

ЛАБОРАТОРНА РОБОТА №4

ДОСЛІДЖЕННЯ МЕТОДІВ НЕКОНТРОЛЬОВАНОГО НАВЧАННЯ

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

Хід роботи:

Завдання 1. Кластеризація даних за допомогою методу k-середніх

Провести кластеризацію даних методом k-середніх. Використовувати файл вхідних даних: data_clustering.txt.

Лістинг програми:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn import metrics

# Завантаження вхідних даних
X = np.loadtxt('data_clustering.txt', delimiter=',')

# кількість кластерів
num_clusters = 5

# Включення вхідних даних до графіка
plt.figure()
plt.scatter(X[:, 0], X[:, 1], marker='o', facecolors='none', edgecolors='black',
            s=80)
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
plt.title('Вхідні дані')
plt.xlim(x_min, x_max)
plt.ylim(y_min, y_max)
plt.xticks(())
plt.yticks(())
plt.show()

# Створення об'єкту KMeans
kmeans = KMeans(init='k-means++', n_clusters=num_clusters, n_init=10)

# Навчання моделі кластеризації KMeans
kmeans.fit(X)

# Визначення кроку сітки
step_size = 0.01

# Відображення точок сітки
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
```

					Житомирська Політехніка.22.121.15.000 – Лр04					
Змн.	Арк.	№ докум.	Підпис	Дата						
Розроб.		Сівченко О. О.			Звіт з лабораторної роботи			Літ.	Арк.	Аркушів
Перевір.		Філіпов В. О.							1	8
Керівник								ФІКТ Гр. ПІ-60[2]		
Н. контр.										
Зав. каф.										

```

x_vals, y_vals = np.meshgrid(np.arange(x_min, x_max, step_size), np.arange(y_min,
y_max, step_size))

# Передбачення вихідних міток для всіх точок сітки
output = kmeans.predict(np.c_[x_vals.ravel(), y_vals.ravel()])

# Графічне відображення областей та виділення їх кольором
output = output.reshape(x_vals.shape)
plt.figure()
plt.clf()
plt.imshow(output, interpolation='nearest', extent=(x_vals.min(), x_vals.max(),
y_vals.min(), y_vals.max()),
           cmap=plt.cm.Paired, aspect='auto', origin='lower')

