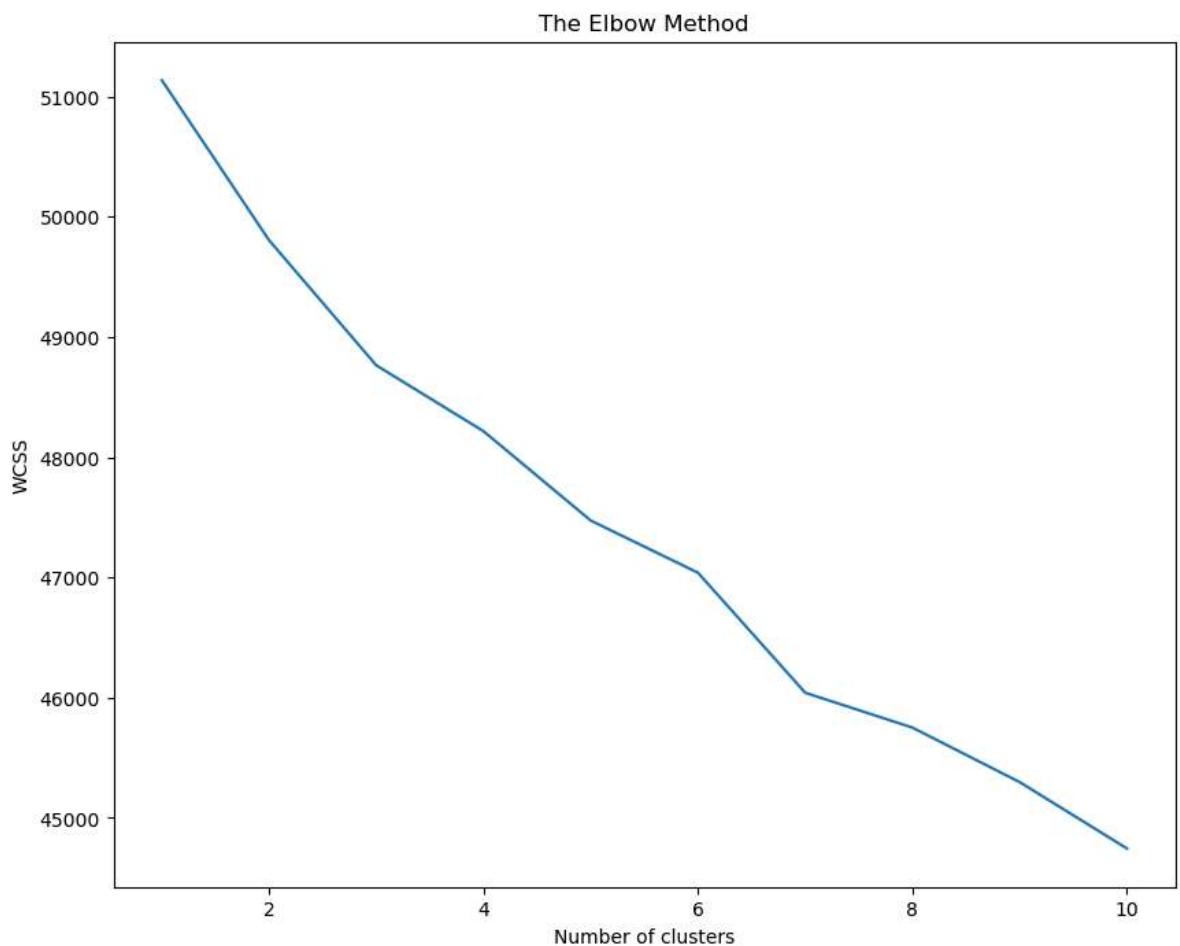
- **Elbow Plotting:** After calculating WCSS for each cluster number, the function plots these values against the number of clusters using matplotlib. The plot is sized 10x8 inches and is labeled appropriately, including a title 'The Elbow Method'. The elbow point in the plot, where the rate of decrease in WCSS sharply shifts, can indicate the optimal number of clusters.

The function is finally called with the dataset `X` and an upper limit of 11 for the number of clusters to be examined. This approach allows for a visual inspection of the WCSS to identify the best cluster count based on the Elbow Method.

In [19]:

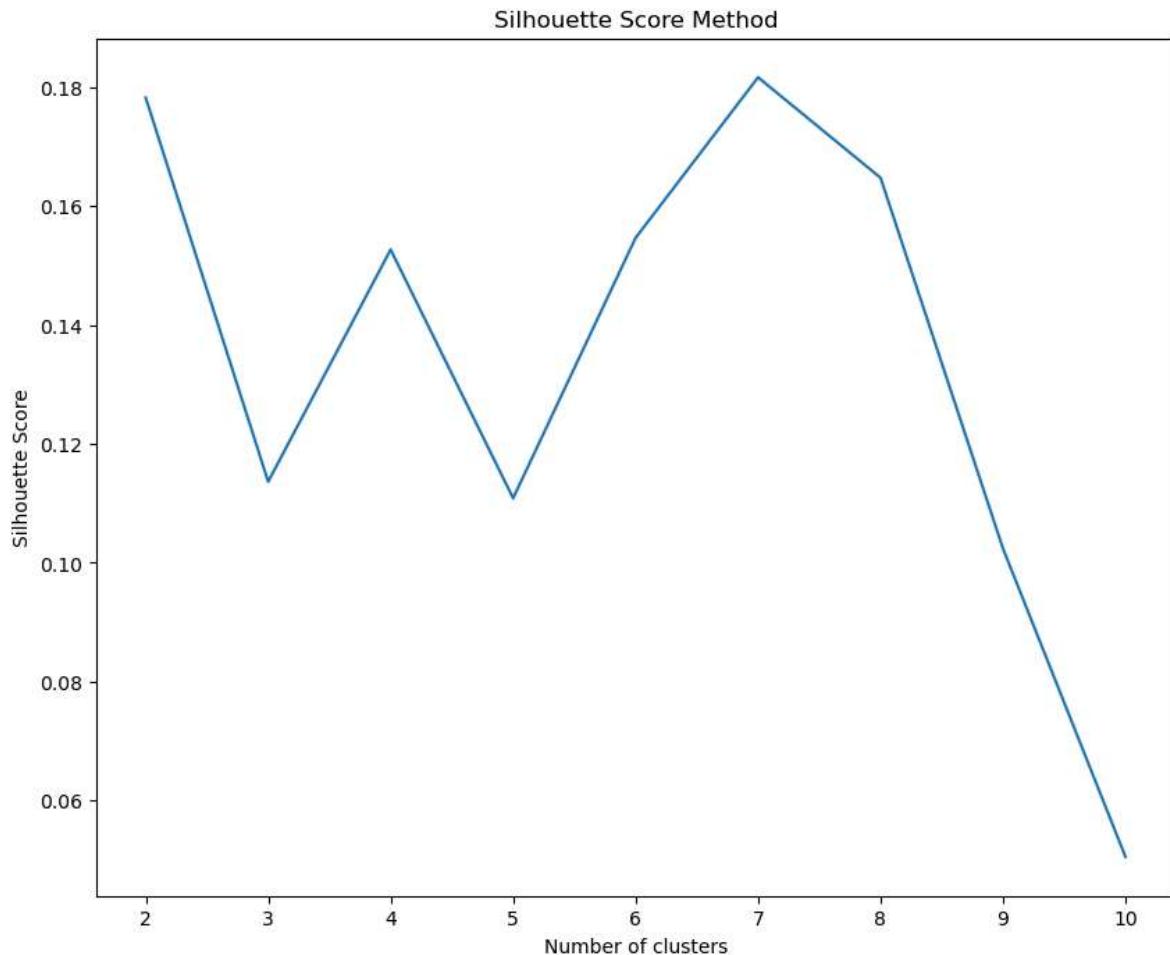
```
1 def KMeansClassFinder(X, n_classes):
2     # Calculating WCSS values for 1 to 10 clusters
3     wcss = []
4     for i in range(1, n_classes):
5         kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_in
6             kmeans.fit(X)
7             wcss.append(kmeans.inertia_)
8
9     # Plotting the results onto a Line graph to observe 'The Elbow'
10    plt.figure(figsize=(10,8))
11    plt.plot(range(1, 11), wcss)
12    plt.title('The Elbow Method')
13    plt.xlabel('Number of clusters')
14    plt.ylabel('WCSS') # Within cluster sum of squares
15    plt.show()
16
17 KMeansClassFinder(X, 11)
```

```
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
```



```
In [20]: 1 def KMeansClassFinder_silhouette_scores(X):
2     silhouette_scores = []
3     for i in range(2, 11): # Silhouette score can't be calculated for one
4         kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_ini
5         kmeans.fit(X)
6         score = silhouette_score(X, kmeans.labels_)
7         silhouette_scores.append(score)
8
9     # Plotting the results
10    plt.figure(figsize=(10,8))
11    plt.plot(range(2, 11), silhouette_scores)
12    plt.title('Silhouette Score Method')
13    plt.xlabel('Number of clusters')
14    plt.ylabel('Silhouette Score')
15    plt.show()
16
17 KMeansClassFinder_silhouette_scores(X)
```

```
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
```



In [21]:

```

1 # Assuming X is your dataset
2 def applyKMeans(X, num_clusters):
3     kmeans = KMeans(n_clusters=num_clusters)
4     k_clusters = kmeans.fit(X)
5     print(set(k_clusters.labels_))
6     return k_clusters
7
8 k_clusters = applyKMeans(X, 10)

```

```

c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1412: 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)
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
{0, 1, 2, 3, 4, 5, 6, 7, 8, 9}

```

In [22]:

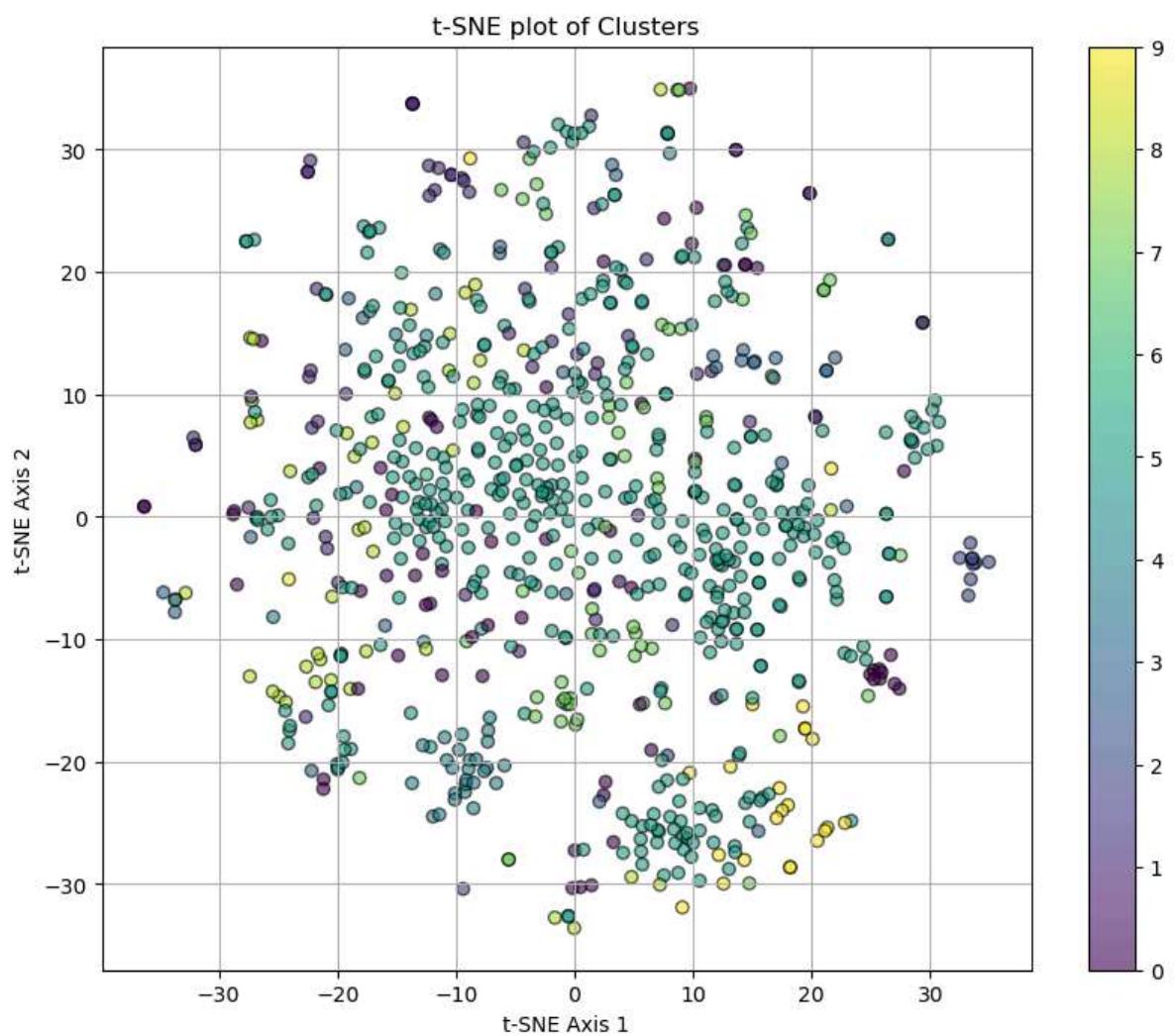
```

1 count_labels(k_clusters.labels_)
2

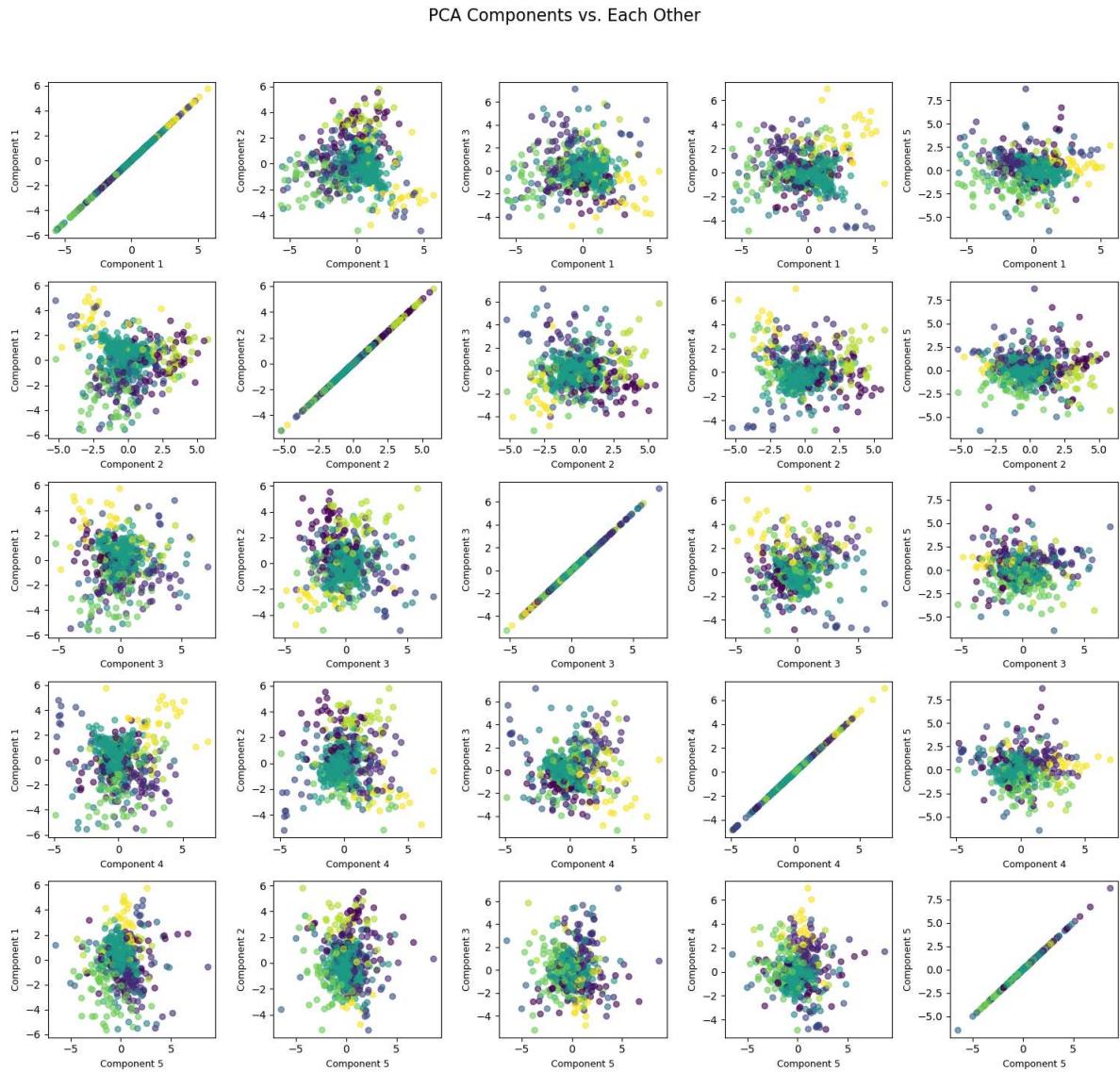
```

