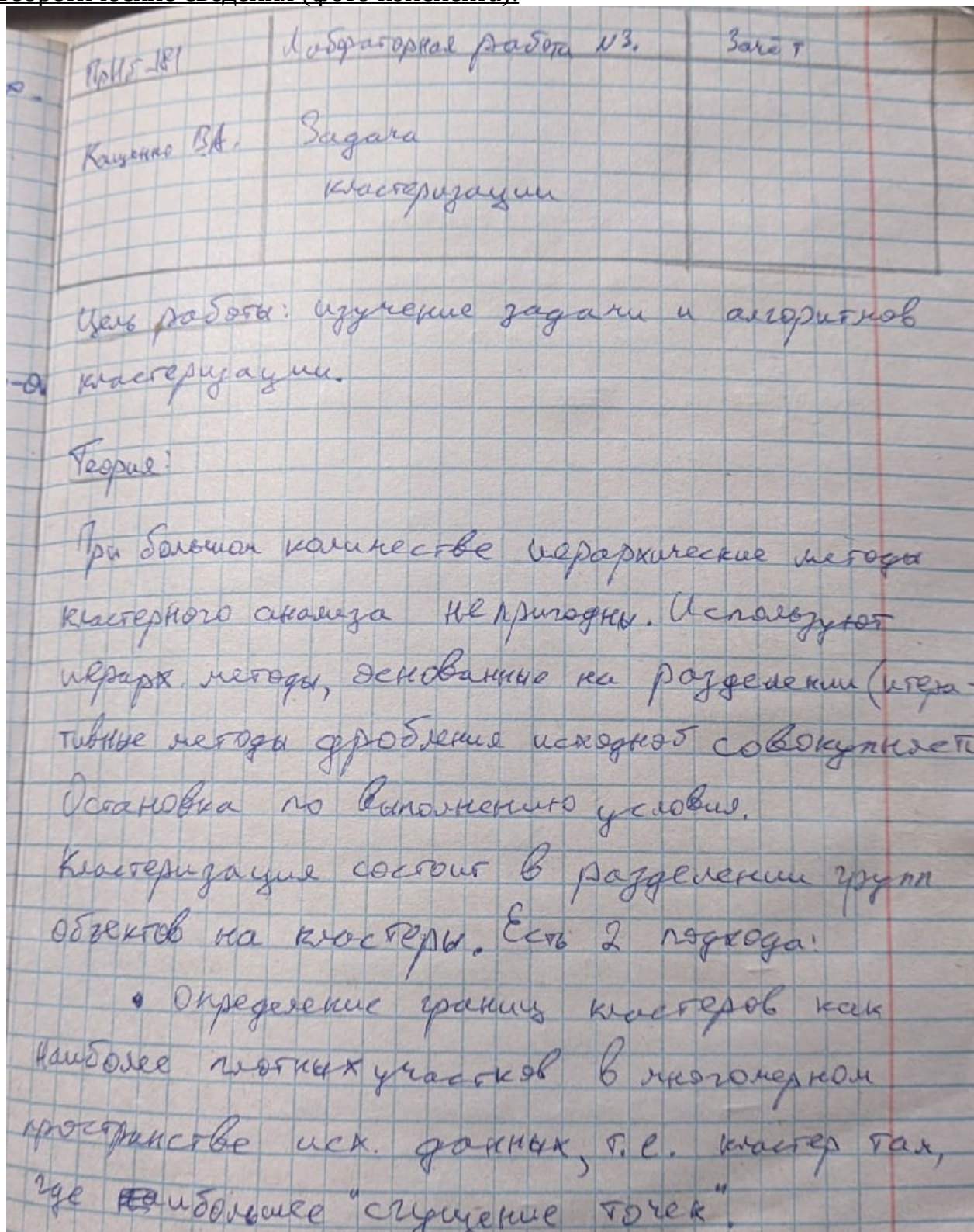


Приб-181	Лабораторная работа №3	Зачёт
Кащенко В. А.	Задача кластеризации.	

Цель работы: Изучить задачу и алгоритмы кластеризации.

Теоретические сведения (фото конспекта):





- минимизация меры различия объектов.

## 1. Иерархич. алгоритмы.

Получены алгоритмы, основанные на поиске оптимального в опр. смысле разбиения ~~точек~~<sup>данных</sup> множества данных на ~~группы (точки)~~ кластеры. Из множества функций используются алгоритмы разбиения. Они пытаются кластеризовать данные т. образом, чтобы целевая ф-я достигла экстремума (минимума).

