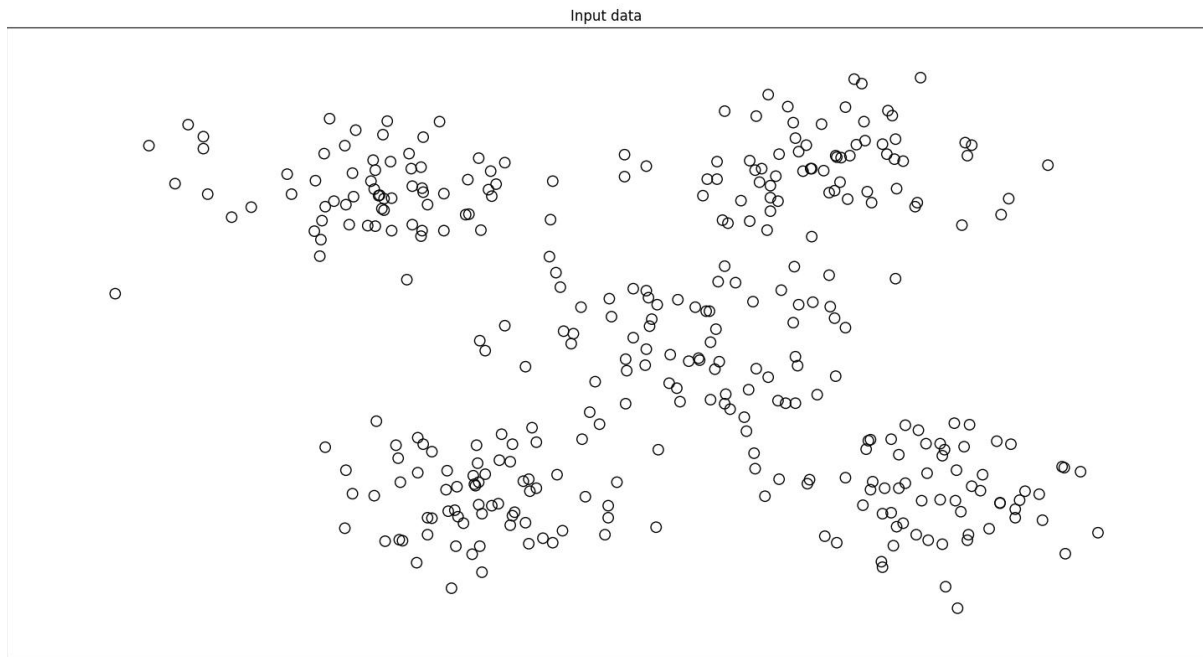


ЛАБОРАТОРНА РОБОТА № 4

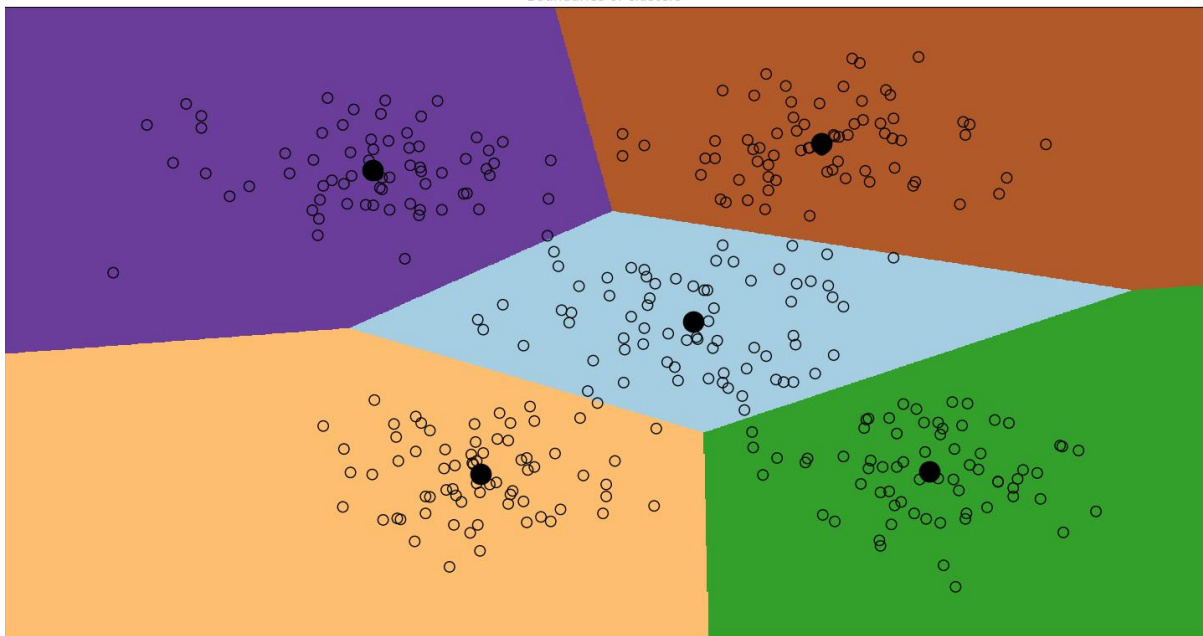
ДОСЛІДЖЕННЯ МЕТОДІВ НЕКОНТРОЛЬОВАНОГО НАВЧАННЯ

Мета роботи: використовуючи спеціалізовані бібліотеки та мову програмування Python дослідити методи неконтрольованої класифікації даних у машинному навчанні.

Завдання 1. Кластеризація даних за допомогою методу k-середніх



Boundaries of clusters



LR_7_task_1.py x

LR_7_task_1.py

```

1  import numpy as np
2  from sklearn.cluster import KMeans
3  import matplotlib.pyplot as plt
4
5  X = np.loadtxt('data_clustering.txt', delimiter=',')
6  num_clusters = 5
7
8  plt.figure()
9  plt.scatter(X[:,0], X[:,1], marker='o', facecolors='none', edgecolors='black', s=80)
10 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
11 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
12 plt.title('Input data')
13 plt.xlim(x_min, x_max)
14 plt.ylim(y_min, y_max)
15 plt.xticks(())
16 plt.yticks(())
17 plt.show()
18
19
20 kmeans = KMeans(init='k-means++', n_clusters=num_clusters, n_init=10)
21 kmeans.fit(X)
22
23 step_size = 0.01
24
25 x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
26 y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
27 x_vals, y_vals = np.meshgrid(np.arange(x_min, x_max, step_size), np.arange(y_min, y_max, step_size))
28
29 output = kmeans.predict(np.c_[x_vals.ravel(), y_vals.ravel()])
30 output = output.reshape(x_vals.shape)
31
32 plt.figure()
33 plt.clf()
34 plt.imshow(output, interpolation='nearest', extent=(x_vals.min(), x_vals.max(), y_vals.min(), y_vals.max()))
35 plt.scatter(X[:,0], X[:,1], marker='o', facecolors='none', edgecolors='black', s=80)
36 cluster_centers = kmeans.cluster_centers_
37 plt.scatter(cluster_centers[:,0], cluster_centers[:,1], marker='o', s=210, linewidths=4, color='black', zorder=5)
38 plt.title('Boundaries of clusters')
39 plt.xlim(x_min, x_max)
40 plt.ylim(y_min, y_max)
41 plt.xticks(())
42 plt.yticks(())
43 plt.show()
44