Out[22]: Counter({5: 458, 0: 68, 7: 63, 1: 61, 8: 43, 4: 38, 3: 30, 9: 24, 2: 11, 6: 3})

```
In [23]: 1 plot_tsne(X, k_clusters.labels_)
```



In [24]: 1 plot_all_pca_components(X, k_clusters.labels_, n_components=5)



Evaluation metrics

In [25]: 1 matices(X, k_clusters.labels_)

Out[25]:

	Metric	Score
0	Silhouette Score	0.150063
1	Calinski-Harabasz Index	12.485185
2	Davies-Bouldin Index	3.426045

Hierarchical Clustering and Dendrogram Visualization

This Python code snippet utilizes the `linkage` method from the `scipy.cluster.hierarchy` module to perform hierarchical clustering using the Ward linkage method. The Ward method minimizes the variance of the clusters being merged, making it well-suited for identifying

clusters that are relatively compact and equal in size.

Function `plot_dendrogram`

- **Parameters:**
 - X : A data matrix where samples are rows and features are columns.
- **Process:**
 1. **Linkage Computation:** The `linkage()` function computes hierarchical clusters using the 'ward' method, which is a common approach for hierarchical agglomerative clustering (HAC). This method is particularly effective for creating clusters by minimizing the total within-cluster variance. At each step, the pair of clusters with the minimum between-cluster distance are merged.
 2. **Dendrogram Plotting:** The resulting linkage array, which contains the hierarchical clustering information, is used to plot a dendrogram. The dendrogram illustrates how each cluster is composed by drawing a U-shaped link between a non-singleton cluster and its children.
- **Visualization Settings:**
 - `orientation='top'` : This places the root at the top, and items are displayed from top to bottom.
 - `distance_sort='descending'` : This sorts the distances between pairs of clusters in descending order, aiding in the visualization of the most distinct clusters.
 - `show_leaf_counts=True` : This option adds the number of items in each node, providing a clear view of the cluster sizes at a glance.
- **Plot Customization:** The plot is adjusted to a size of 10x7 inches to ensure all elements of the dendrogram are clearly visible, making it suitable for detailed analysis or presentations.

The function returns the linkage matrix, which can be used for further analysis or for plotting additional dendograms with different configurations.

Example Usage

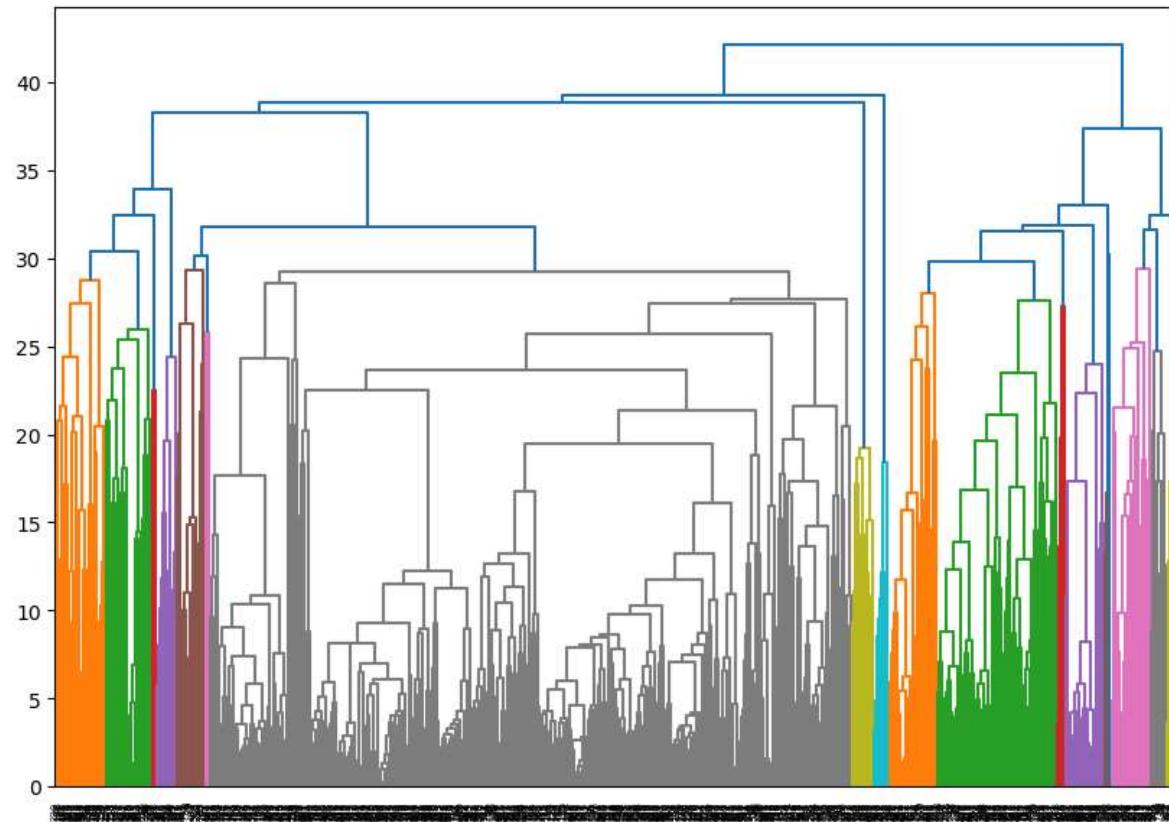
The function `plot_dendrogram` is called with the dataset `X`, and the linkage matrix `linked` is stored for potential further analysis. This methodology provides a profound insight into the data's hierarchical structure, which is crucial for understanding natural groupings within the dataset.

In [26]:

```

1 # Assuming X is your data matrix with samples as rows
2 # Use the 'ward' Linkage method to compute the clusters
3 def plot_dendrogram(X):
4     linked = linkage(X, 'ward')
5
6     # Plot the hierarchical clustering as a dendrogram.
7     plt.figure(figsize=(10, 7))
8     dendrogram(linked,
9                 orientation='top',
10                distance_sort='descending',
11                show_leaf_counts=True)
12    plt.show()
13    return linked
14 linked = plot_dendrogram(X)

```



In [27]:

```

1 def applyFcluster(linked, max_d):
2     fclusters = fcluster(linked, max_d, criterion='distance')
3     print(set(fclusters)) # Prints cluster labels for each point
4     return fclusters
5
6 fclusters = applyFcluster(linked, 35)

```

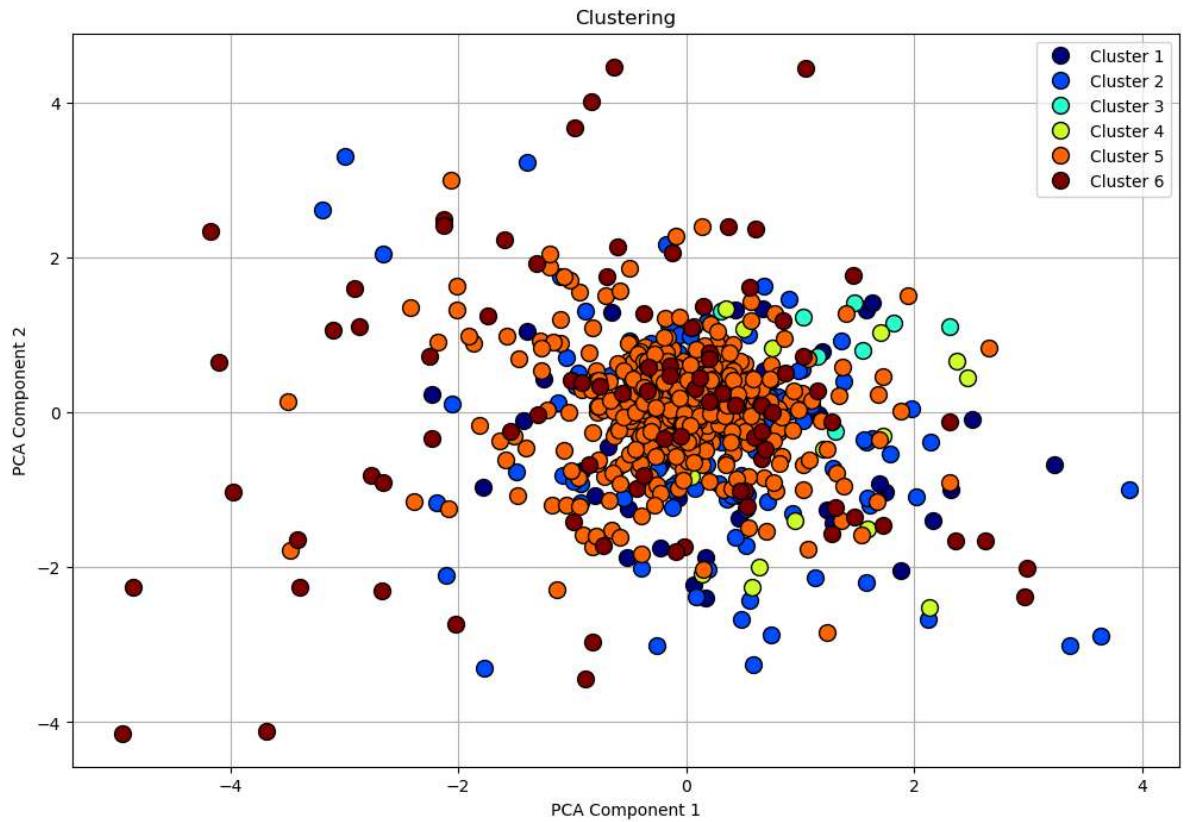
{1, 2, 3, 4, 5, 6}

In [28]:

