

Міністерство освіти і науки України
Львівський національний університет імені Івана Франка
Факультет електроніки та комп'ютерних технологій

Звіт

Про виконання лабораторної роботи №5
З курсу «Системи опрацювання даних»
«Кластеризація даних»

Виконав:
Студент групи Фес-21
Шавало Андрій

Львів-2025

Мета роботи: Ознайомитися з методами кластеризації даних, реалізувати та порівняти алгоритми кластеризації (KMeans, DBSCAN, Agglomerative Clustering) на реальному датасеті, а також оцінити якість отриманих клас

Теоретичні відомості: Ознайомитись з алгоритмом кластеризації k-means та прикладом його реалізації

<https://www.kaggle.com/kushal1996/customer-segmentation-k-meansanalysis/notebook>

Хід роботи

Завдання 1. Підготовка даних о Завантажити набір даних о Виконати первинний аналіз даних: розмірність, типи змінних, наявність пропущених значень. о Провести масштабування (стандартизацію або нормалізацію) числових змінних.

```
1 > import ...
9
10 data = pd.read_csv('size_1000_n_5_sepval_0.1.csv')
11
12 print("Розмірність даних:", data.shape)
13 print("\nПерші 5 рядків:\n", data.head())
14 print("\nОпис даних:\n", data.describe())
15
16 plt.figure(figsize=(10, 6))
17 sns.scatterplot(x='x', y='y', data=data)
18 plt.title('Вхідні дані')
19 plt.show()
20
21 scaler = StandardScaler()
22 scaled_data = scaler.fit_transform(data[['x', 'y']])
[15]
```

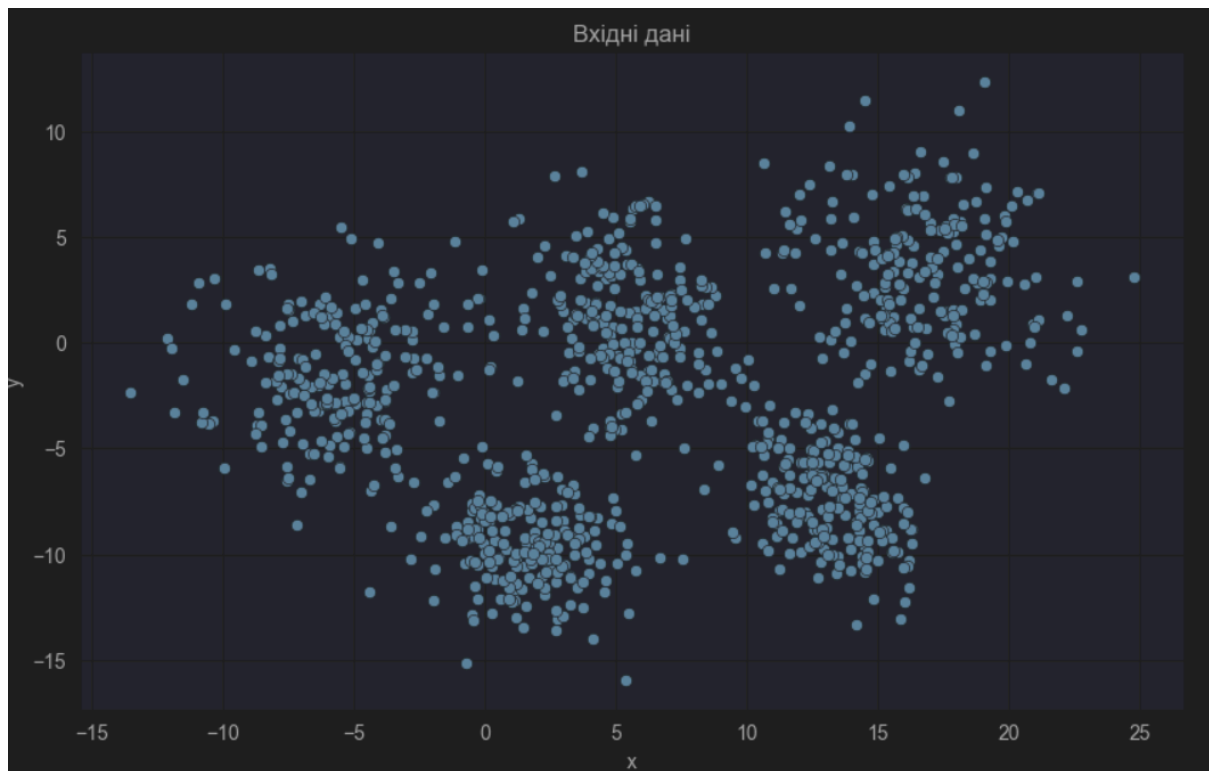
Розмірність даних: (1000, 4)

