



Министерство науки и высшего образования Российской Федерации  
Федеральное государственное бюджетное образовательное учреждение  
высшего образования  
«Московский государственный технический университет имени  
Н.Э. Баумана  
(национальный исследовательский университет)»  
(МГТУ им. Н.Э. Баумана)

---

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

## Отчёт по лабораторной работе №9 по дисциплине "Методы машинного обучения"

Тема Кластерный анализ

Студент Варламова Е. А.

Группа ИУ7-23М

Оценка (баллы) \_\_\_\_\_

Преподаватели Солодовников Владимир Игоревич

# СОДЕРЖАНИЕ

<b>1</b>	<b>Теоретическая часть</b>	<b>3</b>
1.1	Постановка задачи . . . . .	3
1.2	K-means . . . . .	4
1.3	DBSCAN . . . . .	4
<b>2</b>	<b>Практическая часть</b>	<b>6</b>
2.1	Выбор средств разработки . . . . .	6
2.2	ПО . . . . .	6
2.3	Описательный статистический анализ данных . . . . .	11
2.4	Кластерный анализ данных . . . . .	13
2.4.1	DBSCAN: окрестность соседства . . . . .	14
2.4.2	DBSCAN: минимальное количество в кластере . . . . .	14
2.4.3	DBSCAN: алгоритм расчёта расстояния . . . . .	14
2.4.4	kmeans: количество кластеров . . . . .	15
2.4.5	kmeans: алгоритм выбора центроида . . . . .	15
2.5	Оценка работы алгоритмов . . . . .	16
	<b>СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ</b>	<b>19</b>

# 1 | Теоретическая часть

В 1936 году Рональд Фишер опубликовал статью, в которой он представил набор данных об ирисах, состоящий из трех различных видов: *Iris setosa*, *Iris versicolor* и *Iris virginica*. Для каждого вида ириса были измерены четыре характеристики: длина чашелистика, ширина чашелистика, длина лепестка и ширина лепестка.

Фишер использовал этот набор данных для демонстрации метода дискриминантного анализа, который позволяет разделить объекты на классы на основе их признаков. Он показал, что можно эффективно классифицировать ирисы по видам, используя комбинацию этих четырех признаков.

Ирисы Фишера стали одним из самых популярных наборов данных в машинном обучении и статистике. Они часто используются для тестирования алгоритмов классификации и кластеризации, а также для демонстрации методов визуализации данных.

Целью данной лабораторной работы является применение алгоритмов кластеризации для решения задачи кластеризации ирисов Фишера. Для этого необходимо решить следующие задачи:

- описать методы кластеризации;
- привести особенности реализации ПО, решающего поставленную задачу;
- провести сравнение результатов работы алгоритмов.

## 1.1 Постановка задачи

- Провести описательный статистический анализ данных. Рассчитать оценки мат. ожидания, доверительный интервал, среднеквадратичное отклонение. Построить графики и гистограммы.
- Построить корреляционную матрицу. Проверить значимость корреляции.

- Провести кластерный анализ данных. Использовать не менее двух алгоритмов кластеризации (например: иерархический, К-средних, DBSCAN). Варьировать различные значения гиперпараметров и тип расстояний.
- Оценить работу алгоритмов с использованием внешних и внутренних мер оценки качества. Определить оптимальное количество кластеров и их структуру.

## 1.2 K-means

Цель алгоритма: разделить данные на  $k$  кластеров при условии наличия набора данных  $X$  с  $n$  объектами и  $p$  признаками.

Шаги алгоритма:

1. Инициализация: выбрать  $k$  случайных точек из  $X$  в качестве начальных центроидов  $c_1, c_2, \dots, c_k$ .
2. Назначение: для каждого объекта  $x_i$  в  $X$  присвоить его ближайшему центроиду  $c_j$ , то есть (для метрики Евклидово расстояние):

$$j = \operatorname{argmin}_j \|x_i - c_j\|^2 \quad (1.1)$$

3. Обновление: для каждого кластера вычислить новый центроид как среднее всех точек в этом кластере.
4. Повторение: повторять шаги 2 и 3, пока центроиды не перестанут существенно меняться или не будет достигнуто максимальное число итераций.
5. Результат: присвоить каждую точку  $x_i$  кластеру, которому принадлежит ее ближайший центроид.