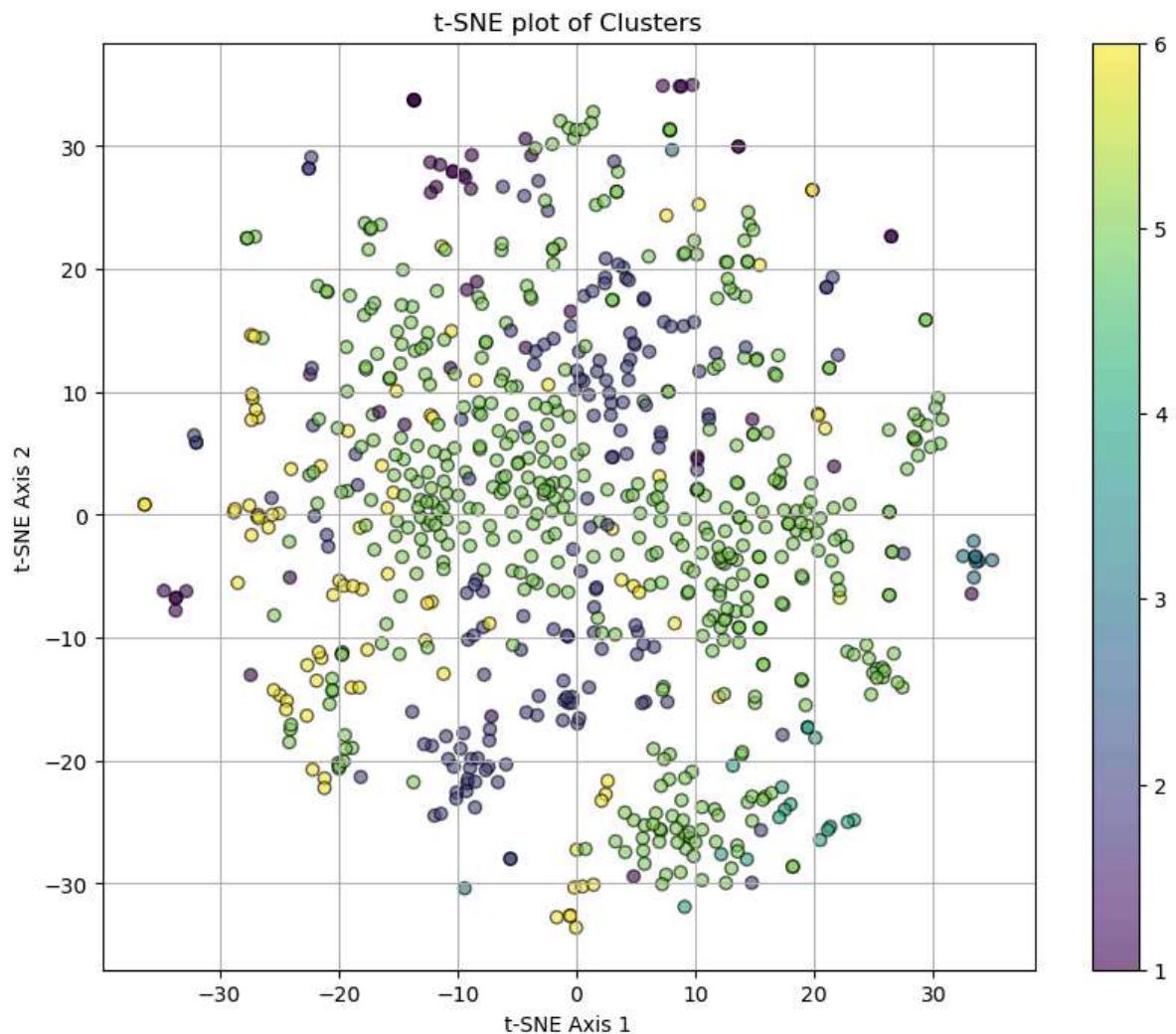
1 count_labels(fclusters)

Out[28]: Counter({5: 482, 2: 159, 6: 87, 1: 45, 4: 16, 3: 10})

```
In [29]: 1 plot_clusters(X, fclusters)
          2
```



In [30]: 1 plot_tsne(X, fclusters)



Evaluation metric

In [31]: 1 matices(X, fclusters)

Out[31]:

	Metric	Score
0	Silhouette Score	0.100648
1	Calinski-Harabasz Index	12.903035
2	Davies-Bouldin Index	3.897753

GCN

The provided Python script is part of a process for setting up graph data in a format suitable for use with machine learning frameworks that deal with graph-based data, like PyTorch Geometric. Here's a breakdown of what each part of the script does:

1. Feature Matrix Creation:

- The script creates a feature matrix where each feature is the degree of a node. It collects the degrees of nodes in a sorted order, converts them into a PyTorch tensor of floats, and adds an extra dimension to make it a column vector.
- This is done with the line `features = torch.tensor([G.degree(node) for node in sorted(G.nodes)], dtype=torch.float).unsqueeze(1)`.

2. Edge Index Tensor:

- An edge index tensor is created from the list of edges in the graph `G`. The tensor is converted to a long data type and then transposed to ensure the shape matches PyTorch Geometric's expectations (two rows where each column is an edge and rows represent source and target nodes respectively).
- This tensor is made contiguous for performance reasons in memory, using `edge_index = torch.tensor(list(G.edges), dtype=torch.long).t().contiguous()`.

3. Graph Data Object:

- A `Data` object (assuming from PyTorch Geometric) is instantiated with the node features `x` and the edge index `edge_index`. This object conveniently packages the graph data for processing in graph neural networks.
- The setup is demonstrated with `data = Data(x=features, edge_index=edge_index)`.

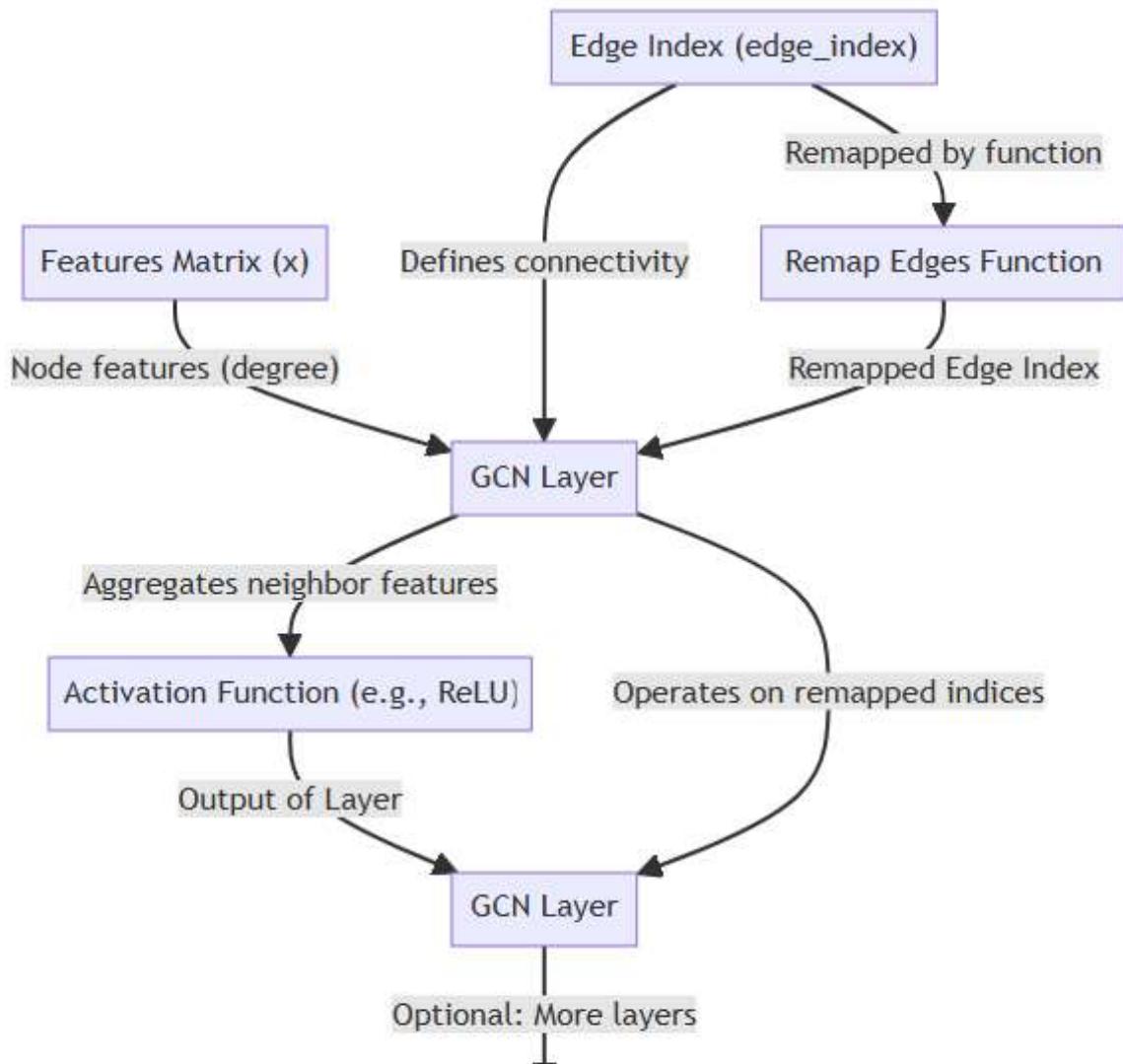
4. Function to Remap Edge Indices:

- The function `remap_edges` takes an edge index tensor and a dictionary mapping old node indices to new ones. It iterates through each edge, updating indices based on the mapping, ensuring the edge indices align with new or reordered node indices.
- This function is critical in cases where node indices may have been changed or reordered and is applied to the graph data with `data.edge_index = remap_edges(data.edge_index, node_to_index)`.

5. Output:

- Finally, the script prints the `data` object to display its structure and contents, which include the node features and the possibly remapped edge indices.

This script is particularly useful for preparing graph data for training graph neural networks, where node features and correct edge connections are crucial.



```
In [32]: 1 # Create a feature matrix with node degrees as features
2 features = torch.tensor([G.degree(node) for node in sorted(G.nodes)], dtype=torch.float)
3
4 # Create an edge index tensor
5 edge_index = torch.tensor(list(G.edges), dtype=torch.long).t().contiguous()
6
7 # Assuming all nodes are part of the same graph and unlabelled
8 data = Data(x=features, edge_index=edge_index)
9
10 # Function to remap edge indices
11 def remap_edges(edge_index, node_to_index):
12     for i in range(edge_index.size(1)):
13         edge_index[0, i] = node_to_index[edge_index[0, i].item()]
14         edge_index[1, i] = node_to_index[edge_index[1, i].item()]
15     return edge_index
16
17 # Assuming 'data.edge_index' is your original edge index tensor
18 data.edge_index = remap_edges(data.edge_index, node_to_index)
19
20
21
22 print(data)
```

Data(x=[799, 1], edge_index=[2, 5358])

```
In [33]: 1 edge_index_max = data.edge_index.max()
2 num_nodes = data.x.size(0) # Assuming 'data.x' represents node features just like edge_index
3
4 print(f"Maximum node index in 'edge_index': {edge_index_max}")
5 print(f"Number of nodes based on 'data.x': {num_nodes}")
6
7 if edge_index_max >= num_nodes:
8     raise ValueError("Edge index out of bounds. Make sure all indices are less than or equal to the number of nodes")
9
```

Maximum node index in 'edge_index': 798
Number of nodes based on 'data.x': 799

In [34]:

```

1  class GCNEncoder(torch.nn.Module):
2      def __init__(self, in_channels, out_channels):
3          super(GCNEncoder, self).__init__()
4          self.conv1 = GCNConv(in_channels, 2 * out_channels, cached=True)
5          self.conv2 = GCNConv(2 * out_channels, out_channels, cached=True)
6
7      def forward(self, x, edge_index):
8          x = F.relu(self.conv1(x, edge_index))
9          return self.conv2(x, edge_index)
10
11 class InnerProductDecoder(torch.nn.Module):
12     """Decoding node embeddings to reconstruct the graph."""
13     def forward(self, z, edge_index, sigmoid=True):
14         value = (z[edge_index[0]] * z[edge_index[1]]).sum(dim=1)
15         return torch.sigmoid(value) if sigmoid else value
16
17 class GraphAutoencoder(torch.nn.Module):
18     def __init__(self, in_channels, out_channels):
19         super(GraphAutoencoder, self).__init__()
20         self.encoder = GCNEncoder(in_channels, out_channels)
21         self.decoder = InnerProductDecoder()
22
23     def forward(self, x, edge_index):
24         z = self.encoder(x, edge_index)
25         return self.decoder(z, edge_index)
26
27 # Load your graph data
28
29 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
30 # device = 'cpu'
31 num_features = data.num_features
32 out_features = 64 # Dimensionality of the output embeddings
33 model = GraphAutoencoder(num_features, out_features).to(device)
34 data = data.to(device) # Assuming 'data' is your input dataset that has x
35 optimizer = Adam(model.parameters(), lr=0.01)
36
37 def train():
38     model.train()
39     optimizer.zero_grad()
40     z = model.encoder(data.x.to(device), data.edge_index.to(device))
41     edge_index, link_labels = get_link_labels(data.edge_index, data.num_no
42     link_logits = model.decoder(z, edge_index.to(device))
43     link_labels = link_labels.to(device)
44     loss = F.binary_cross_entropy(link_logits, link_labels)
45     loss.backward()
46     optimizer.step()
47     return loss.item()
48
49 def get_link_labels(pos_edge_index, num_nodes):
50     # Remove self-loops if they exist
51     pos_edge_index, _ = remove_self_loops(pos_edge_index)
52     # Add self-loops back to preserve the main diagonal if needed
53     pos_edge_index, _ = add_self_loops(pos_edge_index)
54     # Sample negative edges
55     neg_edge_index = negative_sampling(pos_edge_index, num_nodes, num_neg_
56
57     # Concatenate positive and negative samples