Алгоритмы кластеризации, основанные на методе разбиения (будут рассмотрены).

Базовые понятия;

— обучающее множество  $M$  (входное), кот. разбивается;

— метрика расстояния:

$$d_A^2(m_i, c^{(i)}) = (m_i - c^{(i)})^2 = (m_i - c^{(i)})^t A (m_i - c^{(i)}),$$



где матрица  $A$  - способ вычисления расстояний

Пример для единичной матрицы используется расстояние по Евклиду:

- вектор центров кластеров  $C$ ;
- матрица разбиения по кластерам  $U$ ;
- целевая ф-я:  $J = J(M, d, C, U)$ ;
- набор ограничений.

Метрика расстояний - довольно важное понятие кластеризации. С помощью расстояний между входными векторами определяются сходства и различия. Есть много способов вычислить расстояние: Евклидова, Хэмминга, Манхэттенского и др. Выбор метода зависит от природы исследуемых объектов и влияет на результат.

Целевая ф-я, кот. даёт решение задачи кластеризации. Алгоритм кластеризации - последовательность действий поиска минимума.



Матрица разбиения —  $R$ -ая иерархия,  
классификации

Это таблица, где  $i$ -строка содержит  
основание матрицы и получается итогов.  
разбиение. Центр кластера — вектор,  
степень принадлежности которого заданному  
классу максимальна. Как правило, центров  
нет в исходном  $m$ -ве.

Набор ограничений связан с условиями,  
налагаемые на значения элементов матрицы  
принадлежности. Они определяются  
алгоритмом кластеризации.

## 2. Алгоритм К-средних (K-means)

Он распространён и также носит имя быстрого  
кластерного анализа. Описание:  
работы Хартигана и Вонга (Hartigan and Wong, 1979).  
В отл. от иерарх. методов, кот. не требуют предв.  
предположений отл. числа кластеров, где возм-ны



использование кебх-но иметь гипотезу о наиб.  
вероятной кол-ве кластеров.

Алгоритм К-средних строит К-кластеров,  
расположенных на возможно больших расстояниях  
друг от друга. Основное тут задат — количество  
предположений (гипотез) о тн. числе кластеров, при  
кот. они должны быть различны настолько, насколько  
это возможно.

Выбор числа К может базироваться на р-ах  
предметов исследования, теоретических соображениях и  
интуиции. Общ. идея — <sup>число</sup> К кластеров составле-  
<sup>кластеров</sup> <sup>еще</sup> так, что средние в кластере макс-но возможно  
отличаются друг от друга.

К-means хорошо работает, если данные можно  
разделить на компактные, примерно сферич.  
группы.

Данный алгоритм — прообраз всех алгоритмов  
четкой кластеризации, и его организация позволяет  
лучшему пониманию принципов, заложенных в более  
сложные алгоритмы.



Базовые понятия:

— обучающее множество  $M = \{m_i\}_{i=1}^d$  ( $d$  — кол-во точек/векторов данных);

— метрика расстояний:

$$d_A^2(m_i, c^{(i)}) = \|m_i - c^{(i)}\|^2 = (m_i - c^{(i)})^T (m_i - c^{(i)});$$

— вектор центров кластеров  $C = \{c^{(i)}\}_{i=1}^c$ , где

$$c^{(i)} = \frac{\sum_{j=1}^d u_{ij} m_j}{\sum_{j=1}^d u_{ij}}, \quad 1 \leq i \leq c$$

— матрица разбиения  $U = \{u_{ij}\}$ , где

$$u_{ij}^{(1)} = \begin{cases} 1, & \text{при } d(m_j, c_k^{(1)}) = \min_{1 \leq k \leq c} d(m_j, c_k^{(1)}) \\ 0, & \text{в ост. случаях.} \end{cases}$$

— целевая ф-я:

$$J(M, U, C) = \sum_{i=1}^c \sum_{j=1}^d u_{ij} d_A^2(m_j, c^{(i)})$$

— набор ограничений:

$$u_{ij} \in \{0, 1\}$$

$$\sum_{i=1}^c u_{ij} = 1; \quad 0 < \sum_{i=1}^c u_{ij} < d, \text{ коэф. опред-ся}$$



что  $\forall$  вектор данных может принадлежать  
одному кластеру и не принадлежать остальным  
В  $k$  кластере не менее 1 точки, но не более  
общ. кол-ва точек.

Конструктивно алгоритм — итерационная процедура.