# Відображення вхідних точок
plt.scatter(X[:, 0], X[:, 1], marker='o', facecolors='none', edgecolors='black',
s=80)

# Відображення центрів кластерів
cluster_centers = kmeans.cluster_centers_
plt.scatter(cluster_centers[:, 0], cluster_centers[:, 1], marker='o', s=210,
           linewidths=4, color='black', zorder=12,
           facecolors='black')
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
plt.title('Межі кластерів')
plt.xlim(x_min, x_max)
plt.ylim(y_min, y_max)
plt.xticks(())
plt.yticks(())
plt.show()

```

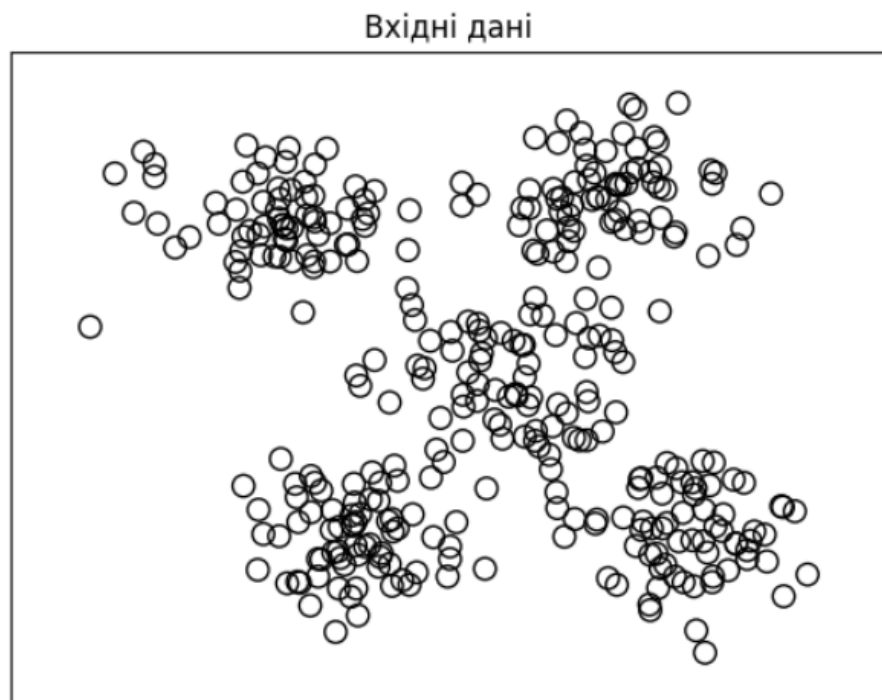


Рисунок 1.1 – Результат роботи програми у вигляді графіку

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр04	Арк.
		Філіпов В. О.				2
Змн.	Арк.	№ докум.	Підпис	Дата		

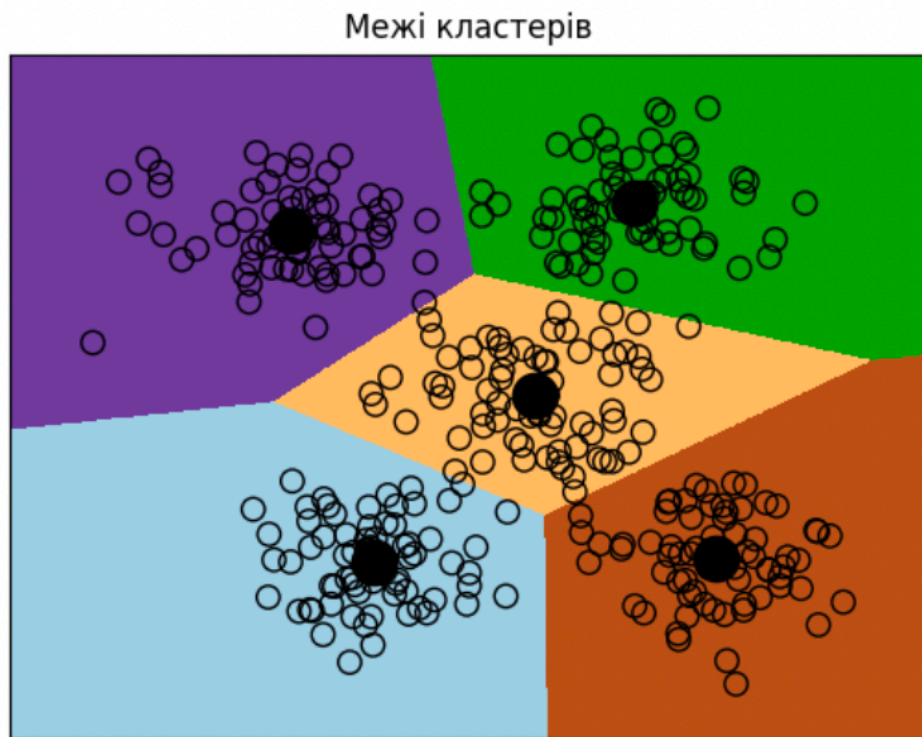


Рисунок 1.2 – Результат роботи програми у вигляді графіку

Метод k -середніх (k -means) - це добре відомий алгоритм кластеризації. Його використання передбачає, що кількість кластерів заздалегідь відома. Далі ми сегментуємо дані до підгруп, застосовуючи різні атрибути даних. Ми починаємо з того, що фіксуємо кількість кластерів та, виходячи з цього, класифікуємо дані. Основна ідея полягає в оновленні положень центроїдів (центрів тяжіння кластеру, або головні точки) на кожній ітерації. Ітеративний процес продовжується до тих пір, поки всі центроїди не займуть оптимального положення.

Завдання 2. Кластеризація K -середніх для набору даних Iris

Виконайте кластеризацію K -середніх для набору даних Iris, який включає три типи (класи) квітів ірису (Setosa, Versicolour і Virginica) з чотирма атрибутами: довжина чашолистка, ширина чашолистка, довжина пелюстки та ширина пелюстки. У цьому завданні використовуйте `sklearn.cluster.KMeans` для пошуку кластерів набору даних Iris.

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр04	Арк.
		Філіпов В. О.				3
Змн.	Арк.	№ докум.	Підпис	Дата		

Лістинг програми:

```
import sklearn
from sklearn.svm import SVC
from sklearn.metrics import pairwise_distances_argmin
from sklearn.datasets import load_iris
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
import numpy as np

iris = load_iris()
X = iris['data']
y = iris['target']

# Налаштування класу KMeans
KMeans(n_clusters=8, init='k-means++', n_init=0, max_iter=300,
       tol=0.0001, verbose=0, random_state=None, copy_x=True, algorithm='auto')

# Створення об'єкту KMeans та визначення кількості кластерів
kmeans = KMeans(n_clusters=3)

# Навчання моделі кластеризації KMeans
kmeans.fit(X)

# Передбачення вихідних міток для всіх точок
y_kmeans = kmeans.predict(X)

# Відображення центрів кластерів
plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, s=50, cmap='viridis')
centers = kmeans.cluster_centers_
plt.scatter(centers[:, 0], centers[:, 1], c='black', s=200, alpha=0.5)

# Знаходження кластерів
def find_clusters(X, n_clusters, rseed=2):
    # Випадково обираємо кластери
    rng = np.random.RandomState(rseed)
    i = rng.permutation(X.shape[0])[:n_clusters]
    centers = X[i]
    while True:
        # Призначаємо мітки на основі найближчого центру
        labels = pairwise_distances_argmin(X, centers)
        # Знаходимо нові центри за допомогою точок
        new_centers = np.array([X[labels == i].mean(0)
                               for i in range(n_clusters)])
        # Перевіряємо центри на збіжність
        if np.all(centers == new_centers):
            break
        centers = new_centers
    return centers, labels

# Знаходження центрів кластерів та зображення їх на графіку
centers, labels = find_clusters(X, 3)
plt.scatter(X[:, 0], X[:, 1], c=labels, s=50, cmap='viridis')

# Передбачення вхідних точок кластерів та зображення їх на графіку
labels = KMeans(3, random_state=0).fit_predict(X)
plt.scatter(X[:, 0], X[:, 1], c=labels, s=50, cmap='viridis')
plt.show()
```

		Сівченко О. О.			Житомирська Політехніка. 22.121.15.000 – Пр04	Арк.
		Філіпов В. О.				4
Змн.	Арк.	№ докум.	Підпис	Дата		

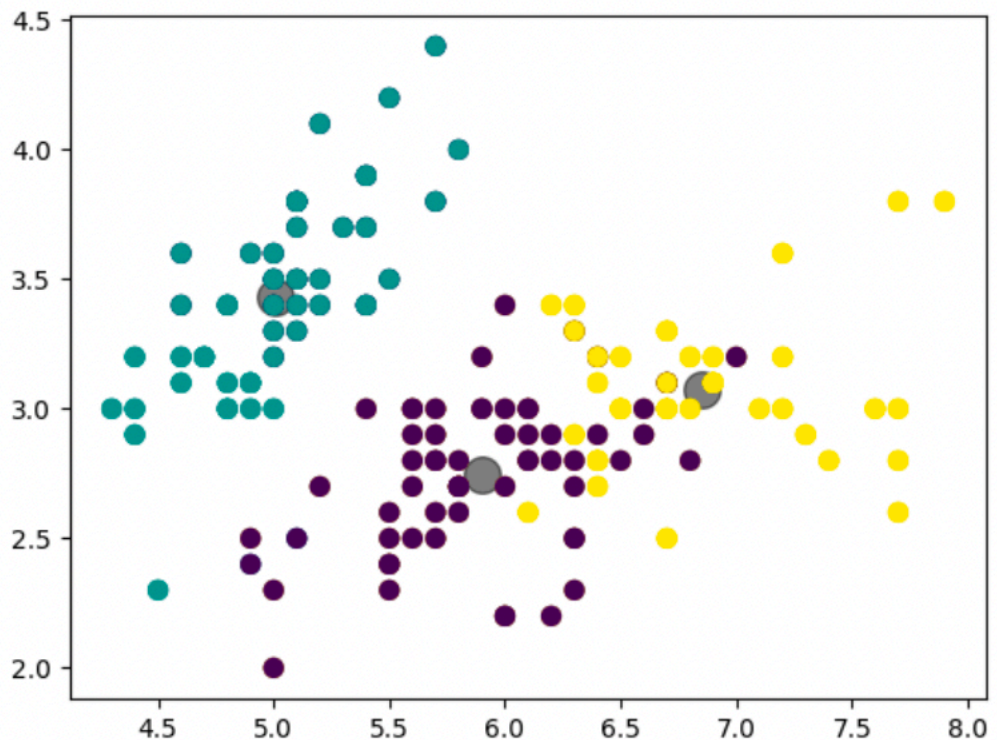


Рисунок 1.3 – Результат роботи програми у вигляді графіку

Завдання 3. Оцінка кількості кластерів з використанням методу зсуву середнього

Відповідно до рекомендацій, напишіть програму та оцініть максимальну кількість кластерів у заданому наборі даних за допомогою алгоритму зсуву середнього. Для аналізу використовуйте дані, які містяться у файлі data_clustering.txt.

Лістинг програми:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import MeanShift, estimate_bandwidth
from itertools import cycle

# Завантаження вхідних даних
X = np.loadtxt('data_clustering.txt', delimiter=',')

# Оцінка ширини вікна для X
bandwidth_X = estimate_bandwidth(X, quantile=0.1, n_samples=len(X))

# Кластеризація даних методом зсуву середнього
meanshift_model = MeanShift(bandwidth=bandwidth_X, bin_seeding=True)
meanshift_model.fit(X)

# Витягування центрів кластерів
```