```

```
58     edge_index = torch.cat([pos_edge_index, neg_edge_index], dim=1)
59
60     # Create Link Labels: 1 for existing edges, 0 for negative samples
61     link_labels = torch.cat([torch.ones(pos_edge_index.size(1), dtype=torch.float),
62                             torch.zeros(neg_edge_index.size(1), dtype=torch.float)])
63
64
65
66     # Example of a training Loop
67     for epoch in range(200):
68         loss = train()
69         print(f'Epoch: {epoch+1}, Loss: {loss:.4f}')
70
```

```
Epoch: 1, Loss: 40.1179
Epoch: 2, Loss: 31.4031
Epoch: 3, Loss: 17.9443
Epoch: 4, Loss: 16.0224
Epoch: 5, Loss: 16.1675
Epoch: 6, Loss: 15.1585
Epoch: 7, Loss: 9.8795
Epoch: 8, Loss: 5.1327
Epoch: 9, Loss: 6.9878
Epoch: 10, Loss: 7.0480
Epoch: 11, Loss: 3.8040
Epoch: 12, Loss: 2.1236
Epoch: 13, Loss: 2.0728
Epoch: 14, Loss: 2.3655
Epoch: 15, Loss: 1.5440
Epoch: 16, Loss: 1.0220
Epoch: 17, Loss: 1.2008
Epoch: 18, Loss: 1.0054
Epoch: 19, Loss: 0.8858
Epoch: 20, Loss: 0.7852
```

```
In [35]: 1 model.eval() # Set the model to evaluation mode to disable dropout layers
2 with torch.no_grad(): # Disable gradient computation
3     embeddings_gcn = model.encoder(data.x.to(device), data.edge_index.to(device))
4
```

```
In [36]: 1 print(len(embeddings_gcn))
```

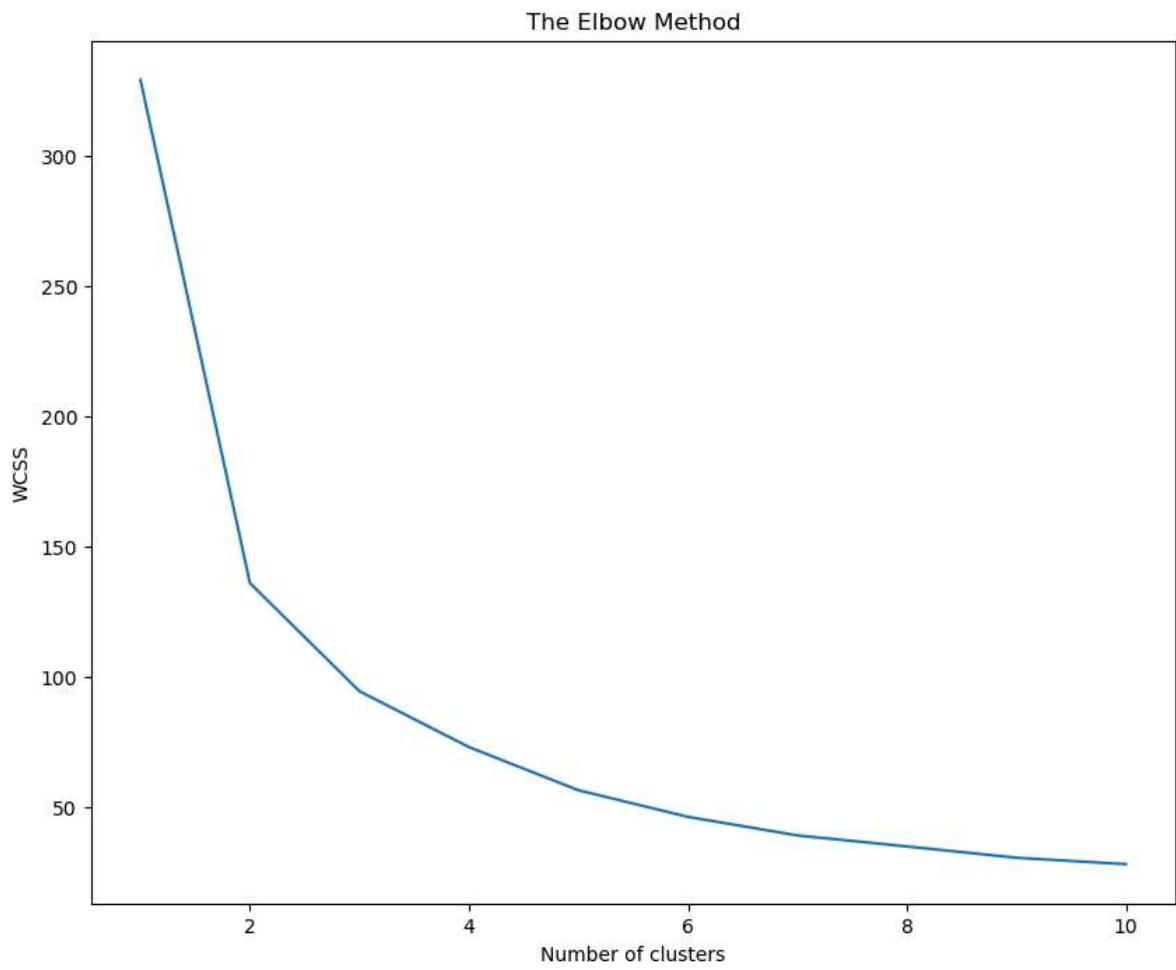
799

K-Means

In [37]:

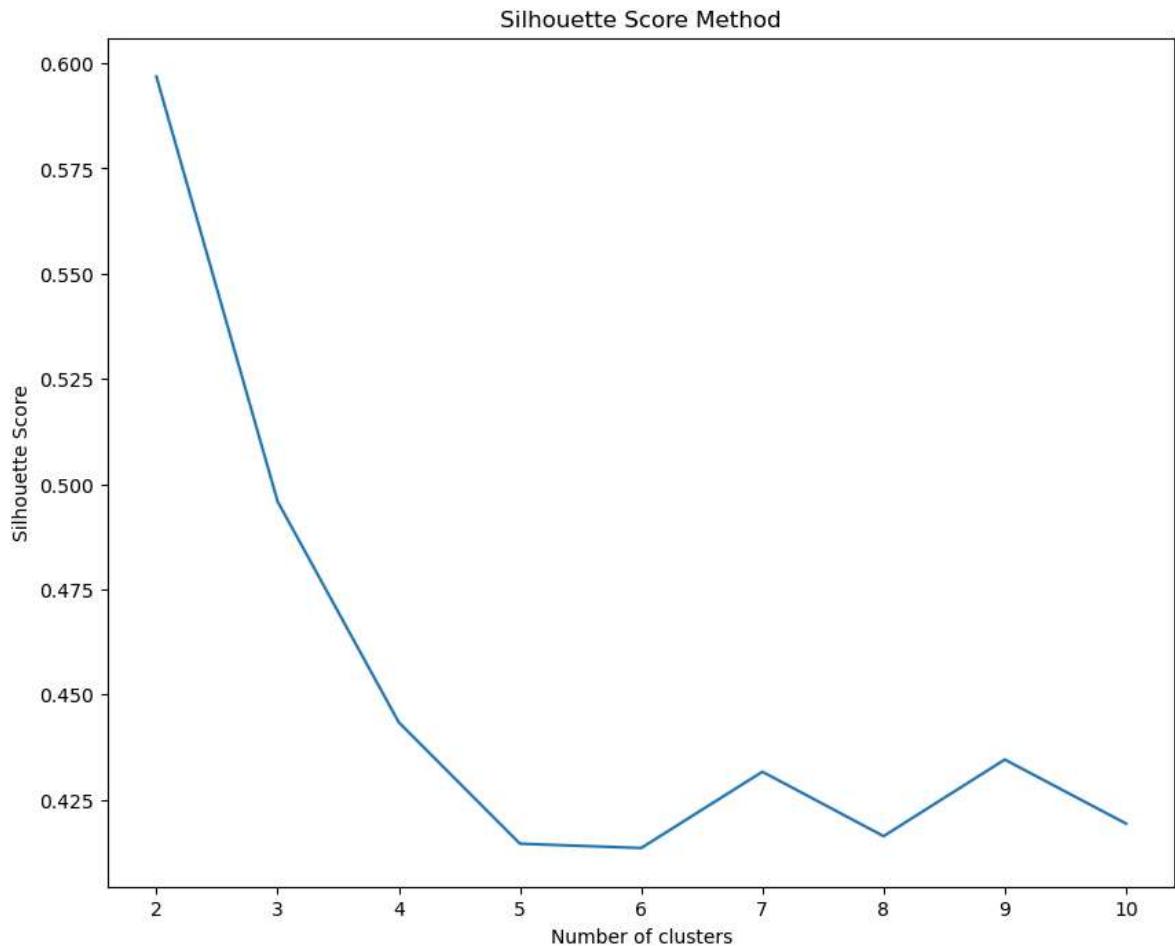
```
1 embeddings_np_gcn = embeddings_gcn.cpu().detach().numpy()
2
3 KMeansClassFinder(embeddings_np_gcn, 11)
4
```

```
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
```



```
In [38]: 1 KMeansClassFinder_silhouette_scores(embeddings_np_gcn)
```

```
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
```



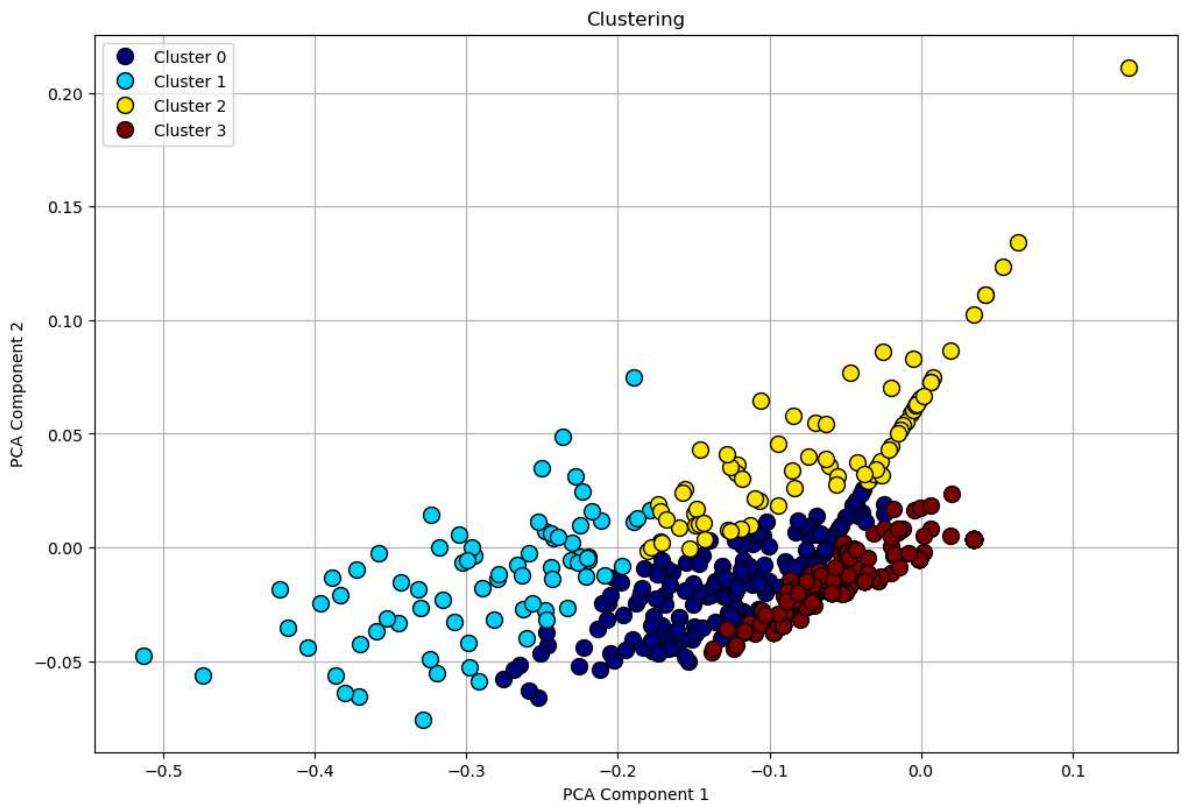
```
In [39]: 1 gcn_kclusters = applyKMeans(embeddings_np_gcn, 4)
          2
```

```
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1412: 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)
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
{0, 1, 2, 3}
```