## 1.3 DBSCAN

Цель алгоритма: разделить данные на кластеры, разделяя плотные области от разреженных при наличии набора данных  $X$  с  $n$  объектами и  $p$  признаками, радиуса окрестности  $\epsilon$  и минимальном количестве точек в кластере  $\minPts$ .

Шаги алгоритма:

1. Инициализация: пометить все точки как не посещенные.
2. Посещение точек: для каждой непосещенной точки  $x_i$  найти все точки в окрестности  $\epsilon$  от  $x_i$ . Если количество точек в окрестности больше или равно  $\minPts$ , создать новый кластер. Иначе пометить  $x_i$  как шум.

3. Расширение кластеров: для каждой точки в созданном кластере посетить все ее непосещенные соседи в окрестности  $\epsilon$ . Если количество соседей больше или равно  $minPts$ , добавить их в кластер.
4. Повторять шаги 2 и 3, пока все точки не будут посещены.

## 2 | Практическая часть

### 2.1 Выбор средств разработки

В качестве языка программирования был использован язык Python, поскольку этот язык кроссплатформенный и для него разработано огромное количество библиотек и модулей, решающих разнообразные задачи.

В частности, имеются библиотеки, включающие в себя алгоритмы кластеризации, в частности sklearn [1].

Для создания графиков была выбрана библиотека matplotlib [2], доступная на языке Python, так как она предоставляет удобный интерфейс для работы с данными и их визуализации.

### 2.2 ПО

В листинге 2.1 представлен код, решающий задачу кластеризации двумя алгоритмами.

Листинг 2.1: Код кластеризации

```
1 import pandas as pd
2 import numpy as np
3 import scipy.stats as stats
4 import matplotlib.pyplot as plt
5 from sklearn.cluster import KMeans, DBSCAN
6 from sklearn.decomposition import PCA
7 from sklearn.model_selection import train_test_split
8 from sklearn.preprocessing import LabelEncoder
9 import umap
10 import seaborn as sns
11 import imageio
12 from sklearn.model_selection import train_test_split
13 from sklearn.metrics import accuracy_score, precision_score, recall_score,
    f1_score, confusion_matrix
14 from sklearn.metrics import mean_absolute_error, mean_squared_error,
    r2_score, adjusted_rand_score
15 n_cluster = 4
16 cols = [ "sepal_length", "sepal_width", "petal_length", "petal_width"]
```

```

17 def draw_description(data):
18     confidence_level = 0.95
19
20     for col in data.columns:
21         column_data = data[col]
22
23         mean = np.mean(column_data)
24
25         confidence_interval = stats.t.interval(0.95, len(column_data)-1, loc
26             =mean, scale=stats.sem(column_data))
27
28         std_dev = np.std(column_data)
29
30         plt.hist(column_data, bins='auto', alpha=0.7, color='skyblue',
31             edgecolor='black', linewidth=1.2)
32         plt.axvline(mean, color='r', linestyle='dashed', linewidth=2, label=
33             f'                : {mean:.2f}')
34         plt.axvline(mean - std_dev, color='g', linestyle='dashed', linewidth
35             =2, label=f'                - std: {mean - std_dev:.2f
36             }')
37         plt.axvline(mean + std_dev, color='g', linestyle='dashed', linewidth
38             =2, label=f'                + std: {mean + std_dev:.2f
39             }')
40         plt.axvspan(confidence_interval[0], confidence_interval[1], color='
41             gray', alpha=0.3, label='
42             (95%)')
43         plt.xlabel('                ')
44         plt.ylabel('                ')
45         plt.title('                ')
46         plt.legend()
47         plt.savefig(f"{col}.png")
48         plt.clf()
49
50 def draw_corr(data):
51     corr=data.corr()
52     plt.figure(figsize=(30,20))
53     sns.heatmap(corr, annot=True, cmap="Reds")
54     plt.title('                ', fontsize=20)
55     plt.savefig("correlation.png")
56     plt.clf()
57
58 def read_data(fn):
59     return pd.read_csv(fn)
60
61 def preprocess_data(df):
62     return df[cols].dropna(how='any')
63
64 def count_elbow(df):
65     inertia = []
66     n = 15