```

Завдання 2. Кластеризація К-середніх для набору даних Iris

```
LR_7_task_2.py x
LR_7_task_2.py
1 from sklearn.svm import SVC
2 from sklearn.metrics import pairwise_distances_argmin
3 from sklearn.datasets import load_iris
4 from sklearn.cluster import KMeans
5 import numpy as np
6 import matplotlib.pyplot as plt
7
8 # Load the iris dataset
9 iris = load_iris()
10 X = iris['data']
11 y = iris['target']
12
13 # Initialize KMeans with corrected parameters
14 kmeans = KMeans(n_clusters=8, init='k-means++', n_init=10, max_iter=300, tol=0.0001, verbose=0, random_stat
15
16 # Fit the model
17 kmeans.fit(X)
18
19 # Predict the clusters
20 y_kmeans = kmeans.predict(X)
21
22 # Plot the clusters
23 plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, s=50, cmap='viridis')
24 centers = kmeans.cluster_centers_
25 plt.scatter(centers[:, 0], centers[:, 1], c='black', s=200, alpha=0.5)
26 plt.show()
27
28 # Define a function to find clusters
29 def find_clusters(X, n_clusters, rseed=2):
30     rng = np.random.RandomState(rseed)
31     i = rng.permutation(X.shape[0])[:n_clusters]
32     centers = X[i]
33     while True:
34         labels = pairwise_distances_argmin(X, centers)
35         new_centers = np.array([X[labels == i].mean(0) for i in range(n_clusters)])
36         if np.all(centers == new_centers):
37             break
38         centers = new_centers
39     return centers, labels
40
41 # Find clusters using the custom function
42 centers, labels = find_clusters(X, 3)
43 plt.scatter(X[:, 0], X[:, 1], c=labels, s=50, cmap='viridis')
44 plt.show()
45
46 centers, labels = find_clusters(X, 3, rseed=0)
47 plt.scatter(X[:, 0], X[:, 1], c=labels, s=50, cmap='viridis')
48 plt.show()
49
50 # Use KMeans to predict clusters
51 labels = KMeans(3, random_state=0).fit_predict(X)
52 plt.scatter(X[:, 0], X[:, 1], c=labels, s=50, cmap='viridis')
53 plt.show()
```