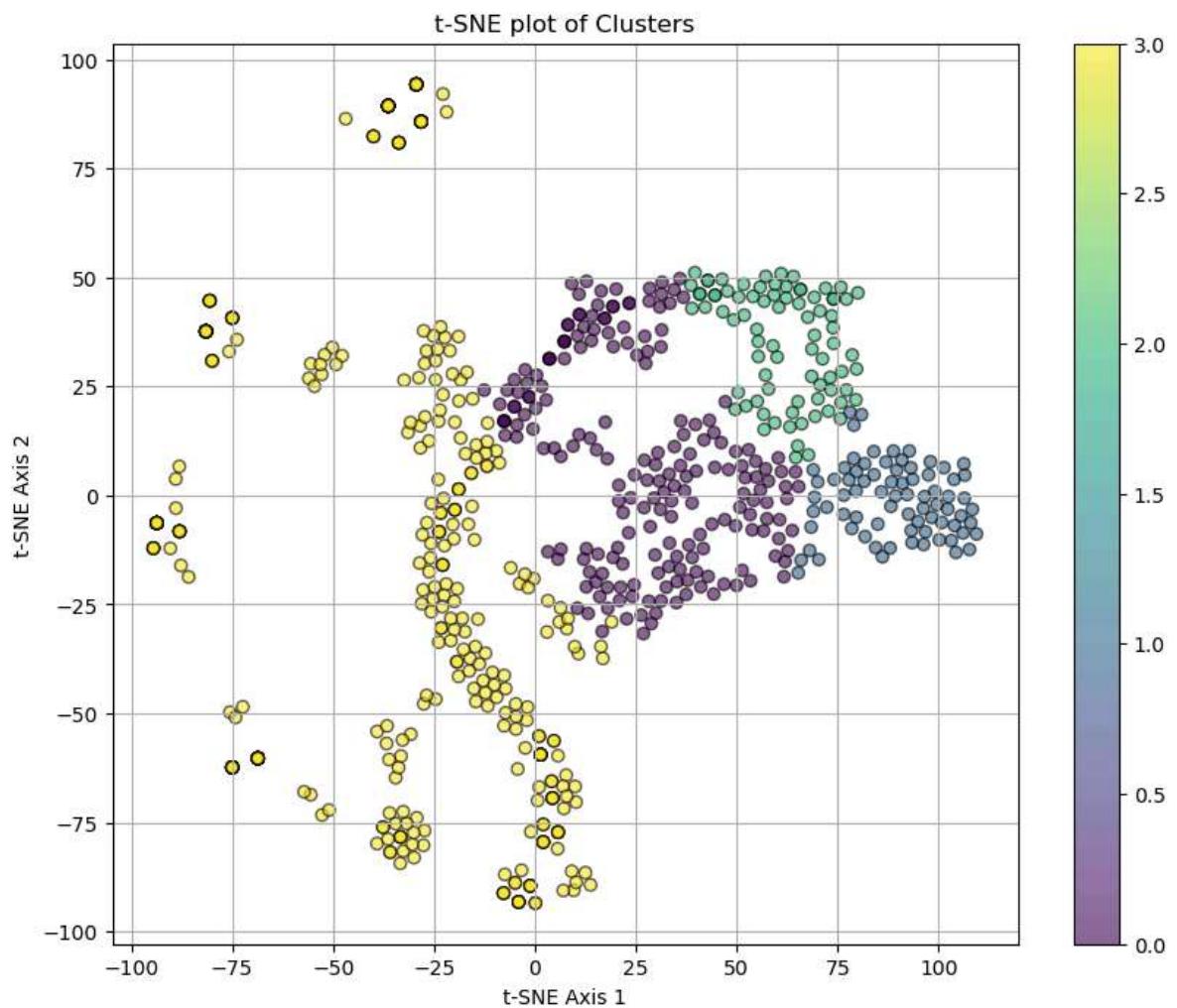
```
In [40]: 1 count_labels(gcn_kclusters.labels_)
```

```
Out[40]: Counter({3: 439, 0: 203, 2: 81, 1: 76})
```

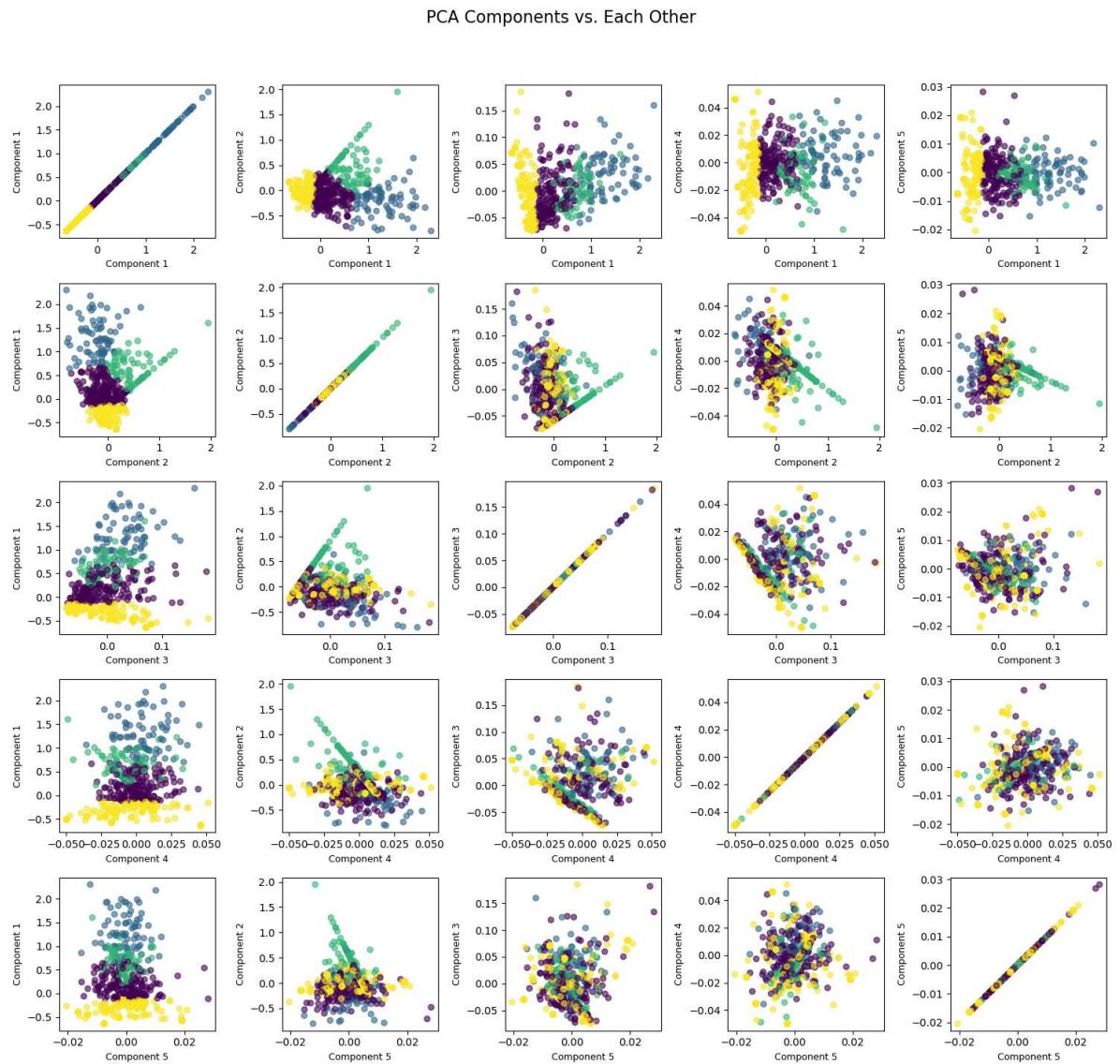
```
In [41]: 1 plot_clusters(embeddings_np_gcn, gcn_kclusters.labels_)
```



```
In [42]: 1 plot_tsne(embeddings_np_gcn, gcn_kclusters.labels_)
```



In [43]: 1 plot_all_pca_components(embeddings_np_gcn, gcn_kclusters.labels_, n_compor



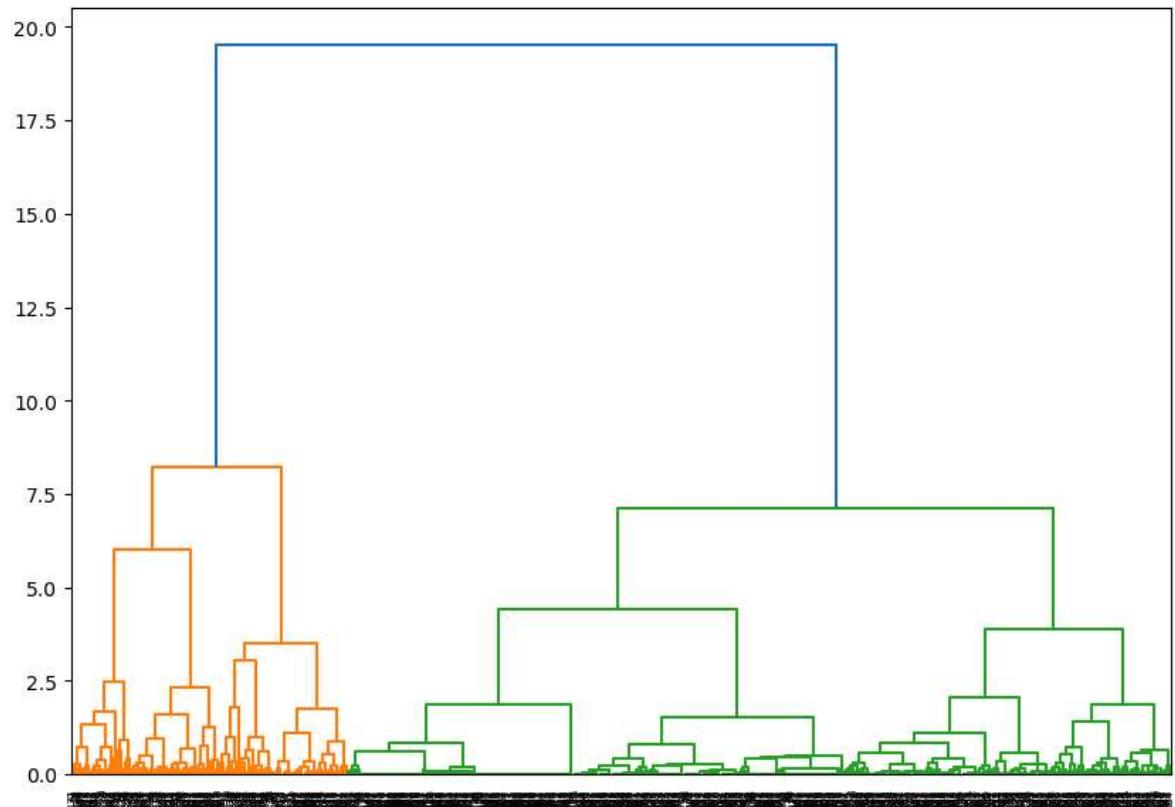
In [44]: 1 matices(X, gcn_kclusters.labels_)

Out[44]:

	Metric	Score
0	Silhouette Score	0.071329
1	Calinski-Harabasz Index	4.046779
2	Davies-Bouldin Index	10.459836

Hierarchical Clustering

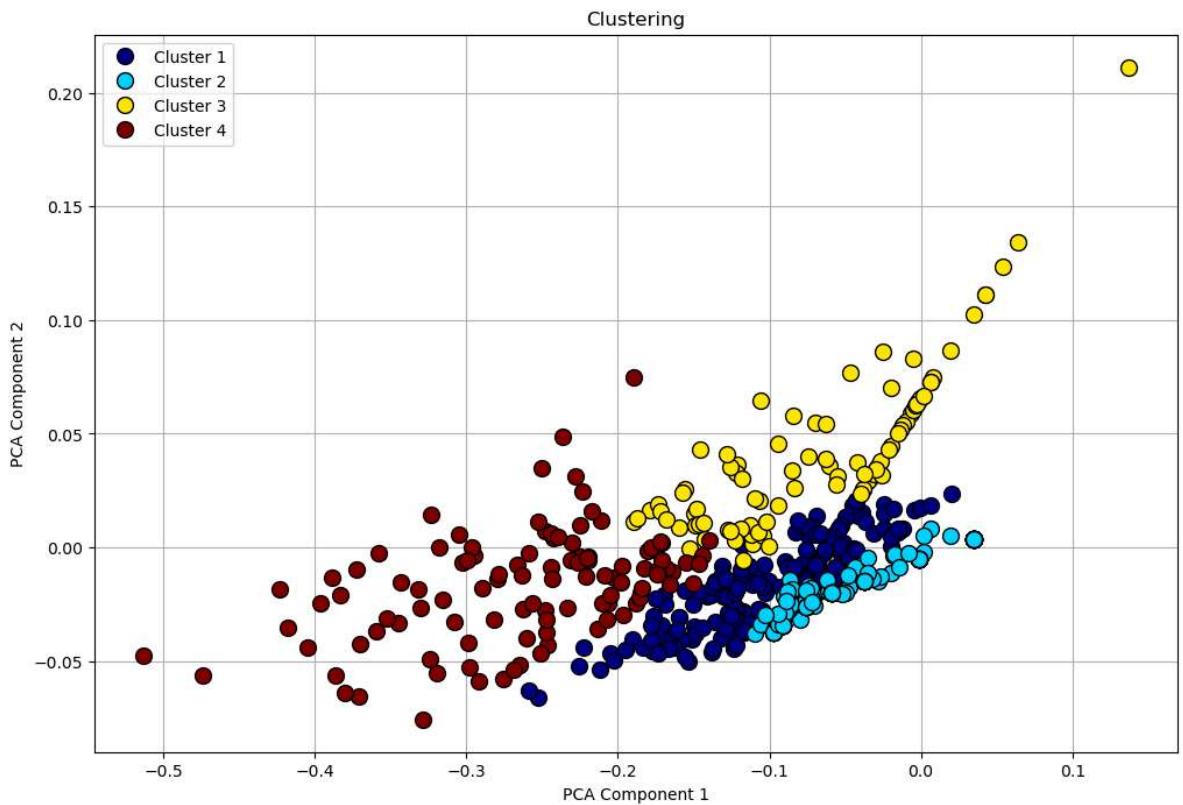
```
In [45]: 1 linked_gcn = plot_dendrogram(embeddings_np_gcn)
```



