

**МИНОБРНАУКИ РОССИИ**  
**САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ**  
**ЭЛЕКТРОТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ**  
**«ЛЭТИ» ИМ. В.И. УЛЬЯНОВА (ЛЕНИНА)**  
**Кафедра математического обеспечения и применения ЭВМ**

**ОТЧЕТ**  
**по лабораторной работе №6**  
**по дисциплине «Статистические методы обработки**  
**экспериментальных данных»**  
**Тема: Кластерный анализ. Метод k-средних.**

Студент гр. 8383

Бабенко Н.С.

Студент гр. 8383

Сахаров В.М.

Преподаватель

Середа А.-В.И.

Санкт-Петербург

2022

## Цель работы

Освоение основных понятий и некоторых методов кластерного анализа, в частности, метода k-means.

## Основные теоретические положения

Задача кластерного анализа заключается в том, чтобы разбить множества исследуемых объектов и признаков на однородные в соответствующем понимании группы. К характеристикам кластера относятся:

1. *Центр кластера* – это среднее геометрическое место точек, принадлежащих кластеру, в пространстве данных.
2. *Радиус кластера* – максимальное расстояние точек, принадлежащих кластеру, от центра кластера.
3. Кластеры могут быть *перекрывающимися*. В этом случае невозможно при помощи используемых процедур однозначно отнести объект к одному из двух или более кластеров. Такие объекты называют спорными.
4. *Спорный объект* – это объект, который по мере сходства может быть отнесен к более, чем одному кластеру.
5. *Размер кластера* может быть определен либо по радиусу кластера, либо по среднеквадратичному отклонению объектов для этого кластера. Объект относится к кластеру, если расстояние от объекта до центра кластера меньше радиуса кластера. Если это условие выполняется для двух и более кластеров, объект является спорным.

Существуют различные способы нормировки данных:

$$z = \frac{(x - \bar{x})}{\sigma}; z = \frac{x}{\bar{x}}; z = \frac{x}{x_{max}}; z = \frac{(x - \bar{x})}{x_{max} - x_{min}}$$

Расстоянием (метрикой) между объектами  $a$  и  $b$  пространстве параметров называется такая величина  $d_{ab}$ , которая удовлетворяет аксиомам:

1.  $d_{ab} > 0$ , если  $a \neq b$ ,
2.  $d_{ab} = 0$ , если  $a = b$ ;
3.  $d_{ab} = d_{ba}$ ;
4.  $d_{ab} + d_{bc} \geq d_{ac}$ .

Мерой близости (сходства) называется величина  $\mu_{ab}$ , имеющая предел и возрастающая с возрастанием близости объектов и удовлетворяющая условиям:

$$\mu_{ab} \text{ непрерывна; } \mu_{ab} = \mu_{ba}; 0 \leq \mu_{ab} \leq 1.$$

Суть метода k-средних заключается в том, что он стремится минимизировать суммарное квадратичное отклонение точек кластеров от центров этих кластеров:

$$V = \sum_{i=1}^k \sum_{x_j \in S_i} (x_j - \mu_i)^2$$

Центроиды выбираются в тех местах, где визуально скопление точек выше. Алгоритм разбивает множество элементов векторного пространства на заранее известное число кластеров  $k$ . Основная идея заключается в том, что на каждой итерации пересчитывается центр масс для каждого кластера, полученного на предыдущем шаге, затем векторы разбиваются на кластеры вновь в соответствии с тем, какой из новых центров оказался ближе по выбранной метрике. Алгоритм завершается, когда на какой-то итерации не происходит изменения центра масс кластеров.

Задание количества кластеров является сложным вопросом. Если нет разумных соображений на этот счет, рекомендуется первоначально создать 2 кластера, затем 3, 4, 5 и так далее, сравнивая полученные результаты.

Возможны две разновидности метода. Первая предполагает пересчет центра кластера после каждого изменения его состава, а вторая – лишь после завершения цикла.

После завершения многомерной классификации необходимо оценить полученные результаты. Для этой цели используются специальные характеристики – функционалы качества. Наилучшим разбиением считается такое, при котором достигается экстремальное (минимальное или максимальное) значение выбранного функционала качества.

В качестве таких функционалов могут быть использованы:

1. Сумма квадратов расстояний до центров кластеров

$$F_1 = \sum_{k=1}^K \sum_{i=1}^{N_k} d^2(X_i^{(k)}, X^{(k)}) \Rightarrow \min$$

2. Сумма внутрикластерных расстояний между объектами

$$F_2 = \sum_{k=1}^K \sum_{X_i, X_j \in S_k} d^2(X_i, X_j) \Rightarrow \min$$

3. Сумма внутрикластерных дисперсий

$$F_3 = \sum_{k=1}^K \sum_{i=1}^{N_k} \sigma_{ij}^2 \Rightarrow \min$$

Здесь  $\sigma$  - дисперсия  $j$ -й переменной в  $k$ -м кластере.

Оптимальным следует считать разбиение, при котором сумма внутрикластерных (внутригрупповых) дисперсий будет минимальной.