1. Инициализация начального разбиения (random,  
for example), выбор точности  $\delta$  (используется  
в условии завершения алгоритма), инициализация  
номеров итерации  $l = 0$ .

2. Центры кластеров:

$$c_i^{(l)} = \frac{\sum_{j=1}^n u_{ij}^{(l-1)} \cdot m_j}{\sum_{j=1}^n u_{ij}^{(l-1)}}, \quad 1 \leq i \leq c$$

3. Обновить матрицу разбиения с целью  
минимизации квадратов ошибок:

$$u_{ij}^{(l)} = \begin{cases} 1, & \text{при } d(m_j, c_i^{(l)}) = \min_{1 \leq k \leq c} d(m_j, c_k^{(l)}) \\ 0, & \text{в ост. случаях} \end{cases}$$

4. Проверить условие  $\|u^{(l)} - u^{(l-1)}\| < \delta$

Осн. недостаток — большой размер пространства  
разбиения



Один из способов устранения недостатка - представление элементов матрицы разбиения числами единичного интервала, т.е. принадлежность должна определяться  $\phi$ -ей принадлежности

Данный подход имеет воплощение в алгоритме нечеткой кластеризации Fuzzy C-Means

3. Кластеризация с библиотекой Python `sklearn.cluster`

Алгоритм бывает 2-х вариантов - класс с методом `fit` для обучения на обучающем множестве и `phi`-я, кот. укажет данные обр. множества возвращает массив целочисл. меток, соотв кластерам. Метод `phi` в классе можно найти в атрибуте `labels`



Входные данные:

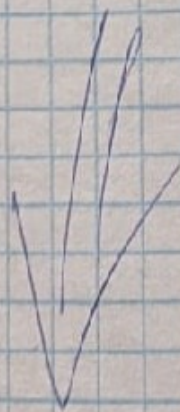
- алгоритмы в качестве вх. данных могут принимать разн. виды матриц.

- Н. метод принимает стандарт. матрицы данных размера  $[n\_samples, n\_features]$

их можно получить с помощью `sklearn.feature\_extractor`.

- Для Affinity Propagation, Spectral Clustering и DBSCAN можно ввести матрицу подобия  $[n\_samples, n\_samples]$

их можно получить с помощью `sklearn.metrics.pairwise`.



методы пакета  
`sklearn`



Название метода	Параметры	Настройка параметров	Варианты исполнения	Геометрия (система)
K-means	Кол-во кластеров	Большое $n$ -samples, среднее $n$ -clusters < 1000 Mini Batch	Универсальный, не слишком много кластеров мощ.	Расстояние между точками
Affinity Propagation	Формирование предположений об образцах	Не масштабируется	Много кластеров, не обязательно разное количество	Расстояние между центрами (поиск ближайшего)
Word hierarchical clustering	Кол-во кластеров или вычисл. расстояние	Большое $n$ -samples и $n$ -clusters	Много кластеров, от 1 до 1000	Расстояние между точками
Agglomerative clustering	Кол-во кластеров, вычисл. расстояние	Большое $n$ -samples и $n$ -clusters	Много кластеров, от 1 до 1000	Расстояние между точками
DBSCAN	Размер ядра	Большое $n$ -samples, сред. $n$ -clusters	Неплохо, но не масштабируется	Расстояние между точками
OPTICS	Мин. кол-во в кластере	М. большое $n$ -samples, большое $n$ -clusters	Неплохо, но не масштабируется	Расстояние между точками
Gaussian mixtures	Многочисленность	Настройка параметров	Плохо, но можно сделать лучше	Расстояние между центрами
Birch	Параметры	Большое $n$ -samples и $n$ -clusters	Большое количество выбросов, отбраковка лишних	Расстояние между точками



Кластеризация с помощью метода кластеризации полезна, когда кластеры имеют иную форму, кроме сферической, и стандартное евклидово расстояние не является правильной метрикой

Ход работы:

1. Подготовлен csv-файл - Banknote-authentication.csv

2. Выполнен алгоритм K-means (Scikit-learn)

3. Для сравнения выполнен алгоритм Mean-Shift (Сдвиг среднего значения)

4. Алгоритм среднего сдвига кластеризации не делает никаких предположений, в отличие от K-means,

мы видим результаты работы параметрического (K-means) и непараметрического алгоритма (Mean-Shift)

Вывод: Успешна задача и алгоритмы кластеризации.