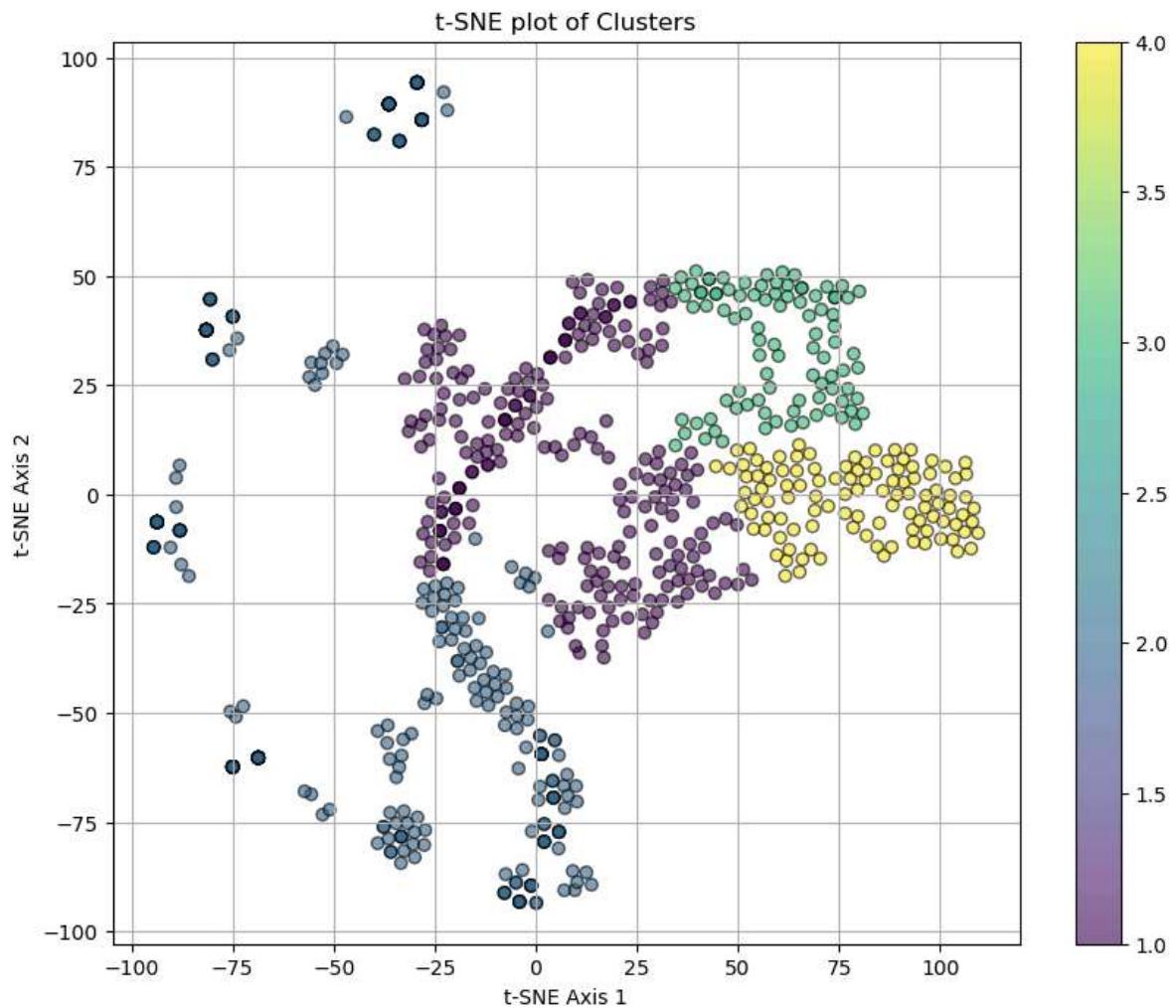
```
In [46]: 1 fclusters_gcn = applyFcluster(linked_gcn, 6.5)  
2
```

```
{1, 2, 3, 4}
```

```
In [47]: 1 plot_clusters(embeddings_np_gcn, fclusters_gcn)
```



In [48]: 1 plot_tsne(embeddings_np_gcn, fclusters_gcn)



In [49]: 1 matices(X, fclusters_gcn)

	Metric	Score
0	Silhouette Score	0.014585
1	Calinski-Harabasz Index	3.726467
2	Davies-Bouldin Index	10.939329

Spectral

The provided Python script performs spectral graph analysis using NetworkX and NumPy libraries. Here's a detailed explanation of each part of the script:

1. Adjacency Matrix Creation:

- The script starts by converting a graph G into its adjacency matrix A using NetworkX's `to_numpy_array` function. The adjacency matrix A represents the connections between nodes, where each entry $A[i, j]$ is 1 if there is an edge between node i and node j , and 0 otherwise.

2. Degree Matrix Calculation:

- Next, it computes the degree matrix D , which is a diagonal matrix where each diagonal entry $D[i, i]$ is the sum of the i -th row of the adjacency matrix A . This represents the number of connections (degree) of each node.

3. Laplacian Matrix Computation:

- The Laplacian matrix L is calculated as $D - A$. The Laplacian matrix is a fundamental matrix in graph theory used for various applications like clustering, network analysis, and dimensionality reduction. It reflects the structure of the graph in terms of its connectivity.

4. Eigenvalues and Eigenvectors:

- The script then calculates the eigenvalues and eigenvectors of the Laplacian matrix using the `eigh` function, which is optimized for symmetric matrices like L . The eigenvalues are sorted in increasing order since the Laplacian matrix is positive semi-definite, meaning all eigenvalues are non-negative.

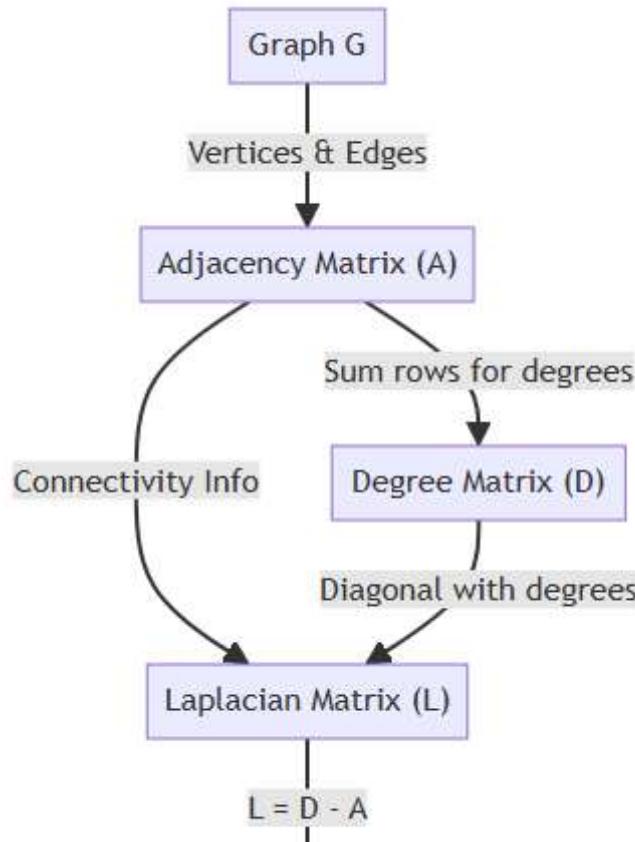
5. Spectral Embedding:

- To reduce the graph into a lower-dimensional space for tasks like visualization or clustering, the script selects a subset of the smallest eigenvalues (skipping the smallest trivial zero eigenvalue which represents the overall connectivity of the graph). It then uses the corresponding eigenvectors to form a `spectral_embedding`.
- The number of dimensions for the embedding is set to 64, but this can be adjusted based on the specific application or the graph's size and complexity.

6. Output:

- The script originally intended to print the smallest eigenvalues and the spectral embedding, though the actual printing of the embedding is commented out in the code provided. The user can uncomment these lines to see the embedding or modify the print statements to focus on different aspects of the output, such as specific eigenvectors or dimensions of the embedding.

This code effectively demonstrates how to perform spectral graph analysis, which is a powerful technique for understanding the structure and properties of graphs through their spectral characteristics.



```
In [50]: 1 is_connected = nx.is_connected(G.to_undirected())
2 print("Is the graph connected?", is_connected)
3
```

Is the graph connected? False

In [51]:

```
1 # Step 1: Adjacency matrix
2 A = nx.to_numpy_array(G)
3
4 # Step 2: Degree matrix
5 D = np.diag(A.sum(axis=1))
6
7 # Step 3: Laplacian matrix
8 L = D - A
9
10 # Step 4: Eigenvalues and eigenvectors
11 eigenvalues, eigenvectors = eigh(L)
12 print("Smallest eigenvalues:", len(eigenvalues)) # adjust range as needed
13
14 # Since the Laplacian matrix is positive semi-definite, the eigenvalues are
15 # We sort the eigenvalues, skipping the first small trivial eigenvalue 0
16 # Assume we want 2 dimensions for our embedding
17 num_dimensions = 64
18 smallest_eigenvalues = np.argsort(eigenvalues)[1:num_dimensions+1] # skip
19
20 # Corresponding eigenvectors are our spectral embedding
21 spectral_embedding = eigenvectors[:, smallest_eigenvalues]
22
23 # print("Spectral Embedding:\n", spectral_embedding)
24
25 # for ele in spectral_embedding:
26 #     print(ele)
```

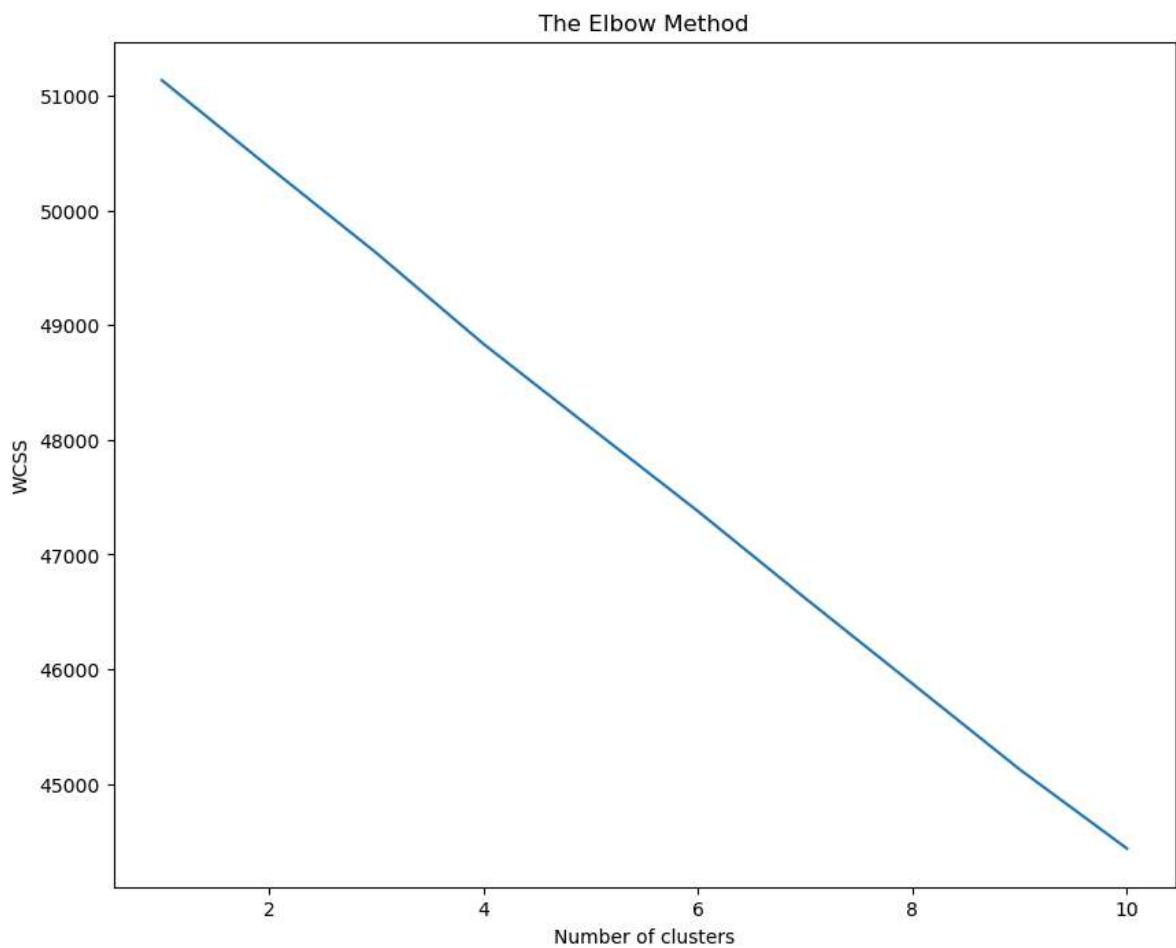
Smallest eigenvalues: 799

K-Means

In [52]:

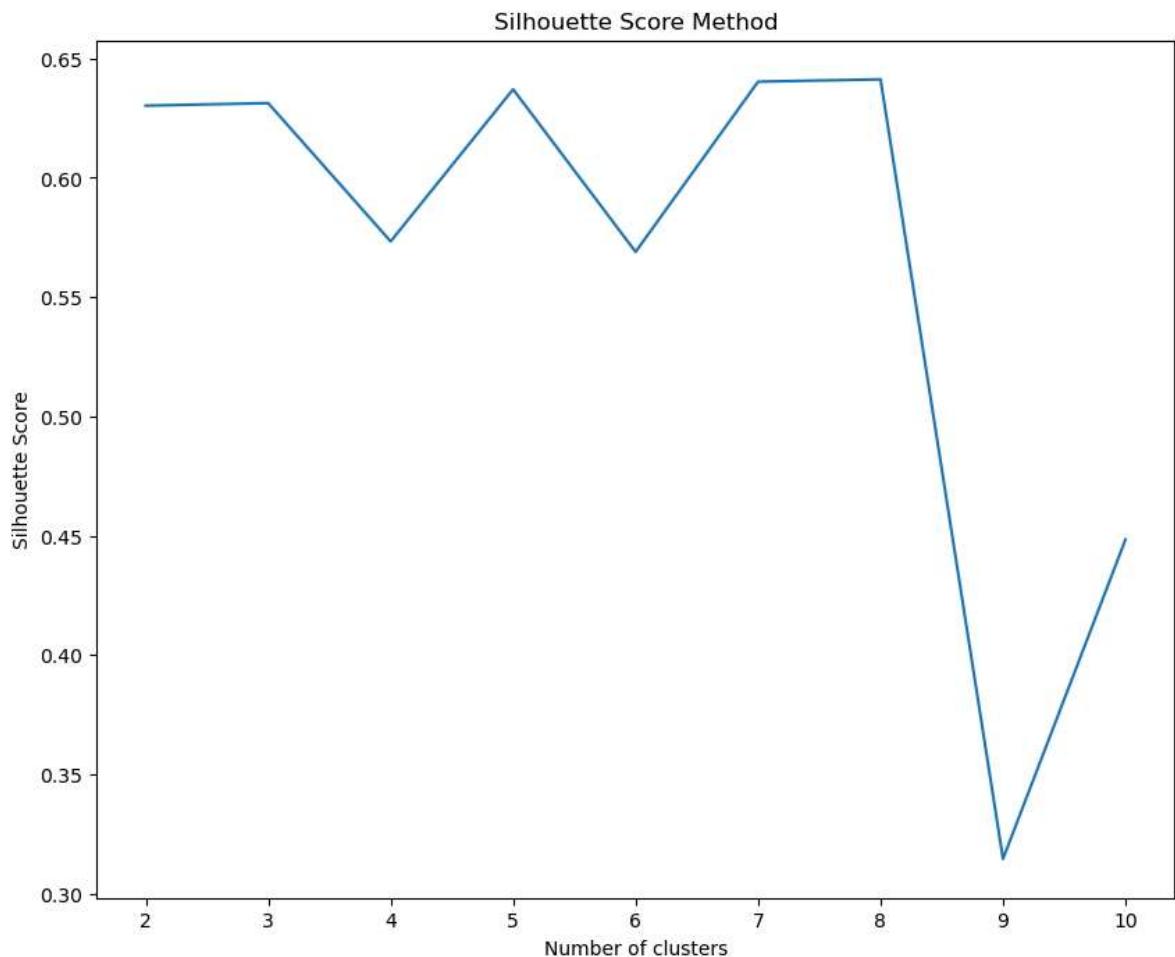
```
1 scaler = StandardScaler()  
2 spectral_embedding = scaler.fit_transform(spectral_embedding)  
3 KMeansClassFinder(spectral_embedding, 11)
```

```
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
```



In [53]: 1 KMeansClassFinder_silhouette_scores(spectral_embedding)

```
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.  
    warnings.warn(  
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
```



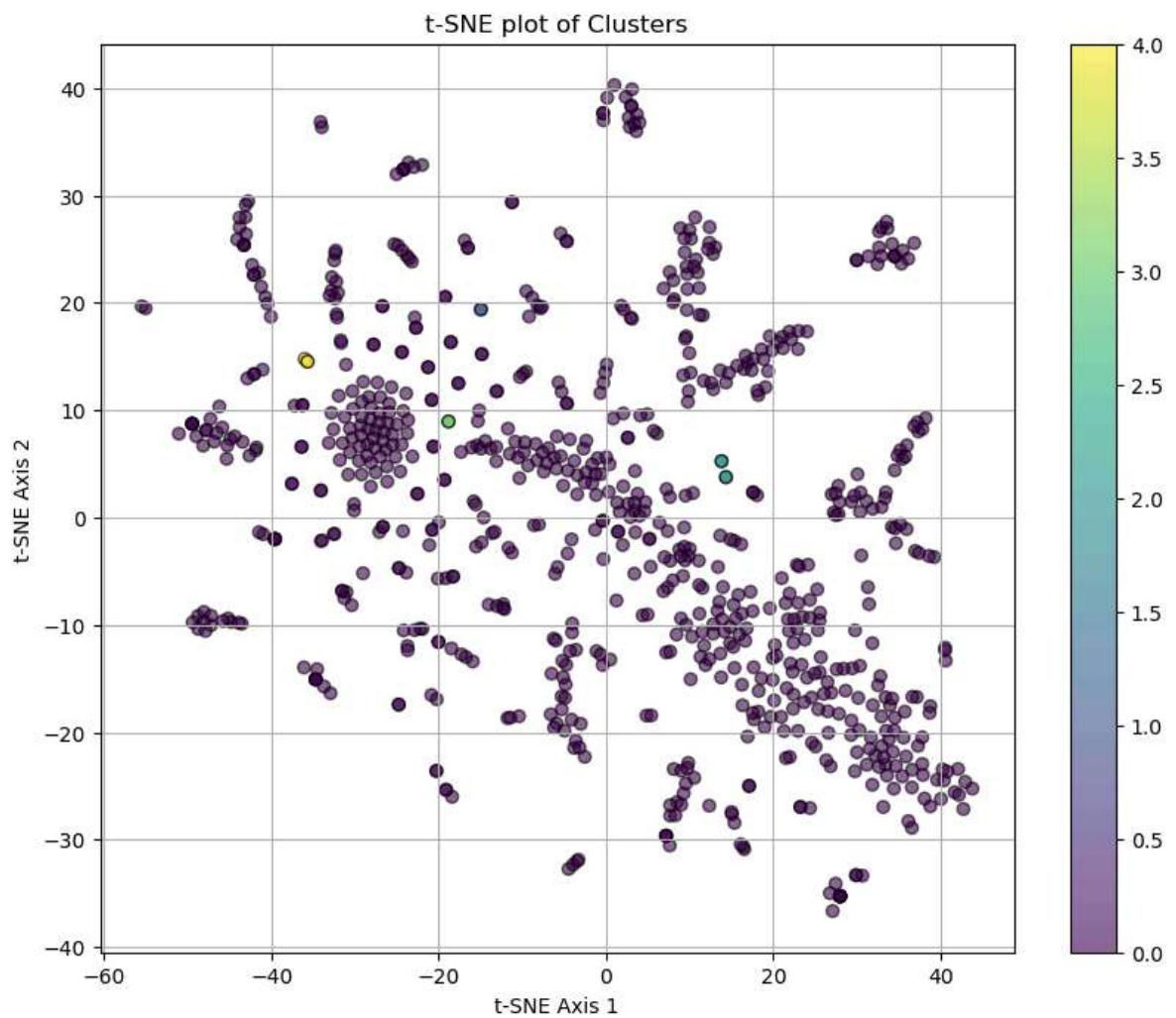
In [54]: 1 spc_clusters = applyKMeans(spectral_embedding, 5)

```
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1412: 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)
c:\Users\gupta\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=4.
    warnings.warn(
{0, 1, 2, 3, 4}
```

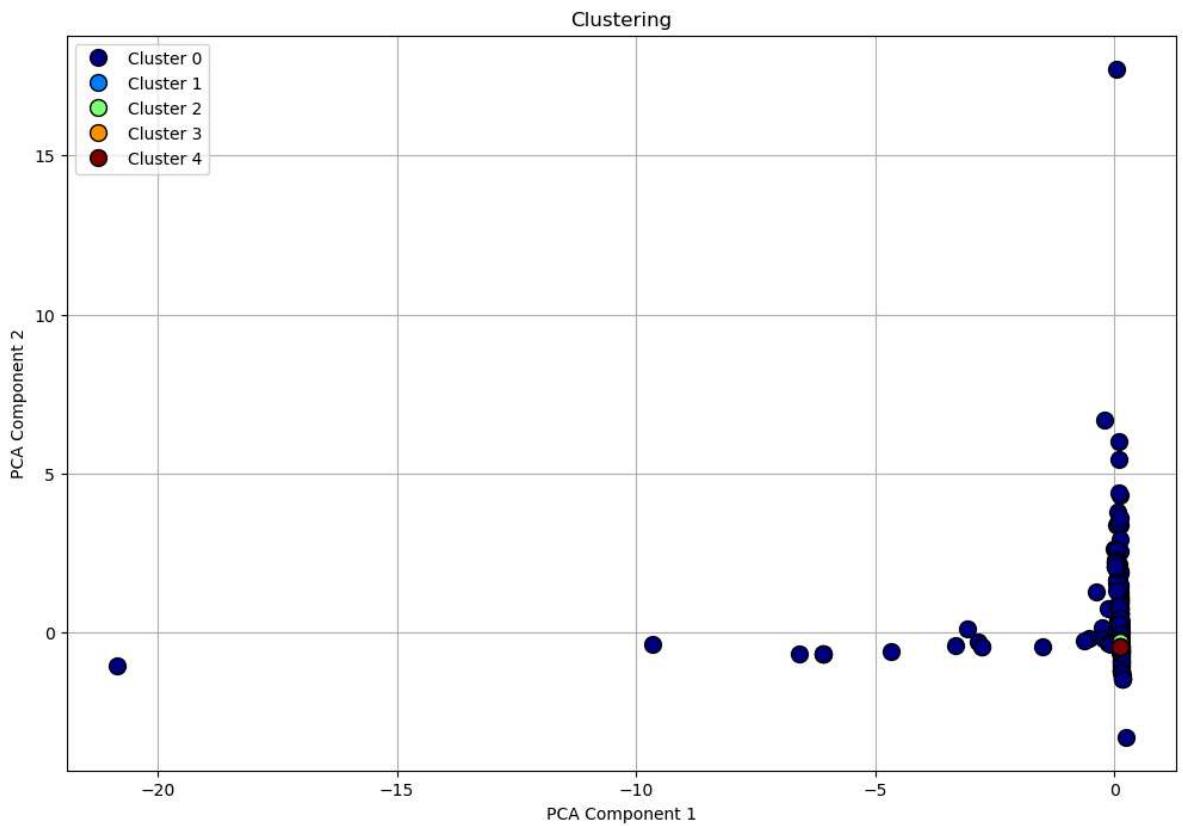
In [55]: 1 count_labels(spc_clusters.labels_)

Out[55]: Counter({0: 788, 2: 4, 4: 3, 1: 2, 3: 2})

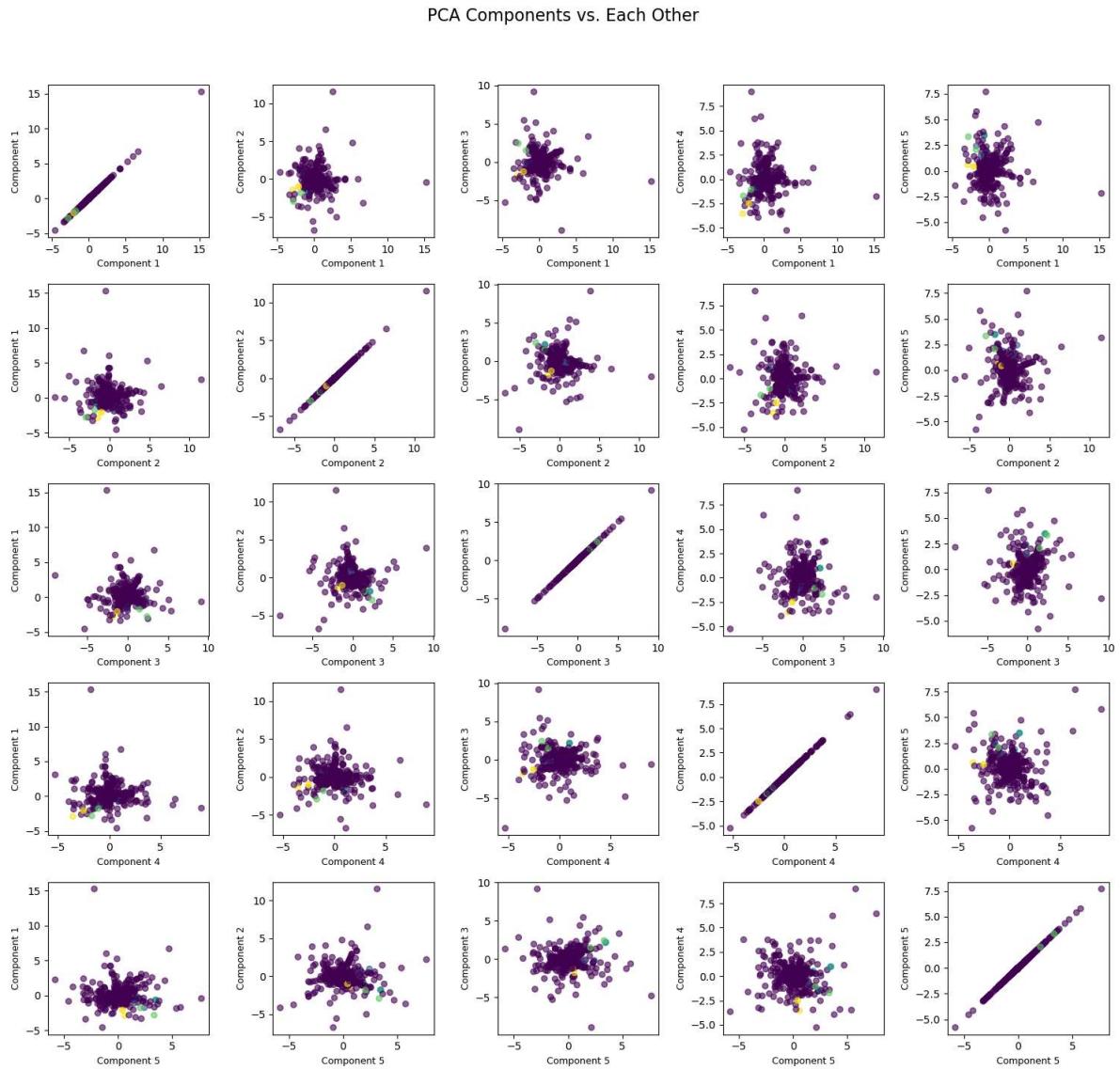
```
In [56]: 1 plot_tsne(spectral_embedding, spc_clusters.labels_)
```



```
In [57]: 1 plot_clusters(spectral_embedding, spc_clusters.labels_)
```