```

```

58     for k in range(1, n):
59         kmeans = KMeans(n_clusters=k)
60         kmeans.fit(df)
61         inertia.append(kmeans.inertia_)
62
63     plt.figure(figsize=(8, 6))
64     plt.plot(range(1, n), inertia, marker='o')
65     plt.xlabel('Number of clusters')
66     plt.ylabel('Inertia')
67     plt.title('Elbow Method for Optimal K')
68     plt.savefig("elbow.png")
69     plt.clf()
70
71 def count_k_means(df):
72     kmeans = KMeans(n_cluster)
73     kmeans.fit(df)
74     return kmeans.labels_
75 def count_dbscan(df):
76     dbscan = DBSCAN(eps=4.2, min_samples=5)
77     return dbscan.fit_predict(df)
78 def count_pca(df):
79     pca = PCA(n_components=2)
80     return pca.fit_transform(df)
81 def count_umap(df):
82     return umap.UMAP(n_components=2).fit_transform(df)
83
84 def plot(title, labels, df_resize, dir_name = "."):
85     plt.figure(figsize=(12, 6))
86     sns.scatterplot(x=df_resize[:, 0], y=df_resize[:, 1], hue=labels,
87                    palette='viridis')
88     plt.title(title)
89     plt.legend()
90     plt.savefig(dir_name + title + ".png")
91     plt.clf()
92
93 def visual(df, f_cluster, f_resize, title, dir_name = "."):
94     df_resize = f_resize(df)
95     labels = f_cluster(df)
96
97     plot(title, labels, df_resize, dir_name = ".")
98
99 def vary_metric_dbscan(df, answers):
100     metric = ["cosine", "euclidean", "l1", "l2", "manhattan"]
101     scores = []
102     for m in metric:
103         dbscan = DBSCAN(eps=4.2, min_samples=5, metric = m)
104         labels = dbscan.fit_predict(df)
105         df_resize = count_umap(df)
106         plot(f"metric_{m}", labels, df_resize, dir_name = "
            research_dbscan_metric/")

```



```

106         scores.append(adjusted_rand_score(answers, labels))
107
108     plt.figure(figsize=(12, 6))
109     plt.plot(metric, scores)
110     plt.title("scores")
111     plt.legend()
112     plt.savefig("research_dbscan_metric/scores.png")
113     plt.clf()
114 def vary_eps_dbscan(df, answers):
115     epss = np.linspace(2, 6, 21)
116     scores = []
117     for eps in epss:
118         dbscan = DBSCAN(eps=eps, min_samples=5)
119         labels = dbscan.fit_predict(df)
120         df_resize = count_umap(df)
121         eps = round(eps * 10)
122         plot(f"metric_{eps}", labels, df_resize, dir_name = "
            research_dbscan_eps/")
123         scores.append(adjusted_rand_score(answers, labels))
124
125     plt.figure(figsize=(12, 6))
126     plt.plot(epss, scores)
127     plt.title("scores")
128     plt.legend()
129     plt.savefig("research_dbscan_eps/scores.png")
130 def vary_min_samples_dbscan(df, answers):
131     scores = []
132     x = []
133     for s in range(5, 50, 5): # 150
134         dbscan = DBSCAN(eps=4.2, min_samples=s)
135         labels = dbscan.fit_predict(df)
136         df_resize = count_umap(df)
137         plot(f"samples_{s}", labels, df_resize, dir_name = "
            research_dbscan_samples/")
138         scores.append(adjusted_rand_score(answers, labels))
139         x.append(s)
140
141     plt.figure(figsize=(12, 6))
142     plt.plot(x, scores)
143     plt.title("scores")
144     plt.legend()
145     plt.savefig("research_dbscan_samples/scores.png")
146 def vary_num_clusters_kmeans(df, answers):
147     scores = []
148     x = []
149     for cluster in range(2, 11, 1):
150         kmeans = KMeans(cluster)
151         kmeans.fit(df)
152         labels = kmeans.labels_
153         df_resize = count_umap(df)