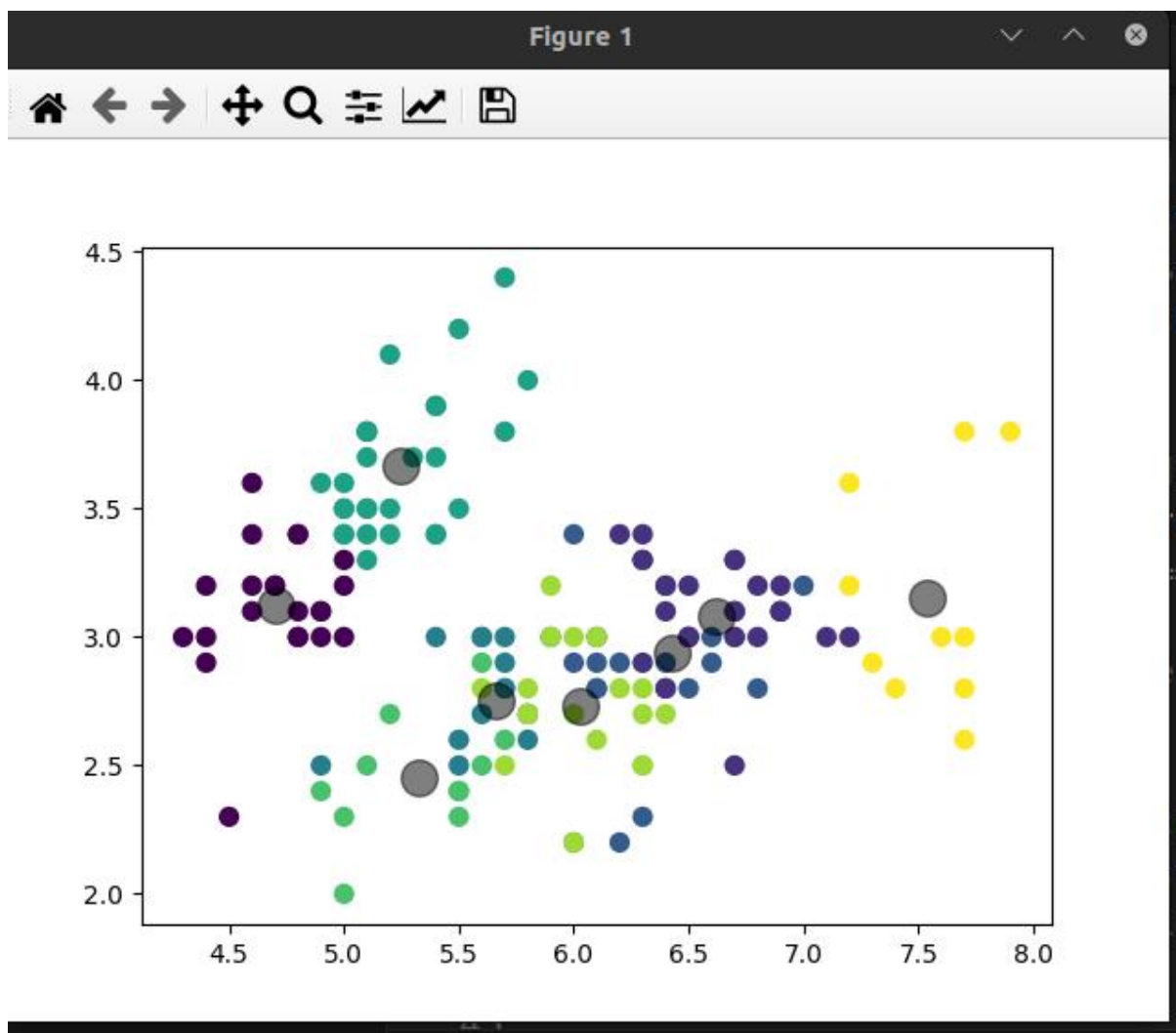
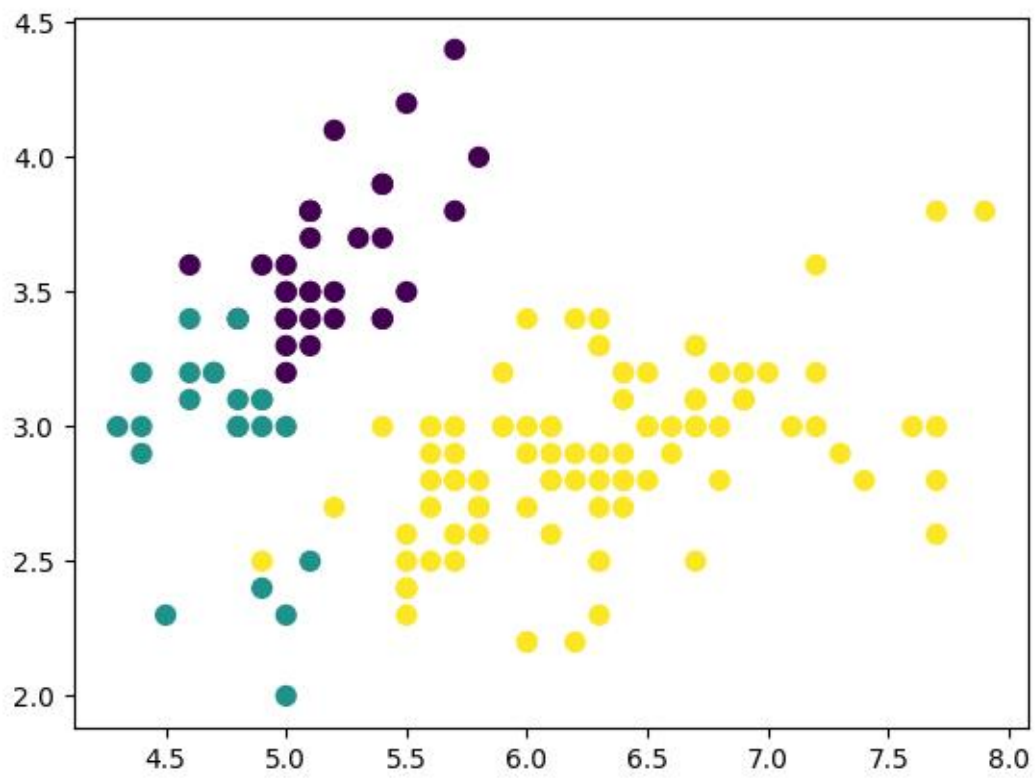