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр04	Арк.
		Філіпов В. О.				5
Змн.	Арк.	№ докум.	Підпис	Дата		

```

cluster_centers = meanshift_model.cluster_centers_
print('\nCenters of clusters: \n', cluster_centers)

# Оцінка кількості кластерів
labels = meanshift_model.labels_
num_clusters = len(np.unique(labels))
print("\nNumber of clusters in input data =", num_clusters)

# Відображення на графіку точок та центрів кластерів
plt.figure()
markers = 'o*xvs'
for i, marker in zip(range(num_clusters), markers):
    # Отображение на графике точек, принадлежащих
    # текущему кластеру
    plt.scatter(X[labels == i, 0], X[labels == i, 1], marker=marker,
                color='black')
    # Відображення на графіку центру кластера
    cluster_center = cluster_centers[i]
    plt.plot(cluster_center[0], cluster_center[1], marker='o',
             markerfacecolor='black', markeredgecolor='black',
             markersize=15)
plt.title('Clusters')
plt.show()

```

```

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

Process finished with exit code 0

```

Рисунок 1.4 – Результат роботи програми

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр04	Арк.
		Філіпов В. О.				
Змн.	Арк.	№ докум.	Підпис	Дата		6

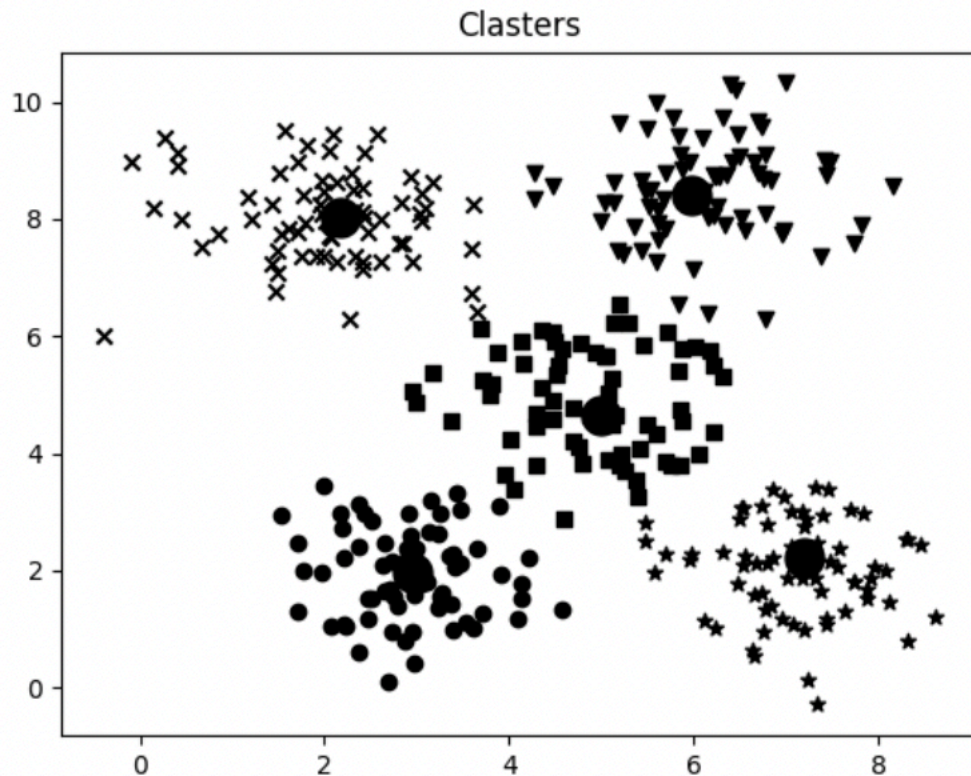


Рисунок 1.5 – Результат роботи програми у вигляді графіку

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

Використовуючи модель поширення подібності, знайти підгрупи серед учасників фондового ринку. У якості керуючих ознак будемо використовувати варіацію котирувань між відкриттям і закриттям біржі. Використовувати файл вхідних даних фондового ринку, що доступний в бібліотеці `matplotlib`. Прив'язки символічних позначень компаній до повних назв містяться у файлі `company_symbol_mapping.json`.

Лістинг програми:

```
import datetime
import json
import numpy as np
import matplotlib.pyplot as plt
from sklearn import covariance, cluster
from mplfinance import quotes_historical_yahoo_ohlc as quotes_yahoo

# Вхідний файл із символічними позначеннями компаній
input_file = 'company_symbol_mapping.json'