```

```

154         plot(f"clusters_{cluster}", labels, df_resize, dir_name = "
            research_kmeans_clusters/")
155     scores.append(adjusted_rand_score(answers, labels))
156     x.append(cluster)
157
158     plt.figure(figsize=(12, 6))
159     plt.plot(x, scores)
160     plt.title("scores")
161     plt.legend()
162     plt.savefig("research_kmeans_clusters/scores.png")
163 def vary_centroid_algo_kmeans(df, answers):
164     scores = []
165     metrics = ['random', 'k-means++']
166     for metric in metrics:
167         kmeans = KMeans(n_cluster, init=metric)
168         kmeans.fit(df)
169         labels = kmeans.labels_
170         df_resize = count_umap(df)
171         plot(f"algo_{metric}", labels, df_resize, dir_name = "
            research_kmeans_algo/")
172         scores.append(adjusted_rand_score(answers, labels))
173
174     plt.figure(figsize=(12, 6))
175     plt.plot(metrics, scores)
176     plt.title("scores")
177     plt.legend()
178     plt.savefig("research_kmeans_algo/scores.png")
179 def count_accuracy(df, answers):
180     kmeans = KMeans(4, init = 'k-means++')
181     kmeans.fit(df)
182     labels_kmeans = kmeans.labels_
183     print("kmeans", adjusted_rand_score(answers, labels_kmeans))
184
185     dbscan = DBSCAN(eps=4.2, min_samples=5)
186     labels_dbscan = dbscan.fit_predict(df)
187     print("dbscan", adjusted_rand_score(answers, labels_dbscan))

```

## 2.3 Описательный статистический анализ данных

Был проведён описательный статистический анализ данных: рассчитаны оценки мат. ожидания, доверительные интервалы, среднеквадратичное отклонение. Построены графики и гистограммы:

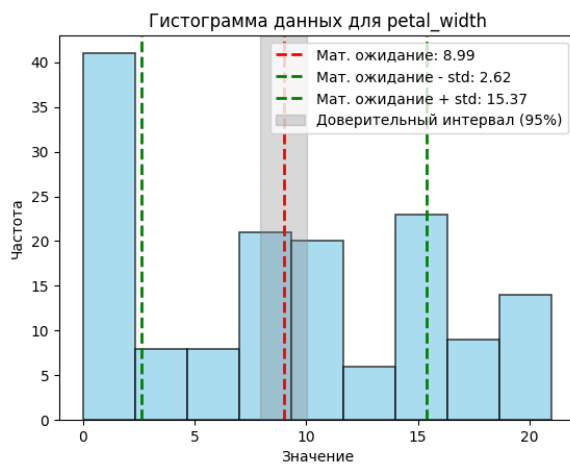


Рис. 2.1: Анализ petal\_width

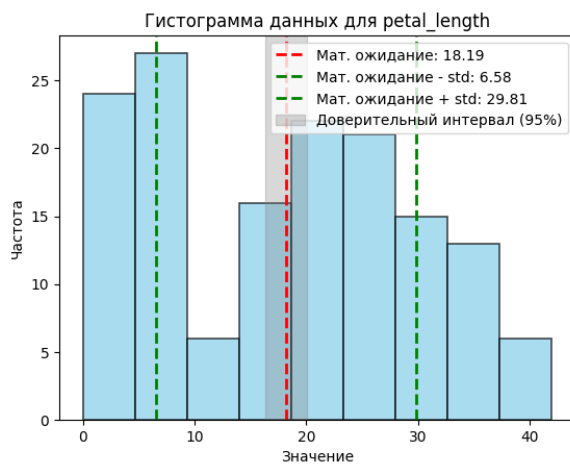


Рис. 2.2: Анализ petal\_length

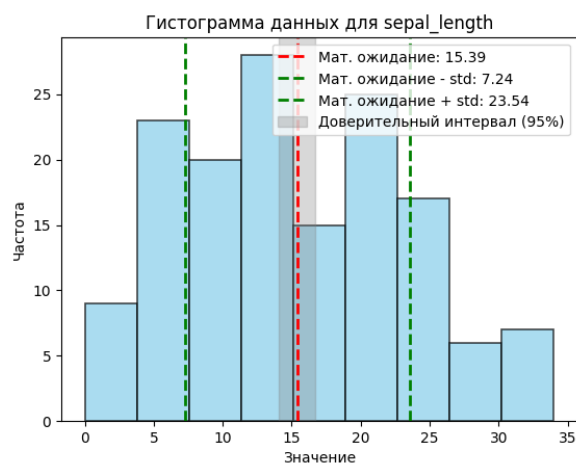


Рис. 2.3: Анализ sepal\_length

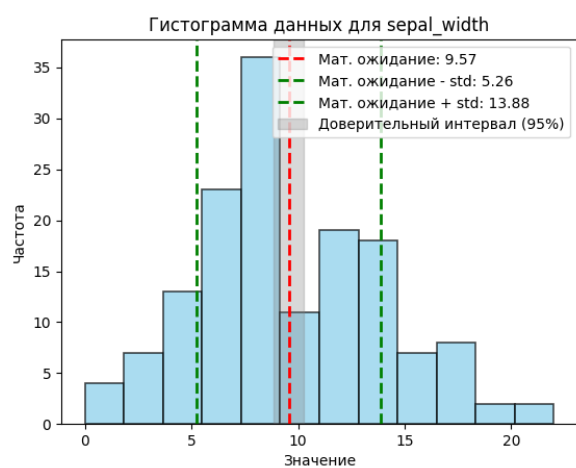


Рис. 2.4: Анализ sepal\_width

Была построена матрица корреляций:

Видно, что в большинстве классов наблюдается слабая корреляция между признаками.

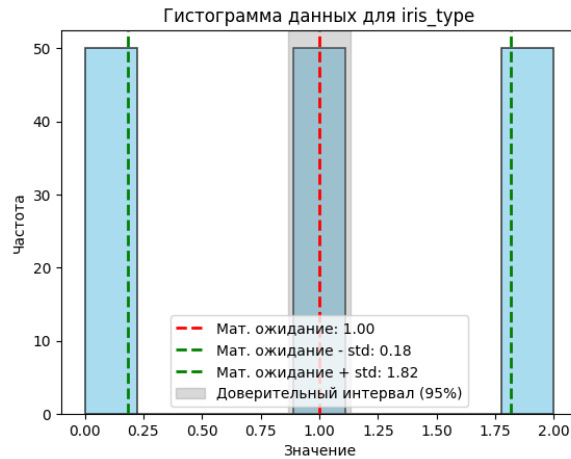


Рис. 2.5: Анализ iris\_type

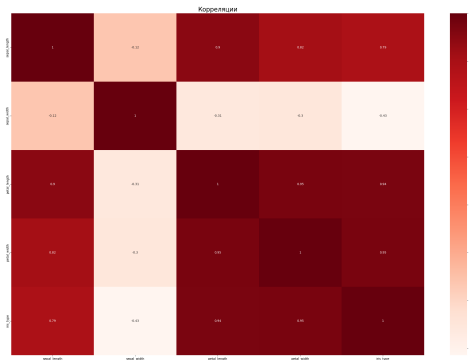


Рис. 2.6: Корреляционная матрица

## 2.4 Кластерный анализ данных

Для кластеризации были использованы 2 алгоритма и для каждого из них были проварьированы основные параметры:

- К-средних (параметры – количество кластеров, алгоритм выбора центроида);
- DBSCAN (параметры – окрестность соседства, минимальное количество в кластере, алгоритм расчёта расстояния).

В качестве метрики качества кластеризации использовалась разница разбиения на кластеры с помощью алгоритма и исходного деления на кластеры в датасете.

### 2.4.1 DBSCAN: окрестность соседства

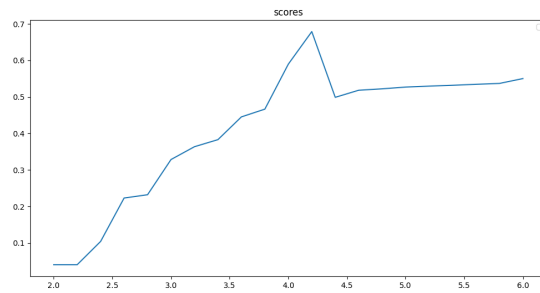


Рис. 2.7: DBSCAN: окрестность соседства

Видно, что лучшее значение 4.2.