Перші 5 рядків:

	Unnamed: 0	x	y	run
0	1	-7.535642	1.787421	ourvCgP
1	2	4.351434	-0.506274	ourvCgP
2	3	-0.776635	-10.419561	ourvCgP
3	4	-4.529642	0.695028	ourvCgP
4	5	-5.107399	4.915450	ourvCgP

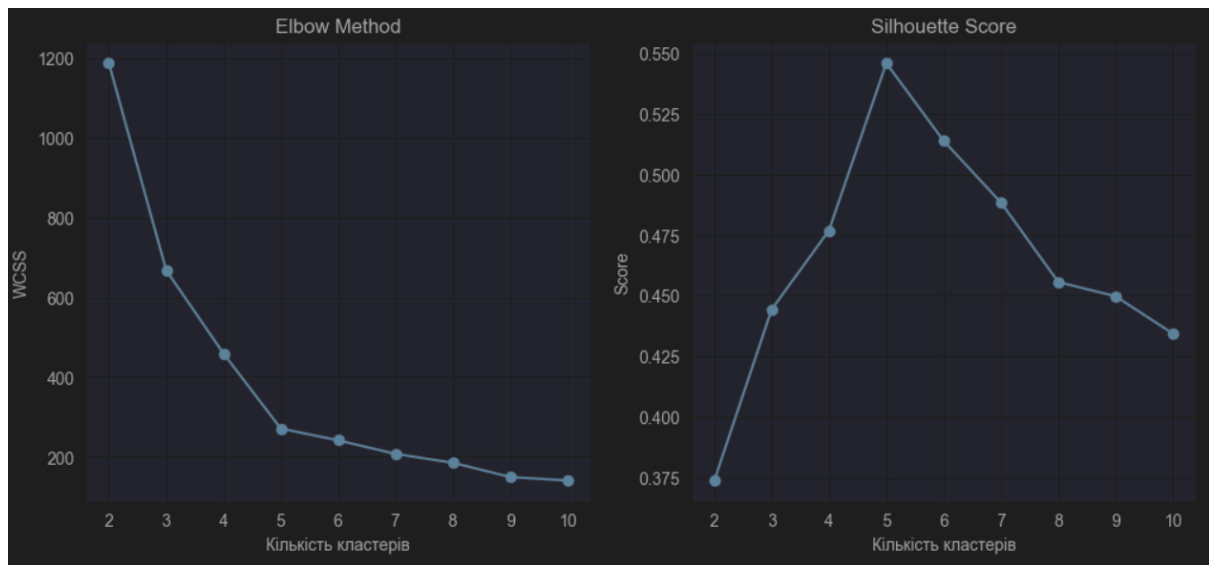
Опис даних:

	Unnamed: 0	x	y
count	1000.000000	1000.000000	1000.000000
mean	500.500000	6.207532	-2.752120
std	288.819436	8.245608	5.545930
min	1.000000	-13.548163	-15.939657
25%	250.750000	0.214693	-7.864758
50%	500.500000	5.446603	-2.192122
75%	750.250000	13.611509	1.606741
max	1000.000000	24.747021	12.345398