# Завантаження прив'язок символів компаній до їх повних назв
```

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр04	Арк.
		Філіпов В. О.				7
Змн.	Арк.	№ докум.	Підпис	Дата		

```

with open(input_file, 'r') as f:
    company_symbols_map = json.loads(f.read())
    symbols, names = np.array(list(company_symbols_map.items())).T

# Завантаження архівних даних котирувань
start_date = datetime.datetime(2003, 7, 3)
end_date = datetime.datetime(2007, 5, 4)
quotes = [quotes_yahoo(symbol, start_date, end_date, asobject=True) for symbol in
symbols]

# Вилучення котирувань, що відповідають
# відкриттю та закриттю біржі
opening_quotes = np.array([quote.open for quote in quotes]).astype(np.float)
closing_quotes = np.array([quote.close for quote in quotes]).astype(np.float)

# Обчислення різниці між двома видами котирувань
quotes_diff = closing_quotes - opening_quotes

# Нормалізація даних
X = quotes_diff.copy().T
X /= X.std(axis=0)

# Створення моделі графа
edge_model = covariance.GraphicalLassoCV()

# Навчання моделі
with np.errstate(invalid='ignore'):
    edge_model.fit(X)

# Створення моделі кластеризації на основі поширення подібності
_, labels = cluster.affinity_propagation(edge_model.covariance_)
num_labels = labels.max()

# Виведення результатів
for i in range (num_labels + 1):
    print("cluster", i+1, "==>", ", ".join(names(labels == i)))

```

```

from mplfinance import quotes_historical_yahoo_ohlc as quotes_yahoo
ImportError: cannot import name 'quotes_historical_yahoo_ohlc' from 'mplfinance' (/Library/Frameworks/Python.framework
Process finished with exit code 1

```

Рисунок 1.6 – Результат роботи програми

Функцію `quotes_historical_yahoo_ochl` було повністю видалено з коду бібліотеки `mpl_finance`, тому завантажити її не вдалося.

Посилання на GitHub: <https://github.com/SoylerProfile/SHI>

Висновки: в ході виконання лабораторної роботи було досліджено методи неконтрольованої класифікації даних у машинному навчанні, використовуючи спеціалізовані бібліотеки та мову програмування Python.

		Сівченко О. О.			Житомирська Політехніка.22.121.15.000 – Лр04	Арк.
		Філіпов В. О.				8
Змн.	Арк.	№ докум.	Підпис	Дата		