### Постановка задачи:

Выполнить кластеризацию методом K-means на примере датасета о банкнотах (banknote-authentication.csv).

Это датасет данных об изображениях подлинных и поддельных образцов банкнот. Для оцифровки использовалась промышленная камера, которая обычно используется для проверки печати. Окончательные изображения имеют размер 400x 400 пикселей. За счет линзы объекта и расстояния до исследуемого объекта были получены полутоновые изображения с разрешением около 660 dpi. Применялось пакетное вейвлет-преобразование (WPT) для проверки подлинности защищенных документов, в частности банкнот.

### Практическая часть (код программы):

```
import matplotlib as mpl
import matplotlib.pyplot as plt
from matplotlib import rcParams
import pandas as pd
from sklearn.cluster import KMeans, MeanShift, estimate_bandwidth

# Reading csv file
df = pd.read_csv('banknote-authentication.csv')

# Deleting class column
df = df.drop('class', axis=1)
processed_df = df.to_numpy()[:, :]

# K-means clustering
cluster_num = 2
k_means = KMeans(init='k-means++', n_clusters=cluster_num, n_init=12)
k_means.fit(processed_df)
k_means_labels = k_means.labels_

# Mean-shift clustering
bandwidth = estimate_bandwidth(processed_df, quantile=0.183, n_samples=1000)
m_shift = MeanShift(bandwidth=bandwidth)
m_shift.fit(processed_df)
m_shift_labels = m_shift.labels_
df['ms_cluster'] = m_shift_labels

# Setting up mpl params
mpl.rcParams['font.family'] = 'Times New Roman'
mpl.rcParams['axes.unicode_minus'] = False
rcParams.update({'figure.autolayout': True})

fig = plt.figure(figsize=(12, 8))

colors = list(map(lambda x: '#002bff' if x == 1 else '#f7ab05', k_means_labels))

ax1 = fig.add_subplot(1, 2, 1, projection = '3d')
```



```

#ax1.view_init(azim=0, elev=90)
ax1.scatter(processed_df[:, 0], processed_df[:, 1], processed_df[:, 2], c=colors, alpha=0.8)
ax1.set_zlabel('кратность изображения')
ax1.set_xlabel('дисперсия изображения')
ax1.set_ylabel('асимметрия изображения')
ax1.set_title('K-Means')

colors = list(map(lambda x: '#002bff' if x == 1 else '#f7ab05', m_shift_labels))

ax2 = fig.add_subplot(1, 2, 2, projection = '3d')
#ax2.view_init(azim=0, elev=90)
ax2.scatter(processed_df[:, 0], processed_df[:, 1], processed_df[:, 2], c=colors, alpha=0.8)
ax2.set_zlabel('кратность изображения')
ax2.set_xlabel('дисперсия изображения')
ax2.set_ylabel('асимметрия изображения')
ax2.set_title('Mean-shift')

# images showing and saving
plt.savefig('banknote-authentication_clustering_done.svg')
plt.show()

```

### Этапы выполнения работы:

1. Подключаем необходимые библиотеки (в том числе и sklearn):

```

import matplotlib as mpl
import matplotlib.pyplot as plt
from matplotlib import rcParams
import pandas as pd
from sklearn.cluster import KMeans, MeanShift, estimate_bandwidth

```

2. Далее читаем csv-файл с помощью библиотеки pandas и удаляем столбец class. Также происходит приведение типов данных для корректной работы:

```

# Reading csv file
df = pd.read_csv('banknote-authentication.csv')

# Deleting class column
df = df.drop('class', axis=1)
processed_df = df.to_numpy()[:, :]

```

3. Затем используем алгоритм K-means. Укажем при инициализации **k-means++**. Это означает разумное определение значения K и начальных центральных точек кластера K, что влияет на результаты работы метода. **.fit(processed\_df)** выполняет саму кластеризацию, далее происходит работа с ярлыками:

```

# K-means clustering
cluster_num = 2
k_means = KMeans(init='k-means++', n_clusters=cluster_num, n_init=12)
k_means.fit(processed_df)
k_means_labels = k_means.labels_