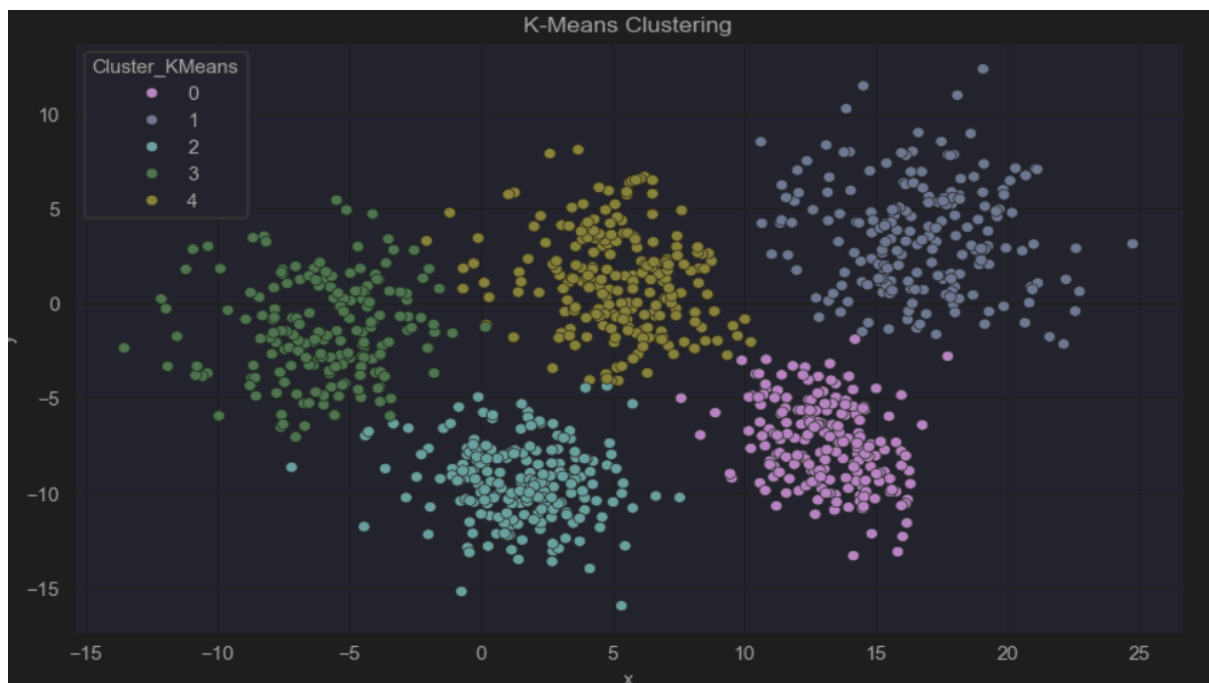
2. Реалізація кластеризації о Визначити оптимальну кількість кластерів за допомогою методу "Elbow Method" або Silhouette Score.

```
1 wcss = []
2 silhouette_scores = []
3 for i in range(2, 11):
4     kmeans = KMeans(n_clusters=i, random_state=42)
5     kmeans.fit(scaled_data)
6     wcss.append(kmeans.inertia_)
7     silhouette_scores.append(silhouette_score(scaled_data, kmeans.labels_))
8
9 plt.figure(figsize=(12, 5))
10 plt.subplot(1, 2, 1)
11 plt.plot(range(2, 11), wcss, marker='o')
12 plt.title('Elbow Method')
13 plt.xlabel('Кількість кластерів')
14 plt.ylabel('WCSS')
15
16 plt.subplot(1, 2, 2)
17 plt.plot(range(2, 11), silhouette_scores, marker='o')
18 plt.title('Silhouette Score')
19 plt.xlabel('Кількість кластерів')
20 plt.ylabel('Score')
21 plt.show()
```