### 2.4.2 DBSCAN: минимальное количество в кластере

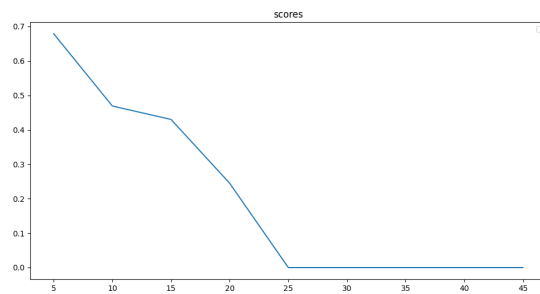


Рис. 2.8: DBSCAN: минимальное количество в кластере

Видно, что лучшее значение 5.

### 2.4.3 DBSCAN: алгоритм расчёта расстояния

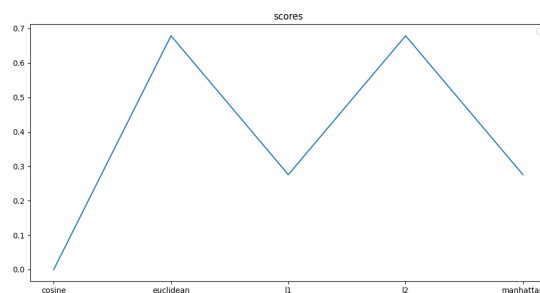


Рис. 2.9: DBSCAN: алгоритм расчёта расстояния

Видно, что лучший алгоритм — евклидово расстояние.

#### 2.4.4 kmeans: количество кластеров

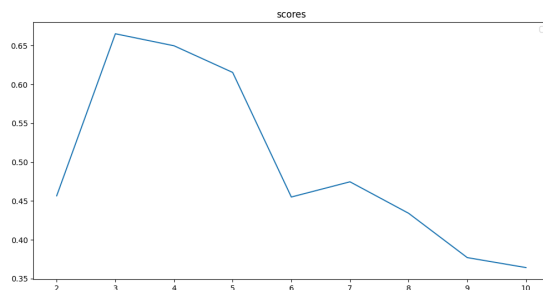


Рис. 2.10: kmeans: количество кластеров

Видно, что лучшее значение – 3. Для определения количества кластеров существует метод локтя. Проверим полученный результат с помощью него.

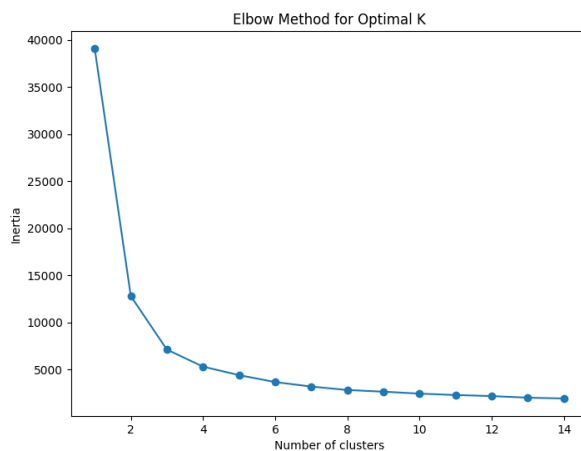


Рис. 2.11: Метод локтя

Видно, что для метода локтя оптимальным является 4-5 кластеров.

#### 2.4.5 kmeans: алгоритм выбора центроида

Видно, что алгоритм выбора центроида не влияет на конечный результат кластеризации.

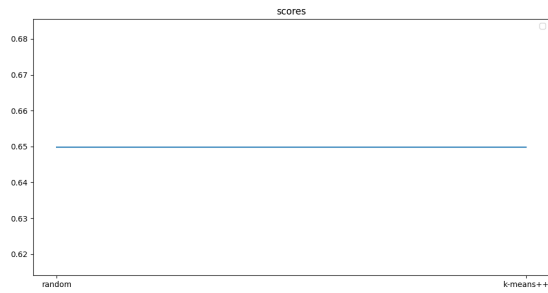


Рис. 2.12: kmeans: алгоритм выбора центроида

## 2.5 Оценка работы алгоритмов

С помощью определённых оптимальных параметров алгоритмов данные были кластеризованы и визуализированы с помощью разных алгоритмов уменьшения размерности данных, а также вычислены точности кластеризации:

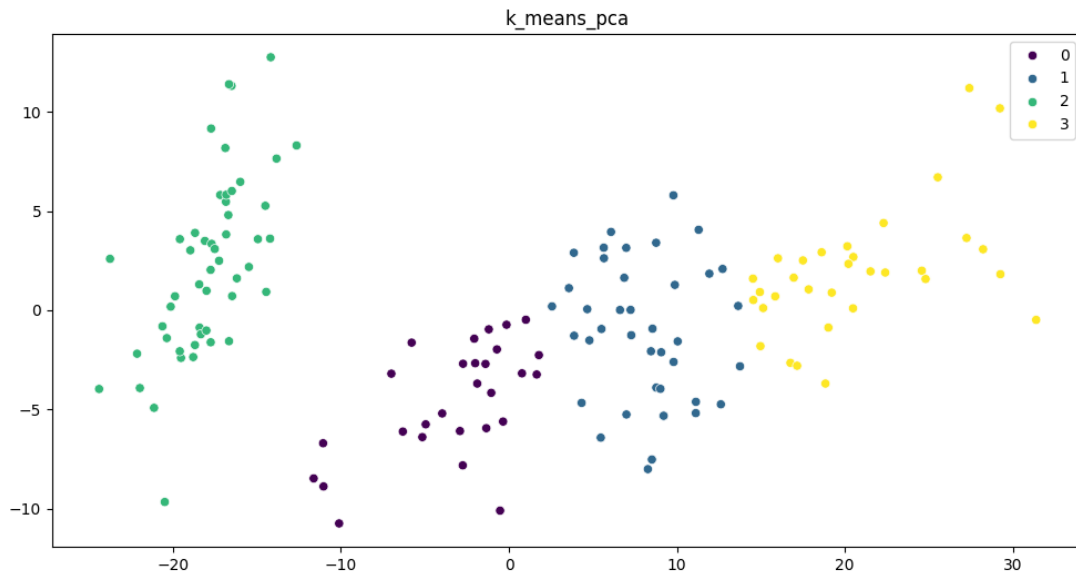


Рис. 2.13: kmeans, pca



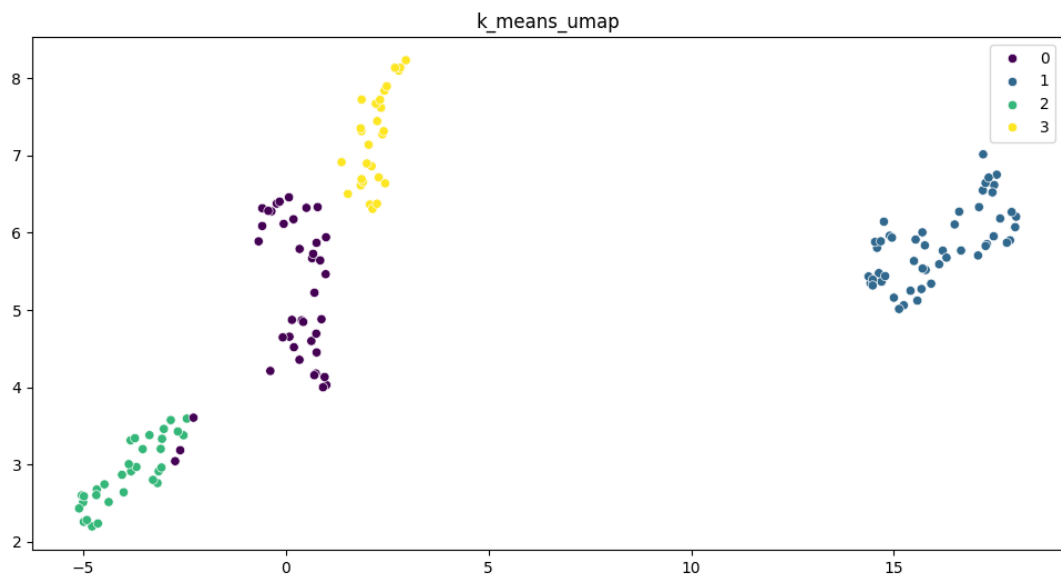


Рис. 2.14: kmeans, UMAP

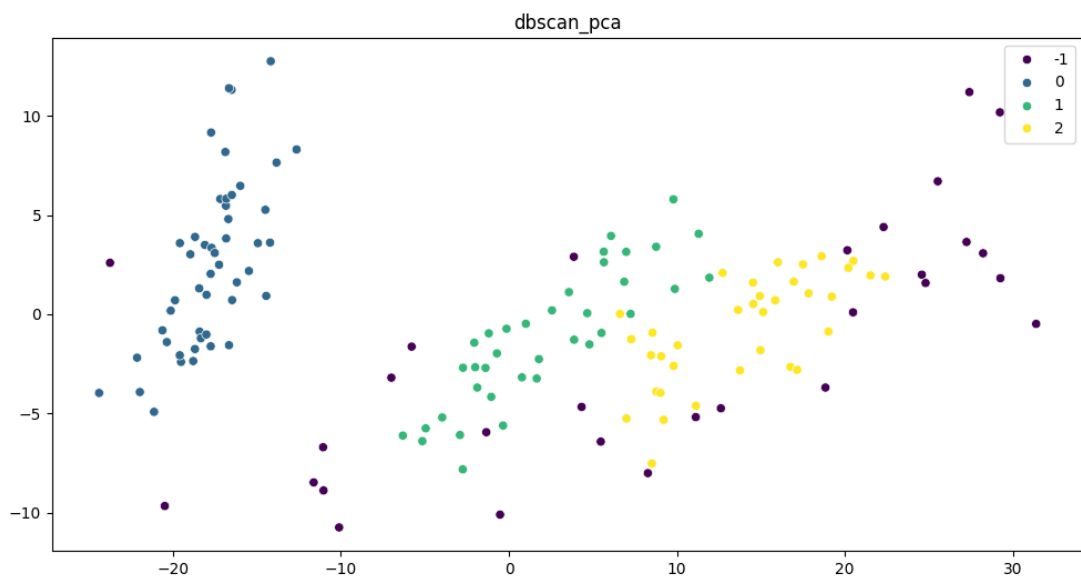


Рис. 2.15: DBSCAN, pca

Точности кластеризации(диапазон:  $[-0.5 - 1]$ , где 0 – случайная кластеризация, 1 – точная):

— kmeans: 0.65.

— DBSCAN: 0.63.

