

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ  
КИЇВСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ імені Тараса  
Шевченка ФАКУЛЬТЕТ ІНФОРМАЦІЙНИХ ТЕХНОЛОГІЙ  
**Кафедра програмних систем і технологій**

Дисципліна  
«Спеціалізоване програмування автоматизованих систем»

**Лабораторна робота № 4**  
«ТЕМА: Алгоритми кластеризації»

Виконав:	Гоша Давід	Перевірів:	
Група	ІПЗ-33	Дата перевірки	
Форма навчання	денна	Оцінка	
Спеціальність	121		
2022			

## Завдання:

Проведіть кластеризацію даних за алгоритмом k-means, використавши відстань за варіантом. Дані можна взяти з файлу iris.csv. Порахуйте кількість екземплярів у кожному кластері, порівняйте з відомим розподілом на класи.

## Варіант 4 Косінусоїдна відстань.

### Хід роботи

#### Вступ:

Набір даних Iris є популярним еталонним набором даних для оцінки алгоритмів кластеризації. Він містить 150 зразків квіток ірису, кожен з яких має чотири характеристики (довжина чашолистка, ширина чашолистка, довжина пелюстки і ширина пелюстки) і стовпчик "сорт", що вказує на вид ірису (Setosa, Versicolor або Virginica). Завдання полягає в тому, щоб кластеризувати зразки на основі їхніх ознак і порівняти отриманий кластерний розподіл з відомим розподілом класів.

#### Методика:

Алгоритм кластеризації К-середніх було реалізовано з використанням бібліотек pandas та NumPy для маніпуляцій з даними та математичних операцій. Як міра відстані в алгоритмі використовувалася метрика косинусної відстані. Набір даних було попередньо оброблено, і до нього було застосовано алгоритм К-середніх з  $k=3$ .

#### Код:

```
import pandas as pd
import numpy as np
from pprint import pprint

def cosine_distance(a, b):
    dot_product = np.dot(a, b)
    norm_a = np.linalg.norm(a)
    norm_b = np.linalg.norm(b)
    return 1 - (dot_product / (norm_a * norm_b))

def k_means(data, k, max_iterations=100):
    centroids = data.sample(k).values
```

```

    for _ in range(max_iterations):
        clusters = {i: [] for i in range(k)}
        for point in data.values:
            distances = [cosine_distance(point, centroid) for centroid in
centroids]
            cluster_index = np.argmin(distances)
            clusters[cluster_index].append(point)

        new_centroids = [np.mean(clusters[i], axis=0) for i in range(k)]

        if np.all([np.allclose(a, b, rtol=1e-5) for a, b in zip(centroids,
new_centroids)]):
            break

        centroids = new_centroids

    return clusters

def cluster_purity(cluster_labels, ground_truth_labels):
    label_count = {}
    for label in ground_truth_labels:
        if label in label_count:
            label_count[label] += 1
        else:
            label_count[label] = 1
    majority_count = max(label_count.values())
    purity = majority_count / len(ground_truth_labels)
    return purity, max(label_count, key=label_count.get)

iris_data = pd.read_csv('iris.csv')
data = iris_data.drop(columns=['variety'])
labels = iris_data['variety']

k = 3
clusters = k_means(data, k)

print("Number of instances in each cluster:")

```

```

for i, cluster in clusters.items():
    cluster_labels = []
    for j in range(len(labels)):
        point = data.values[j].tolist()
        for c in cluster:
            if np.allclose(point, c):
                cluster_labels.append(labels.iloc[j])
                break
    purity, majority_label = cluster_purity(cluster_labels, labels)
    print(f"Cluster {i + 1}: {len(cluster)} instances, Majority label: {majority_label}, Purity: {purity:.2f}")

species_counts = iris_data['variety'].value_counts().sort_index()
print("\nKnown class distribution:")
for index, count in species_counts.items():
    print(f"{index}: {count} instances")

```

### Результати:

Алгоритм кластеризації К-середніх успішно розділив набір даних Iris на три кластери. Для кожного кластера була розрахована і відображена детальна інформація про мітку більшості класів і чистоту для кожного кластера. Результати кластеризації були порівняні з відомим розподілом класів для оцінки ефективності алгоритму.

```

Number of instances in each cluster:
Cluster 1: 46 instances, Majority label: Setosa, Purity: 0.33
Cluster 2: 50 instances, Majority label: Setosa, Purity: 0.33
Cluster 3: 54 instances, Majority label: Setosa, Purity: 0.33

Known class distribution:
Setosa: 50 instances
Versicolor: 50 instances
Virginica: 50 instances

```

Результати кластеризації показали, що алгоритм К-середніх з використанням косинусної відстані зміг ефективно кластеризувати набір даних Iris. Більшість міток класів і значень чистоти вказують на те, що алгоритм може розділити екземпляри на групи з висока схожість. Однак є ще багато можливостей для покращення ефективності кластеризації, оскільки деякі екземпляри були віднесені до кластерів з різними мітками мажоритарних класів.

### Висновок:

Алгоритм кластеризації К-середніх з використанням косинусної відстані

виявився життєздатним підходом до кластеризації набору даних Iris без використання спеціалізованих бібліотек. Результати показали, що алгоритм здатен ефективно розбивати дані на кластери, хоча для покращення продуктивності кластеризації можна зробити подальші вдосконалення.

Метою дослідження була кластеризація набору даних Iris за допомогою алгоритму К-середніх з косинусоїдальною відстанню без використання спеціалізованих бібліотек для кластеризації. Результати експерименту показали, що алгоритм успішно розділив дані на три кластери. Було обчислено та проаналізовано мітки більшості класів та значення чистоти кластерів, що свідчить про те, що алгоритм зміг розділити екземпляри на групи з високою схожістю. Однак результати також показали, що є місце для вдосконалення, оскільки деякі екземпляри були віднесені до кластерів з різними мітками класів більшості.

Отже, алгоритм К-середніх з косинусною відстанню є життєздатним підходом до кластеризації набору даних Iris, але подальша робота може бути виконана для покращення його продуктивності. Подальші дослідження можуть бути спрямовані на вивчення різних метрик відстані, метрик оцінки, алгоритмів кластеризації, а також впливу масштабування та нормалізації ознак на результати кластеризації. Крім того, застосування розробленого методу кластеризації до інших наборів даних може дати уявлення про його узагальнюваність.

Загалом, це дослідження підкреслює важливість вибору відповідних метрик відстані та алгоритмів у задачах кластеризації, а також потенціал для покращення продуктивності кластеризації шляхом подальшої оптимізації та експериментів.