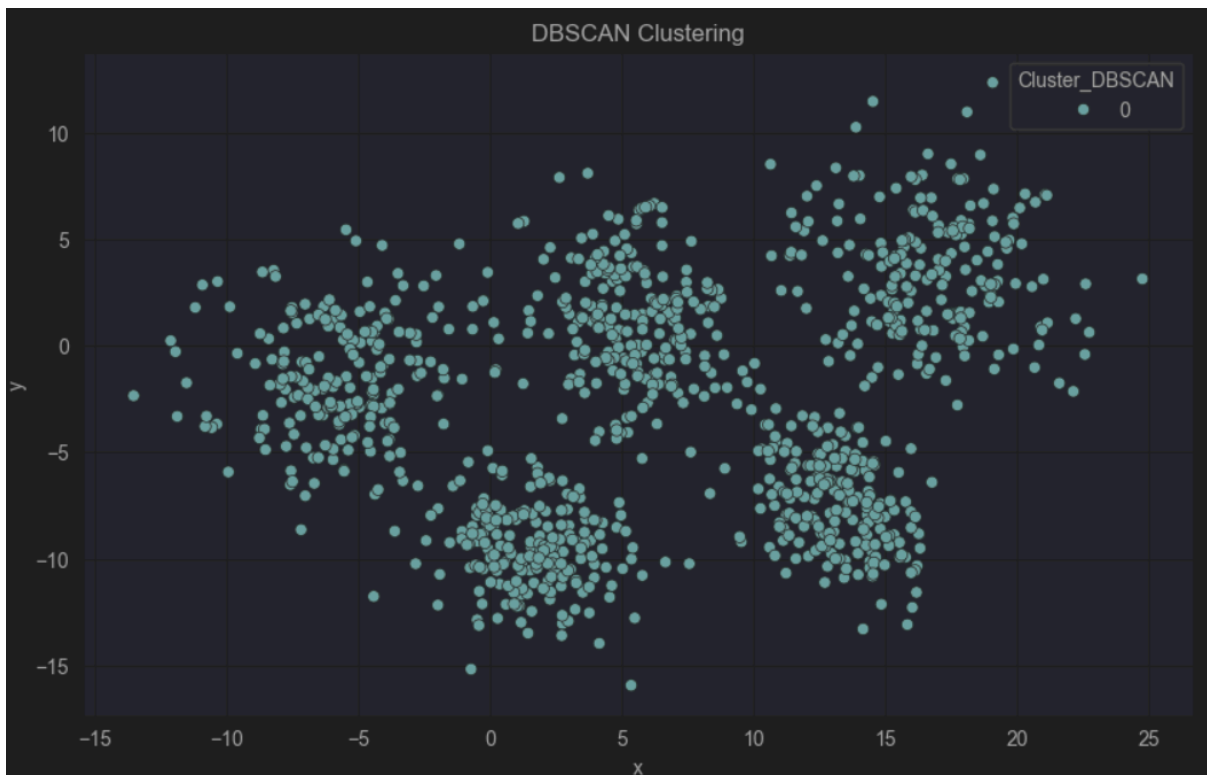
Виконати кластеризацію за допомогою K-Means.

```
1 optimal_k = 5
2 kmeans = KMeans(n_clusters=optimal_k, random_state=42)
3 clusters_kmeans = kmeans.fit_predict(scaled_data)
4
5 data['Cluster_KMeans'] = clusters_kmeans
6
7 plt.figure(figsize=(10, 6))
8 sns.scatterplot(x='x', y='y', hue='Cluster_KMeans', data=data, palette='viridis')
9 plt.title('K-Means Clustering')
10 plt.show()
```