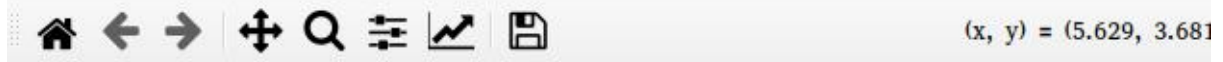
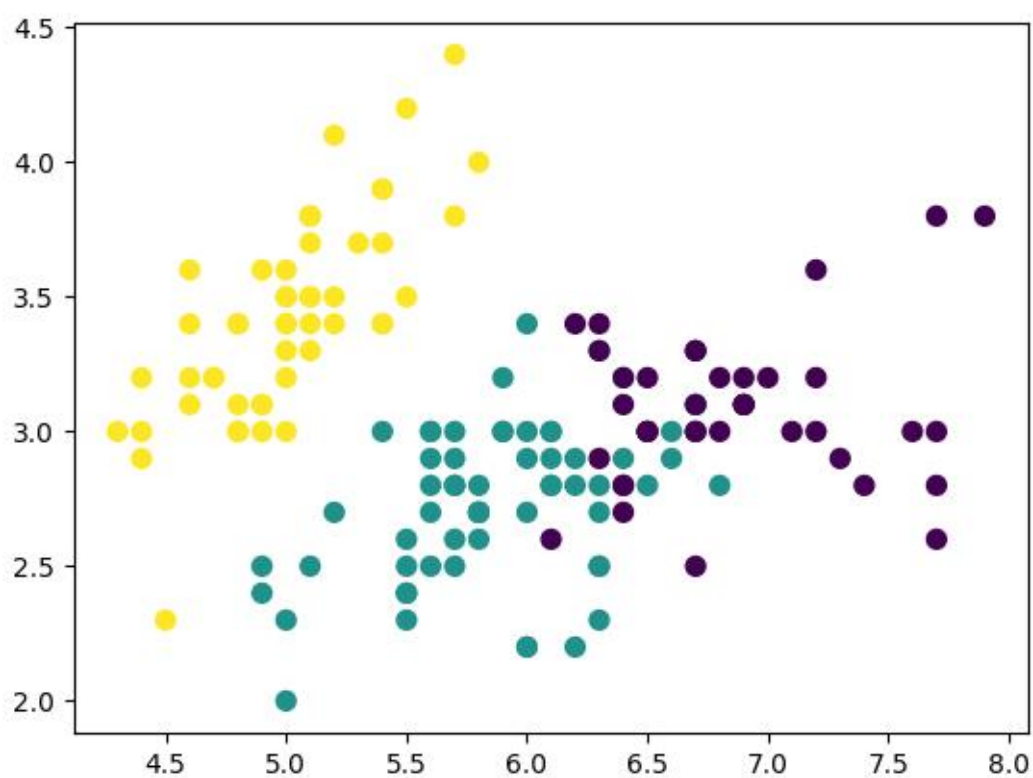
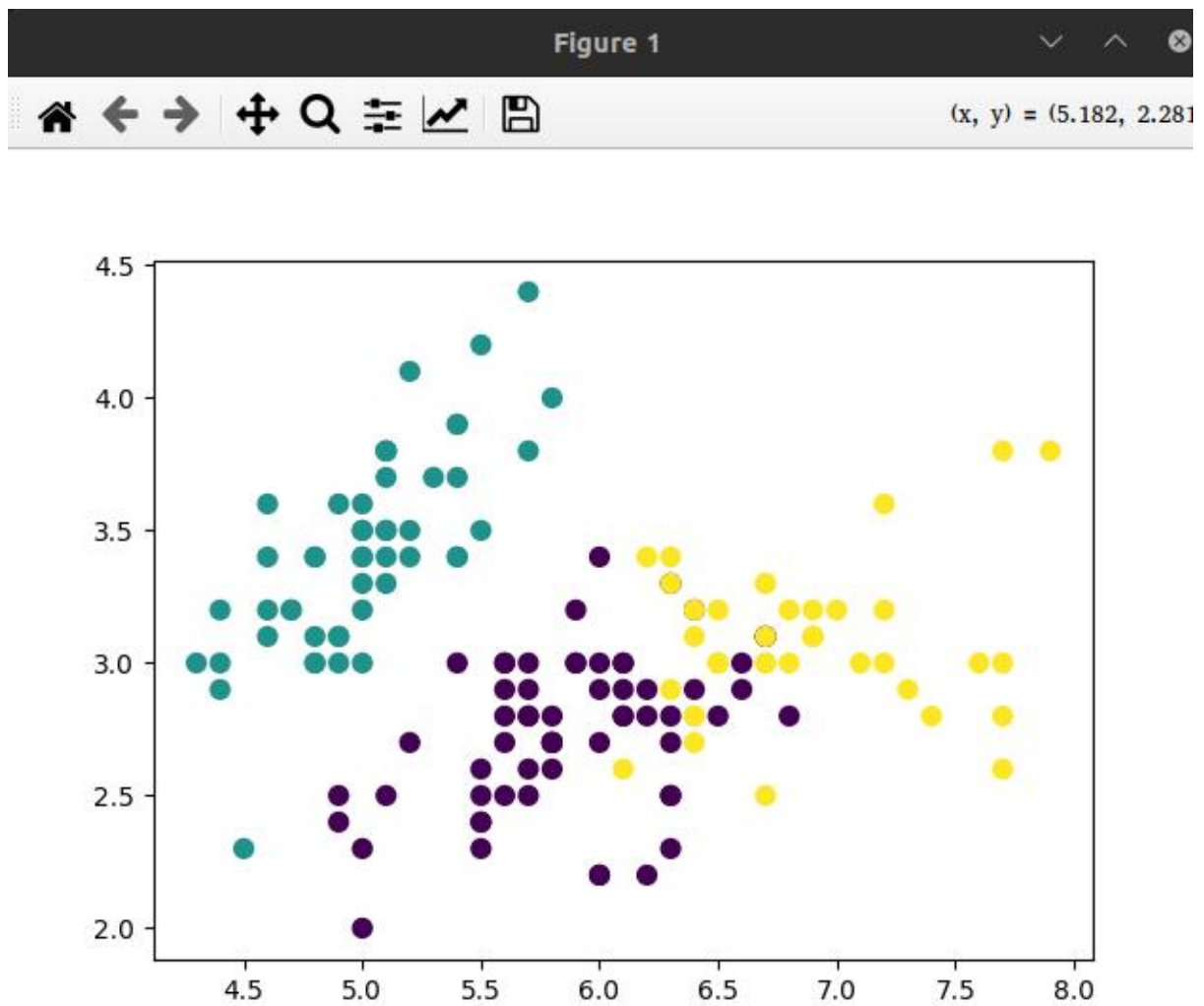