```

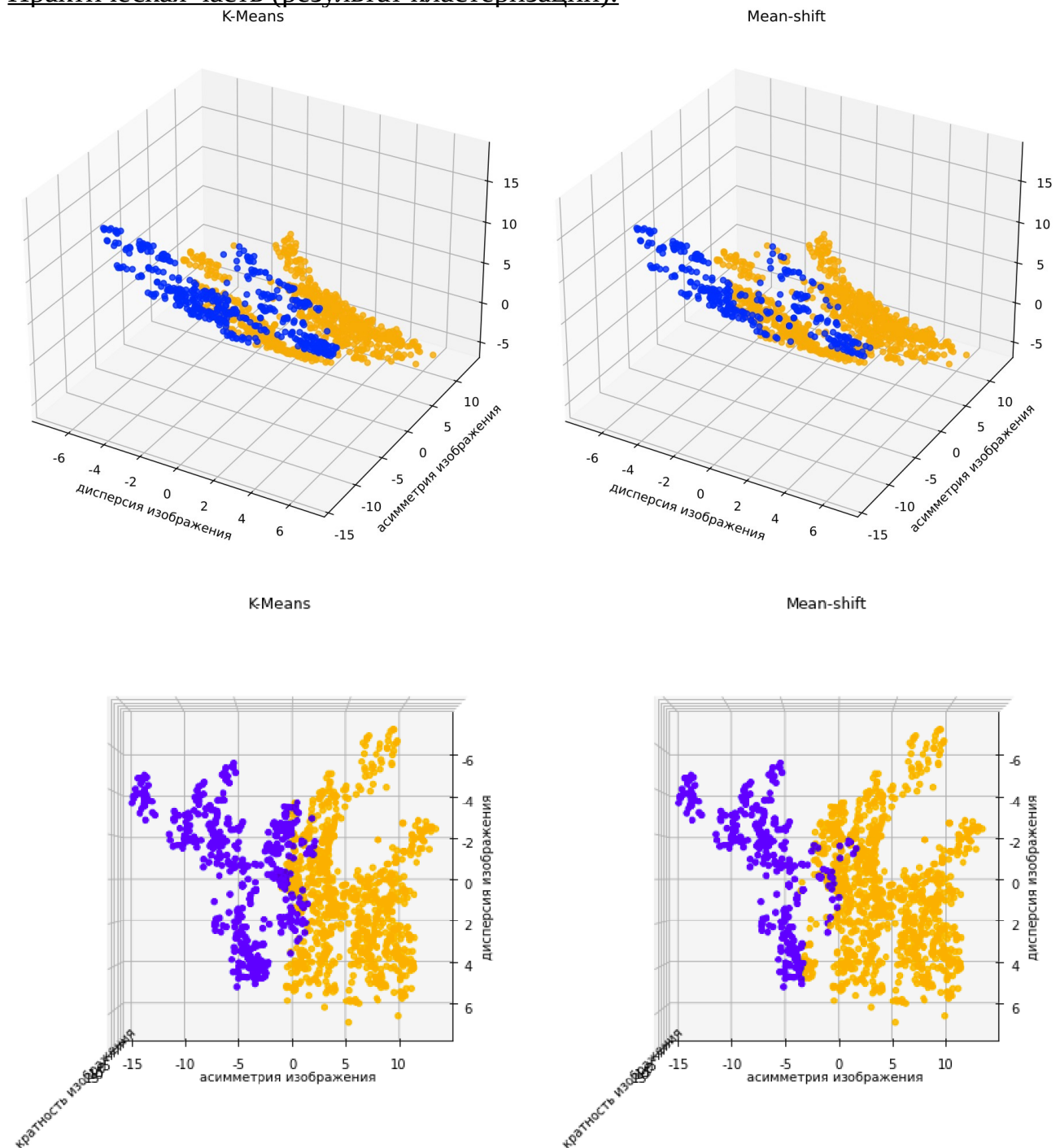


4. Для сравнения используем метод Mean-shift, который характеризуется пропускной способностью текст. В параметрах передаем датасет, медиану попарных расстояний и количество образцов (функция `estimate_bandwidth()` инициализирует поля `bandwidth`), кластеризацию выполняем аналогично k-means:

```
# Mean-shift clustering
bandwidth = estimate_bandwidth(processed_df, quantile=0.183, n_samples=1000)
m_shift = MeanShift(bandwidth=bandwidth)
m_shift.fit(processed_df)
m_shift_labels = m_shift.labels_
df['ms_cluster'] = m_shift_labels_
```

5. Оставшийся код применяется для настройки отображения результата кластеризации.

### Практическая часть (результат кластеризации):





На графиках виден похожий результат выполнения двух алгоритмов, однако метод Mean-shift дает более детальный результат. По полученным сведениям можно отличить поддельные банкноты от настоящих.

Вывод: изучена задача и алгоритмы кластеризации.