

# Лабораторна робота №6. Кластеризація

Виконав студент групи КМ-91мп

Галета М.С.

## Завдання на лабораторну роботу

1. Використовуються набори даних з лабораторної роботи №1.
2. Обравши всі стовпчики, крім останнього:
  - a) Провести розбиття методом K-means на два та три кластери. Проілюструвати результати розбиття графічно.
  - b) Провести розбиття на аналогічну кількість кластерів за допомогою інших двох методів (один з яких використовує аналогічну метрику, а другий – іншу). Результати розбиття проілюструвати графічно, виконати порівняльний аналіз отриманих результатів із результатами попереднього пункту.
  - c) Збудувати дендрограму для перших 100 записів.
3. Проаналізувати отримані результати, враховуючи останній стовпчик початкових даних.
4. Для отримання результатів можна використовувати бібліотеки для мови Python (наприклад, scikit-learn або аналогічні).

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns; sns.set()
5 import plotly.figure_factory as ff
6 from sklearn.cluster import KMeans, AgglomerativeClustering
7 from sklearn.mixture import GaussianMixture
8 from sklearn.decomposition import PCA
9 from sklearn.metrics import mean_squared_error
```

## Зчитування датасету

In [2]:

```
1 df = pd.read_csv('MP-04-Galeta.csv', delimiter=';', names=['x1', 'x2', 'x3', 'x4', 'x5']
2
3 X = df.drop(columns='y')
4 y = df['y']
```

## Центрування і стандартизація даних

$$X_{new} = \frac{X - \mu}{\sigma}$$

In [3]:

```
1 class Scaler:
2     def __init__(self):
3         self.mean = None
4         self.std = None
5
6     def fit(self, X):
7         self.mean = np.mean(X, axis=0, keepdims=True)
8         self.std = np.std(X, axis=0, keepdims=True)
9
10    def transform(self, X):
11        X_new = (X - self.mean)/self.std
12        return X_new
13
14    def fit_transform(self, X):
15        self.mean = np.mean(X, axis=0, keepdims=True)
16        self.std = np.std(X, axis=0, keepdims=True)
17        X_new = (X - self.mean)/self.std
18        return X_new
```

In [4]:

```
1 sc = Scaler()
2 X = pd.DataFrame(columns=['x1', 'x2', 'x3', 'x4', 'x5', 'x6'], data=sc.fit_transform(X
```

## KMeans

In [5]:

```
1 km2 = KMeans(2, random_state=42)
2 km3 = KMeans(3, random_state=42)
3
4 km2.fit(X)
5 km3.fit(X)
```

Out[5]:

```
KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
       n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto',
       random_state=42, tol=0.0001, verbose=0)
```

In [6]:

```

1 km2_res = km2.predict(X)
2 km3_res = km3.predict(X)
3
4 pca_embeddings = PCA(2).fit_transform(X)

```

In [7]:

```

1 f, axes = plt.subplots(1, 2, figsize=(12, 6), sharex=True)
2 sns.scatterplot(x=pca_embeddings[:, 0], y=pca_embeddings[:, 1], hue=km2_res, ax=axes[0])
3 sns.scatterplot(x=pca_embeddings[:, 0], y=pca_embeddings[:, 1], hue=km3_res, ax=axes[1])
4 axes[0].set_title('KMeans 2 clusters')
5 axes[1].set_title('KMeans 3 clusters')
6 plt.show()

```



## AgglomerativeClustering & GaussianMixture

In [8]:

```

1 agg2 = AgglomerativeClustering(2)
2 agg3 = AgglomerativeClustering(3)
3 gaus2 = GaussianMixture(2, random_state=42)
4 gaus3 = GaussianMixture(3, random_state=42)

```

In [9]:

```

1 agg2_res = agg2.fit_predict(X)
2 agg3_res = agg3.fit_predict(X)
3 gaus2_res = gaus2.fit_predict(X)
4 gaus3_res = gaus3.fit_predict(X)