Figure 1







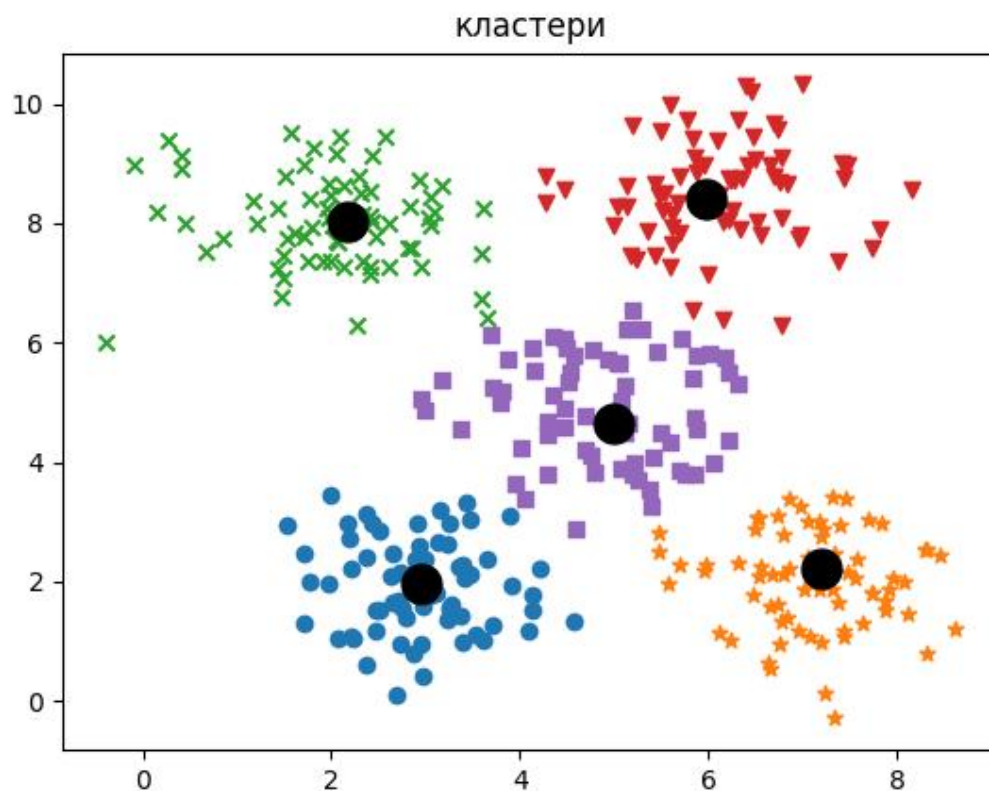
Завдання 3. Оцінка кількості кластерів з використанням методу зсуву середнього

LR_7_task_3.py ×

LR_7_task_3.py

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.cluster import MeanShift, estimate_bandwidth
4 from itertools import cycle
5
6 X = np.loadtxt('data_clustering.txt', delimiter=',')
7 bandwidth = estimate_bandwidth(X, quantile=0.1, n_samples=len(X))
8
9 ms = MeanShift(bandwidth=bandwidth, bin_seeding=True)
10 ms.fit(X)
11
12 cluster_centers = ms.cluster_centers_
13 print('Centers of clusters:\n', cluster_centers)
14
15 labels = ms.labels_
16 num_clusters = len(np.unique(labels))
17 print("Number of clusters in input data =", num_clusters)
18
19 plt.figure()
20 markers = 'o*xvs'
21 for i, marker in zip(range(num_clusters), markers):
22     plt.scatter(X[labels == i, 0], X[labels == i, 1], marker=marker)
23     cluster_centers = ms.cluster_centers_[i]
24     plt.plot(cluster_centers[0], cluster_centers[1], marker='o', markerfacecolor='black', markeredgcolor='black')
25 plt.title('кластеры')
26 plt.show()
27
```

Figure 1




```
~/Desktop/zp/ai/laba7 $ python3 LR_7_task_3.py
Centers of clusters:
[[2.95568966 1.95775862]
 [7.20690909 2.20836364]
 [2.17603774 8.03283019]
 [5.97960784 8.39078431]
 [4.99466667 4.65844444]]
Number of clusters in input data = 5
```

Завдання 4. Знаходження підгруп на фондовому ринку з використанням моделі поширення подібності

LR_7_task_4.py x

LR_7_task_4.py