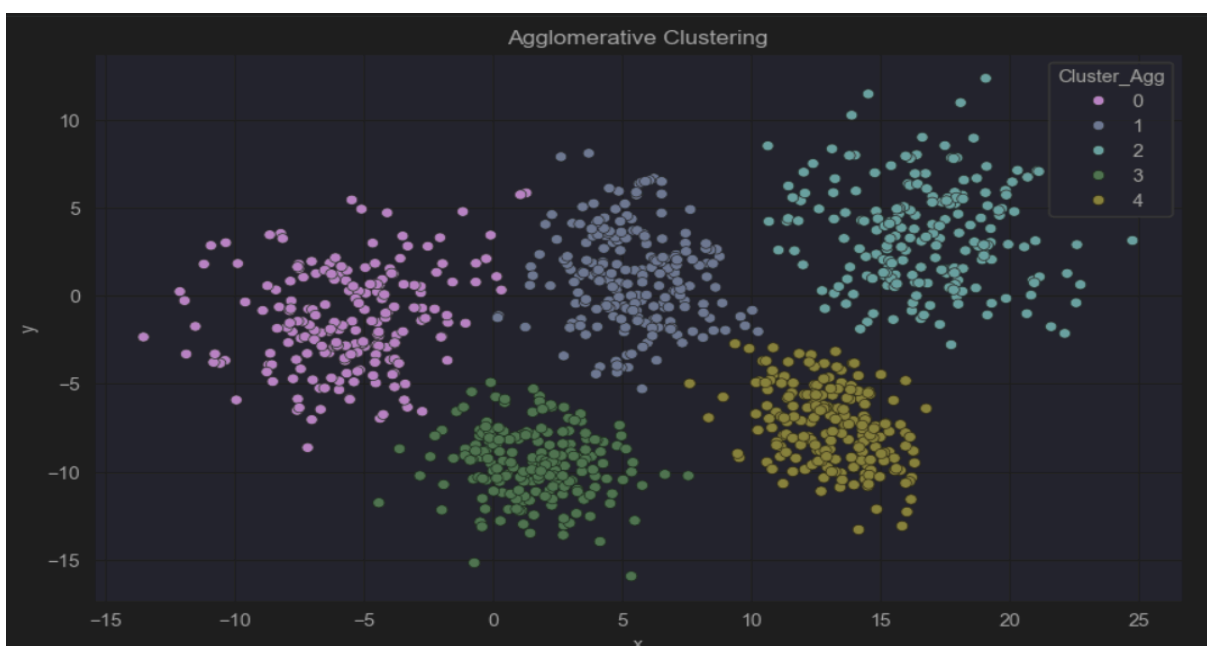
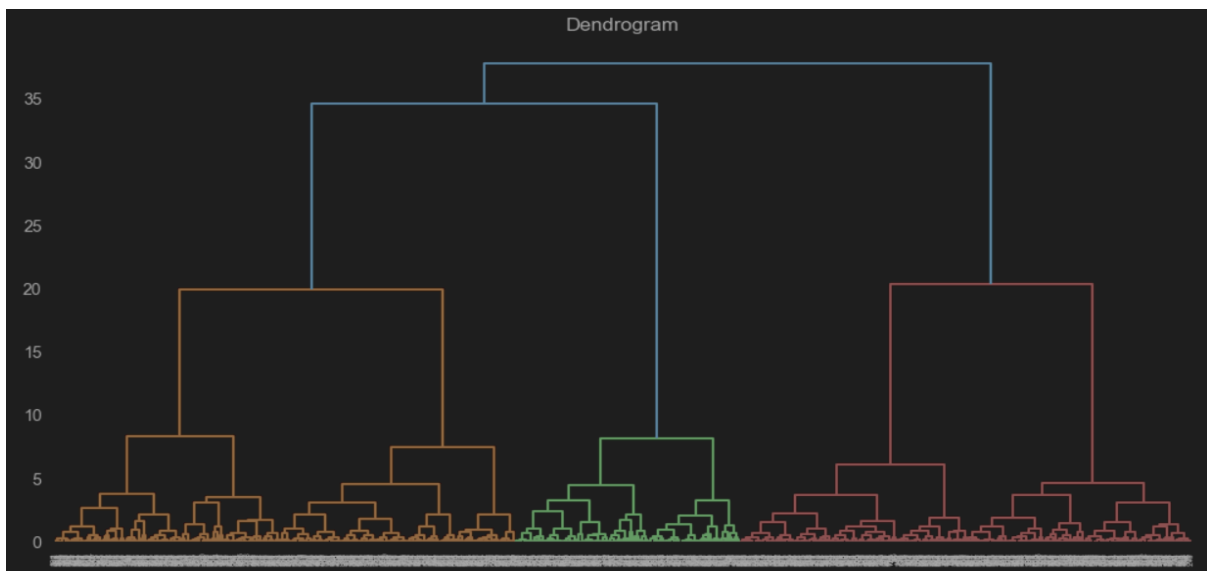
Реалізувати та застосувати DBSCAN для виявлення аномалій.

```
1 dbscan = DBSCAN(eps=0.5, min_samples=5)
2 clusters_dbscan = dbscan.fit_predict(scaled_data)
3
4 data['Cluster_DBSCAN'] = clusters_dbscan
5
6 plt.figure(figsize=(10, 6))
7 sns.scatterplot(x='x', y='y', hue='Cluster_DBSCAN', data=data, palette='viridis')
8 plt.title('DBSCAN Clustering')
9 plt.show()
10
11 n_noise = list(clusters_dbscan).count(-1)
12 print(f"Кількість аномалій (шум): {n_noise}")
```