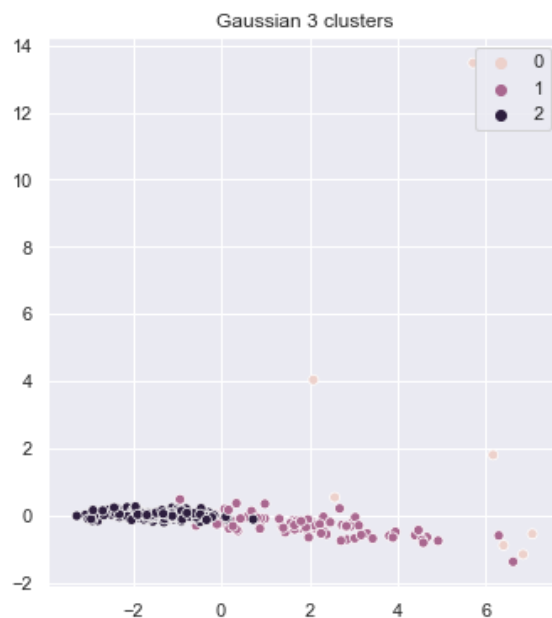
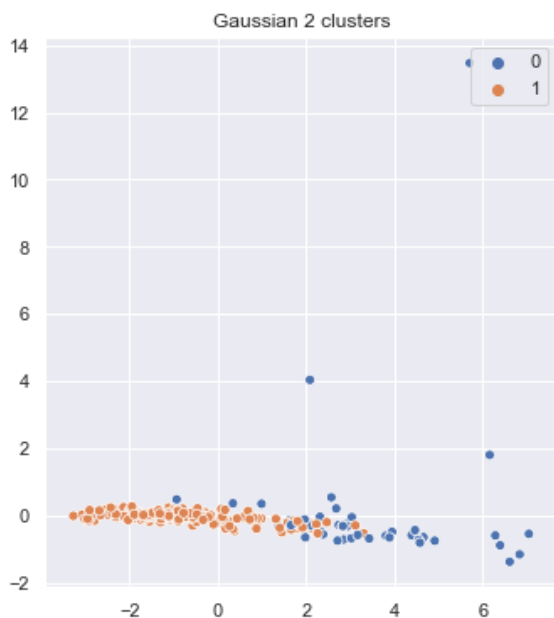
```

In [10]:

```

1 f, axes = plt.subplots(1, 2, figsize=(12, 6), sharex=True)
2 sns.scatterplot(x=pca_embeddings[:, 0], y=pca_embeddings[:, 1], hue=agg2_res, ax=axes[0])
3 sns.scatterplot(x=pca_embeddings[:, 0], y=pca_embeddings[:, 1], hue=agg3_res, ax=axes[1])
4 axes[0].set_title('Agglomerative 2 clusters')
5 axes[1].set_title('Agglomerative 3 clusters')
6
7 f, axes = plt.subplots(1, 2, figsize=(12, 6), sharex=True)
8 sns.scatterplot(x=pca_embeddings[:, 0], y=pca_embeddings[:, 1], hue=gaus2_res, ax=axes[0])
9 sns.scatterplot(x=pca_embeddings[:, 0], y=pca_embeddings[:, 1], hue=gaus3_res, ax=axes[1])
10 axes[0].set_title('Gaussian 2 clusters')
11 axes[1].set_title('Gaussian 3 clusters')
12 plt.show()

```



In [11]:

```

1 res_list = [
2     km2_res, agg2_res, gaus2_res,
3     km3_res, agg3_res, gaus3_res
4 ]

```

In [12]:

```

1 mean_dist = []
2 dist_matrix = np.mean(np.square(X.values - X.values[:, np.newaxis]), axis=2)
3 for res in res_list:
4     df_res = pd.concat((X, pd.DataFrame({'res':res})), axis=1)
5     dist_mean_dict = {}
6     for cluster in np.unique(df_res['res']):
7         df_cluster_index = df_res[df_res['res'] == cluster].drop('res', axis=1).index
8         dist_mean = []
9         for inx in df_cluster_index:
10             dist_mean.append(dist_matrix[df_cluster_index, inx].mean())
11         dist_mean_dict[cluster] = np.mean(dist_mean)
12     mean_dist.append(dist_mean_dict)

```

In [13]:

```

1 for inx, val in enumerate(mean_dist):
2     print(f'Метод {inx}, кількість кластерів {len(val)}, середня внутрішньокласова відс

```

Метод 0, кількість кластерів 2, середня внутрішньокласова відстань {0: 0.39506906746054954, 1: 2.329417115263443}, середня внутрішньокласова відстань по методу 1.3622430913619963

Метод 1, кількість кластерів 2, середня внутрішньокласова відстань {0: 2.474927845553782, 1: 0.48936205043248404}, середня внутрішньокласова відстань по методу 1.4821449479931332

Метод 2, кількість кластерів 2, середня внутрішньокласова відстань {0: 3.4811839979965566, 1: 0.7773732184660165}, середня внутрішньокласова відстань по методу 2.1292786082312865

Метод 3, кількість кластерів 3, середня внутрішньокласова відстань {0: 0.39506906746054954, 1: 1.4215683411463262, 2: 0.0}, середня внутрішньокласова відстань по методу 0.6055458028689585

Метод 4, кількість кластерів 3, середня внутрішньокласова відстань {0: 0.48936205043248404, 1: 1.426270699623292, 2: 0.0}, середня внутрішньокласова відстань по методу 0.6385442500185919

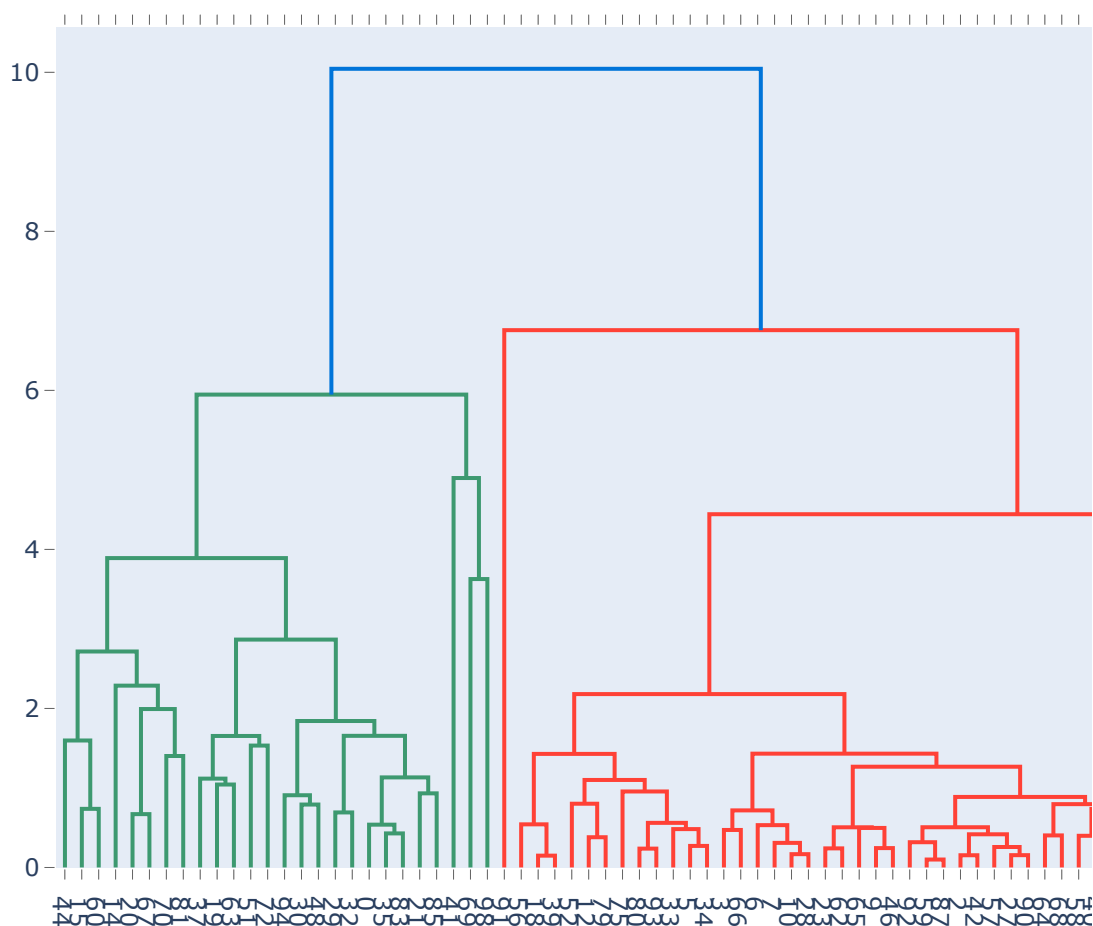
Метод 5, кількість кластерів 3, середня внутрішньокласова відстань {0: 10.045770369590437, 1: 1.247823123317706, 2: 0.30382355147152684}, середня внутрішньокласова відстань по методу 3.8658056814598907

Найкращим виявився метод KMeans, як для двох, так і для трьох кластерів

## Дендрограма для перших 100 записів

In [14]:

```
1 fig = ff.create_dendrogram(X.head(100))  
2 fig.update_layout(width=1000, height=600)  
3 fig.show()
```



### 3. Проаналізувати отримані результати, враховуючи останній стовпчик початкових даних

In [15]:

```
1 res_dict = {  
2     'km2_res': km2_res,  
3     'km3_res': km3_res,  
4     'agg2_res': agg2_res,  
5     'agg3_res': agg3_res,  
6     'gaus2_res': gaus2_res,  
7     'gaus3_res': gaus3_res,  
8     'y': y  
9 }  
10 res_df = pd.DataFrame(res_dict)
```

In [16]:

```
1 res_df.groupby('y')['km2_res', 'agg2_res', 'gaus2_res'].agg([np.mean, np.std])
```

Out[16]:

	km2_res		agg2_res		gaus2_res	
	mean	std	mean	std	mean	std
y						
0.000	0.0	0.0	1	0.0	1	0.0
0.525	0.0	NaN	1	NaN	1	NaN
0.556	0.0	NaN	1	NaN	1	NaN
0.825	1.0	NaN	0	NaN	0	NaN
0.842	1.0	NaN	1	NaN	1	NaN
0.860	1.0	NaN	1	NaN	1	NaN
0.884	0.0	NaN	1	NaN	1	NaN
1.009	0.0	NaN	1	NaN	1	NaN

In [17]:

```
1 res_df.groupby('y')['km3_res', 'agg3_res', 'gaus3_res'].agg(  
2     lambda x: x.value_counts().index[0]  
3 )
```

Out[17]:

	km3_res	agg3_res	gaus3_res
y			
0.000	0	0	2
0.525	0	0	2
0.556	0	0	2
0.825	1	1	0
0.842	1	0	1
0.860	1	0	1
0.884	0	0	2
1.009	0	0	2
1.112	1	1	0