```
1 import datetime
2 import json
3 import numpy as np
4 from sklearn import covariance, cluster
5 import pandas as pd
6 import yfinance as yf
7 from sklearn.impute import SimpleImputer
8
9 # Load company symbols
10 input_file = 'company_symbol_mapping.json'
11 with open(input_file, 'r') as f:
12     company_symbols_map = json.loads(f.read())
13 symbols = list(company_symbols_map.keys())
14
15 # Set date range
16 start_date = "2003-07-03"
17 end_date = "2007-05-04"
18
19 # Load stock data
20 quotes = {}
21 for symbol in symbols:
22     try:
23         print(f>Loading {symbol} ({company_symbols_map[symbol]})...", end='')
24         q = yf.download(symbol, start=start_date, end=end_date)
25         if not q.empty:
26             quotes[symbol] = q
27             print("done.")
28         else:
29             print("no data.")
30     except Exception as e:
31         print(f"error: {e}")
```

LR_7_task_4.py x

LR_7_task_4.py

```
32
33 # Align all data to the same index (date range)
34 aligned_quotes = pd.concat(quotes.values(), keys=quotes.keys(), names=["Symbol", "Date"])
35 aligned_quotes = aligned_quotes.unstack(level=0) # Organize by symbol
36
37 # Extract opening and closing quotes
38 opening_quotes = aligned_quotes['Open'].values # Shape: (time, stocks)
39 closing_quotes = aligned_quotes['Close'].values # Shape: (time, stocks)
40
41 # Compute differences
42 quotes_diff = closing_quotes - opening_quotes
43
44 # Handle missing values
45 imputer = SimpleImputer(strategy='mean') # Replace NaN with column mean
46 quotes_diff_imputed = imputer.fit_transform(quotes_diff)
47
48 # Normalize data
49 X = quotes_diff_imputed.copy()
50 X /= X.std(axis=0) # Normalize each stock's data
51
52 # Graphical Lasso for covariance estimation
53 edge_model = covariance.GraphicalLassoCV(cv=3)
54
55 # Train model
56 with np.errstate(invalid='ignore'):
57     edge_model.fit(X)
58
59 # Perform clustering
60 _, labels = cluster.affinity_propagation(edge_model.covariance_)
61 num_labels = labels.max()
62
63 # Output clustering results
```

```

58
59 # Perform clustering
60 _, labels = cluster.affinity_propagation(edge_model.covariance_)
61 num_labels = labels.max()
62
63 # Output clustering results
64 names = np.array([company_symbols_map[symbol] for symbol in quotes.keys()])
65 for i in range(num_labels + 1):
66     print("Cluster", i+1, "==>", ', '.join(names[labels == i]))
67
Cluster 1 ==> Exxon, Chevron, ConocoPhillips, Valero Energy
Cluster 2 ==> Toyota, Ford, Honda, Boeing, Mc Donalds, Apple, SAP, Caterpillar
Cluster 3 ==> Kraft Foods
Cluster 4 ==> Coca Cola, Pepsi, Kellogg, Procter Gamble, Colgate-Palmolive, Kimberly-Clark
Cluster 5 ==> Time Warner, Comcast, Marriott, Wells Fargo, JPMorgan Chase, AIG, American express, Bank of America, Goldman
Sachs, Xerox, Wal-Mart, Home Depot, Ryder, DuPont de Nemours
Cluster 6 ==> Microsoft, IBM, HP, Amazon, 3M, General Electrics, Cisco, Texas instruments
Cluster 7 ==> Northrop Grumman, Lockheed Martin, General Dynamics
Cluster 8 ==> Walgreen, CVS
Cluster 9 ==> GlaxoSmithKline, Pfizer, Sanofi-Aventis, Novartis

```

Посилання на GitHub: https://github.com/missShevel/SHI_Shevel_Olha_IPZ-21-1/tree/master/Lab7