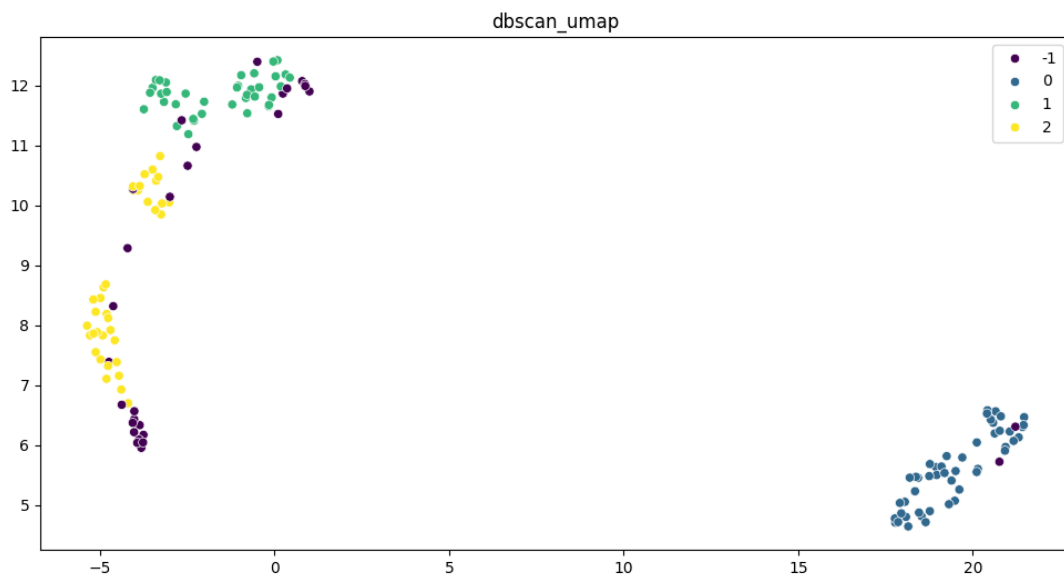


Рис. 2.16: DBSCAN, UMAP

# СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Scikit-learn: Machine learning in Python / F. Pedregosa [и др.]. — 2011.
2. Библиотека визуализации данных matplotlib [Электронный ресурс]. — Режим доступа: URL: <https://matplotlib.org> (дата обращения: 13.12.2023).