In [58]: 1 plot_all_pca_components(spectral_embedding, spc_clusters.labels_, n_compor



In [59]: 1 matices(X, spc_clusters.labels_)

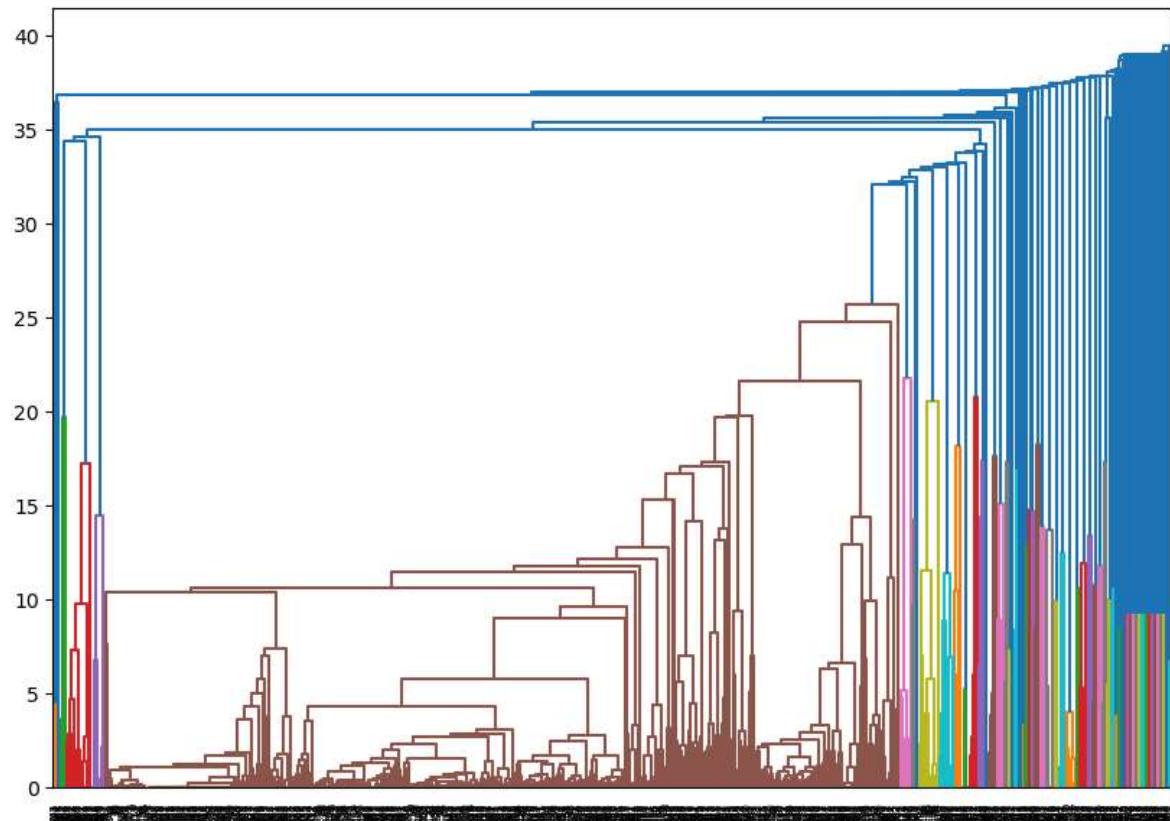
Out[59]:

Metric Score

0	Silhouette Score	-0.296521
1	Calinski-Harabasz Index	1.243626
2	Davies-Bouldin Index	2.344077

Hierarchical Clustering

In [60]: 1 linked_spec = plot_dendrogram(spectral_embedding)



In [61]: 1 spec_fclusters = applyFcluster(linked_spec, 32)

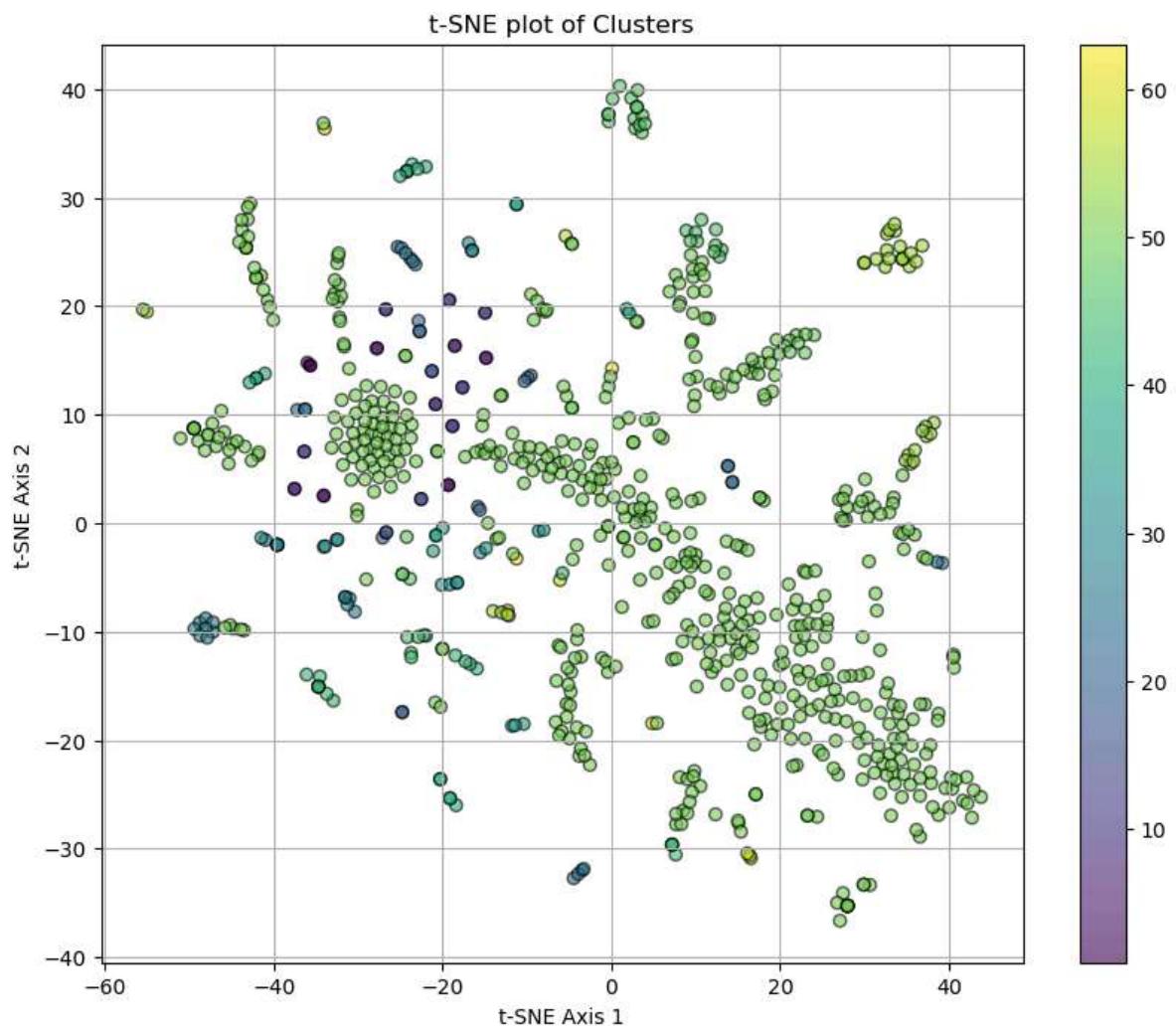
```
{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 2  
2, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 4  
1, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 6  
0, 61, 62, 63}
```

```
In [62]: 1 count_labels(spec_fcclusters)
```

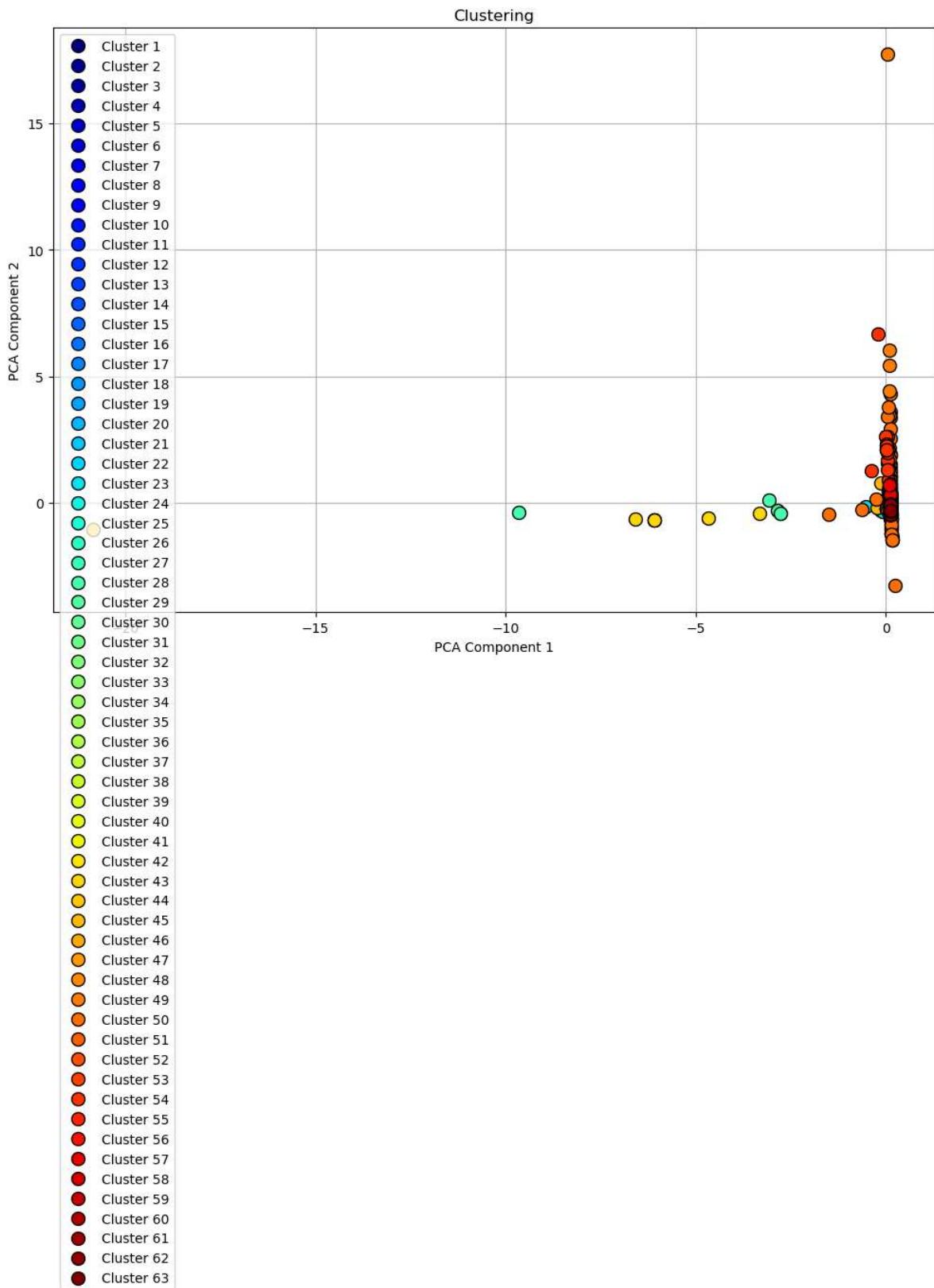
```
Out[62]: Counter({50: 568,
                  54: 18,
                  47: 15,
                  46: 10,
                  53: 10,
                  27: 8,
                  49: 8,
                  43: 8,
                  41: 7,
                  25: 6,
                  30: 6,
                  45: 6,
                  28: 5,
                  42: 5,
                  40: 5,
                  55: 5,
                  22: 4,
                  44: 4,
                  23: 4,
                  32: 4,
                  48: 4,
                  39: 4,
                  31: 4,
                  37: 4,
                  57: 3,
                  35: 3,
                  1: 3,
                  20: 3,
                  29: 3,
                  21: 3,
                  24: 3,
                  17: 3,
                  15: 2,
                  36: 2,
                  18: 2,
                  26: 2,
                  38: 2,
                  19: 2,
                  34: 2,
                  8: 2,
                  7: 2,
                  16: 2,
                  2: 2,
                  14: 2,
                  33: 2,
                  13: 2,
                  12: 2,
                  11: 2,
                  10: 2,
                  6: 2,
                  9: 2,
                  5: 2,
                  4: 2,
                  3: 2,
                  51: 1,
                  62: 1,
                  56: 1,})
```

```
61: 1,  
59: 1,  
60: 1,  
58: 1,  
52: 1,  
63: 1})
```

In [63]: 1 | plot_tsne(spectral_embedding, spec_fclusters)



In [64]: 1 plot_clusters(spectral_embedding, spec_fclusters)



```
In [65]: 1 matices(X, spec_fclusters)
```

```
Out[65]:
```

	Metric	Score
0	Silhouette Score	-0.290648
1	Calinski-Harabasz Index	1.582792
2	Davies-Bouldin Index	2.683453