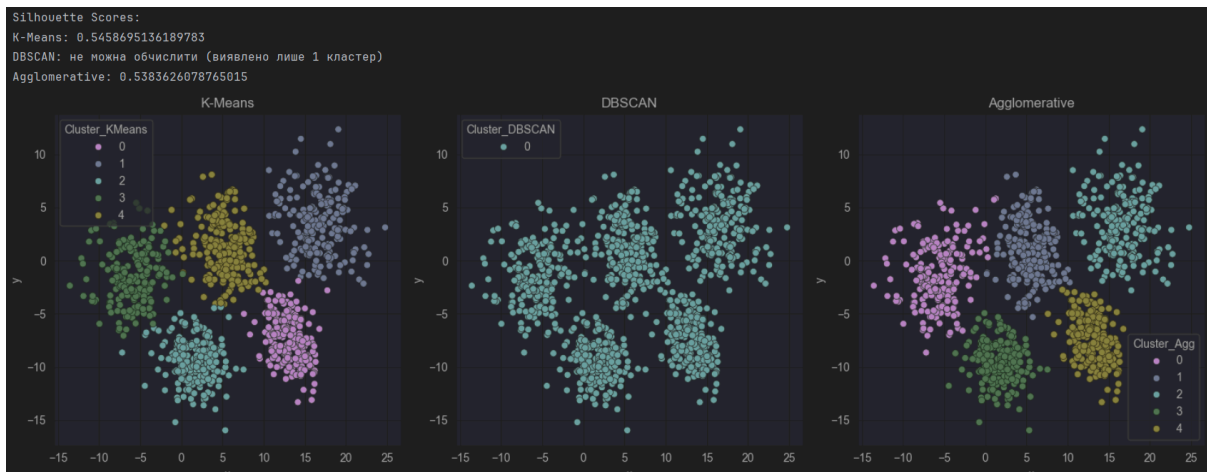
Використати Agglomerative Clustering (ієрархічну кластеризацію) та побудувати дендрограму.

```
1  agg = AgglomerativeClustering(n_clusters=5)
2  clusters_agg = agg.fit_predict(scaled_data)
3  data['Cluster_Agg'] = clusters_agg
4
5  plt.figure(figsize=(12, 6))
6  linked = linkage(scaled_data, method='ward')
7  dendrogram(linked, orientation='top', distance_sort='descending', show_leaf_counts=True)
8  plt.title('Dendrogram')
9  plt.show()
10
11 plt.figure(figsize=(10, 6))
12 sns.scatterplot(x='x', y='y', hue='Cluster_Agg', data=data, palette='viridis')
13 plt.title('Agglomerative Clustering')
14 plt.show()
```



3. Оцінка якості кластеризації о Порівняти результати кластеризації за допомогою метрики Silhouette Score о Візуалізувати отримані кластери (наприклад, scatter plot, t-SNE, PCA). о Охарактеризувати кластери. (Наприклад, кластер тих хто часто робить недорогі покупки)

```
1 print("Silhouette Scores:")
2 print(f"K-Means: {silhouette_score(scaled_data, clusters_kmeans)}")
3
4 unique_dbscan = np.unique(clusters_dbscan)
5 n_dbscan_clusters = len(unique_dbscan) - (1 if -1 in unique_dbscan else 0)
6
7 if n_dbscan_clusters > 1:
8     print(f"DBSCAN: {silhouette_score(scaled_data, clusters_dbscan)}")
9 else:
10    print("DBSCAN: не можна обчислити (виявлено лише 1 кластер)")
11
12 print(f"Agglomerative: {silhouette_score(scaled_data, clusters_agg)}")
13 plt.figure(figsize=(15, 5))
14
15 plt.subplot(1, 3, 1)
16 sns.scatterplot(x='x', y='y', hue='Cluster_KMeans', data=data, palette='viridis')
17 plt.title('K-Means')
18
19 plt.subplot(1, 3, 2)
20 sns.scatterplot(x='x', y='y', hue='Cluster_DBSCAN', data=data, palette='viridis')
21 plt.title('DBSCAN')
22
23 plt.subplot(1, 3, 3)
24 sns.scatterplot(x='x', y='y', hue='Cluster_Agg', data=data, palette='viridis')
25 plt.title('Agglomerative')
26
27 plt.tight_layout()
28 plt.show()
```



Висновок:

У цій лабораторній роботі я ознайомився з методами кластеризації даних та реалізував три алгоритми: K-Means, DBSCAN і Agglomerative Clustering. Я виконав попередню обробку даних, включаючи масштабування, визначив оптимальну кількість кластерів за допомогою "Elbow Method" та Silhouette Score, а також провів аналіз результатів кластеризації.

Порівняння алгоритмів показало, що K-Means добре працює на чітко відокремлених групах даних, DBSCAN виявляє аномалії, але чутливий до параметрів, а Agglomerative Clustering дає гнучкіший розподіл, проте може страждати на великому обсязі даних. Візуалізація кластерів дозволила оцінити якість отриманих розподілів.