### **Постановка задачи**

Дано конечное множество из объектов, представленных двумя признаками (в качестве этого множества принимаем исходную двумерную выборку, сформированную ранее в лабораторной работе №4). Выполнить разбиение исходного множества объектов на конечное число подмножеств (кластеров) с использованием метода k-means. Полученные результаты содержательно проинтерпретировать.

### **Порядок выполнения работы**

- a) Нормализовать множество точек, отобразить полученное множество.
- b) Определить верхнюю оценку количества кластеров по формуле:  $\bar{k} = \lfloor \sqrt{N/2} \rfloor$ , где  $N$  – число точек.
- c) Реализовать алгоритм k-means в двух вариантах:
  - Пересчет центра кластера осуществляется после каждого изменения его состава
  - Пересчет центра кластера осуществляется после того, как будет завершен шаг процедуры

- d) На каждом шаге процедуры разбиения методом k-means вычислять функционалы качества полученного разбиения.
- e) Отобразить полученные кластеры, выделить каждый кластер разным цветом, отметить центроиды.
- f) Содержательно проинтерпретировать полученные результаты.

### Выполнение работы

- Нормирование

Исходная выборка представлена в таблице 1.

Таблица 1

№	<i>nu</i>	<i>E</i>	№	<i>nu</i>	<i>E</i>	№	<i>nu</i>	<i>E</i>	№	<i>nu</i>	<i>E</i>	№	<i>nu</i>	<i>E</i>
<b>1</b>	481	135.2	<b>21</b>	418	131.4	<b>41</b>	513	159.3	<b>61</b>	450	122.3	<b>81</b>	475	143.6
<b>2</b>	445	124.7	<b>22</b>	378	103.8	<b>42</b>	489	149.8	<b>62</b>	468	128.9	<b>82</b>	518	144.4
<b>3</b>	550	147.9	<b>23</b>	521	154.9	<b>43</b>	474	132.5	<b>63</b>	441	122.8	<b>83</b>	566	175.7
<b>4</b>	465	140.9	<b>24</b>	394	117.7	<b>44</b>	379	94.6	<b>64</b>	460	140.7	<b>84</b>	464	131.3
<b>5</b>	566	168.5	<b>25</b>	504	145.3	<b>45</b>	472	135.6	<b>65</b>	480	117.7	<b>85</b>	394	112.1
<b>6</b>	497	147.3	<b>26</b>	440	126.7	<b>46</b>	544	169.6	<b>66</b>	429	112.9	<b>86</b>	480	146.1
<b>7</b>	478	136.6	<b>27</b>	465	114.8	<b>47</b>	507	142.4	<b>67</b>	457	126.4	<b>87</b>	321	86.1
<b>8</b>	521	139.6	<b>28</b>	418	109.3	<b>48</b>	409	116.7	<b>68</b>	464	143.2	<b>88</b>	502	132.5
<b>9</b>	352	84.9	<b>29</b>	418	118.6	<b>49</b>	498	164.0	<b>69</b>	431	125.0	<b>89</b>	460	122.4
<b>10</b>	422	117.9	<b>30</b>	465	127.7	<b>50</b>	468	142.0	<b>70</b>	424	119.0	<b>90</b>	458	104.7
<b>11</b>	506	153.5	<b>31</b>	447	117.5	<b>51</b>	593	187.4	<b>71</b>	502	137.2	<b>91</b>	362	111.7
<b>12</b>	443	122.9	<b>32</b>	433	131.5	<b>52</b>	523	152.6	<b>72</b>	465	140.7	<b>92</b>	503	148.5
<b>13</b>	434	140.4	<b>33</b>	460	136.8	<b>53</b>	478	126.6	<b>73</b>	492	137.5	<b>93</b>	446	144.0
<b>14</b>	422	108.6	<b>34</b>	382	98.8	<b>54</b>	438	122.2	<b>74</b>	446	128.4	<b>94</b>	421	115.1
<b>15</b>	569	157.4	<b>35</b>	532	160.6	<b>55</b>	423	115.9	<b>75</b>	482	136.4	<b>95</b>	407	110.5
<b>16</b>	439	119.2	<b>36</b>	482	148.2	<b>56</b>	408	110.0	<b>76</b>	510	140.6	<b>96</b>	448	137.7
<b>17</b>	437	129.4	<b>37</b>	472	122.6	<b>57</b>	386	105.8	<b>77</b>	434	122.3	<b>97</b>	490	139.9
<b>18</b>	461	138.6	<b>38</b>	532	158.7	<b>58</b>	428	130.3	<b>78</b>	623	195.7	<b>98</b>	482	141.2
<b>19</b>	351	89.0	<b>39</b>	473	137.9	<b>59</b>	560	169.8	<b>79</b>	468	141.2	<b>99</b>	463	129.2
<b>20</b>	390	91.4	<b>40</b>	525	148.3	<b>60</b>	483	130.3	<b>80</b>	471	119.7	<b>100</b>	459	145.4

Нормализация была выполнена по методу минимакс:

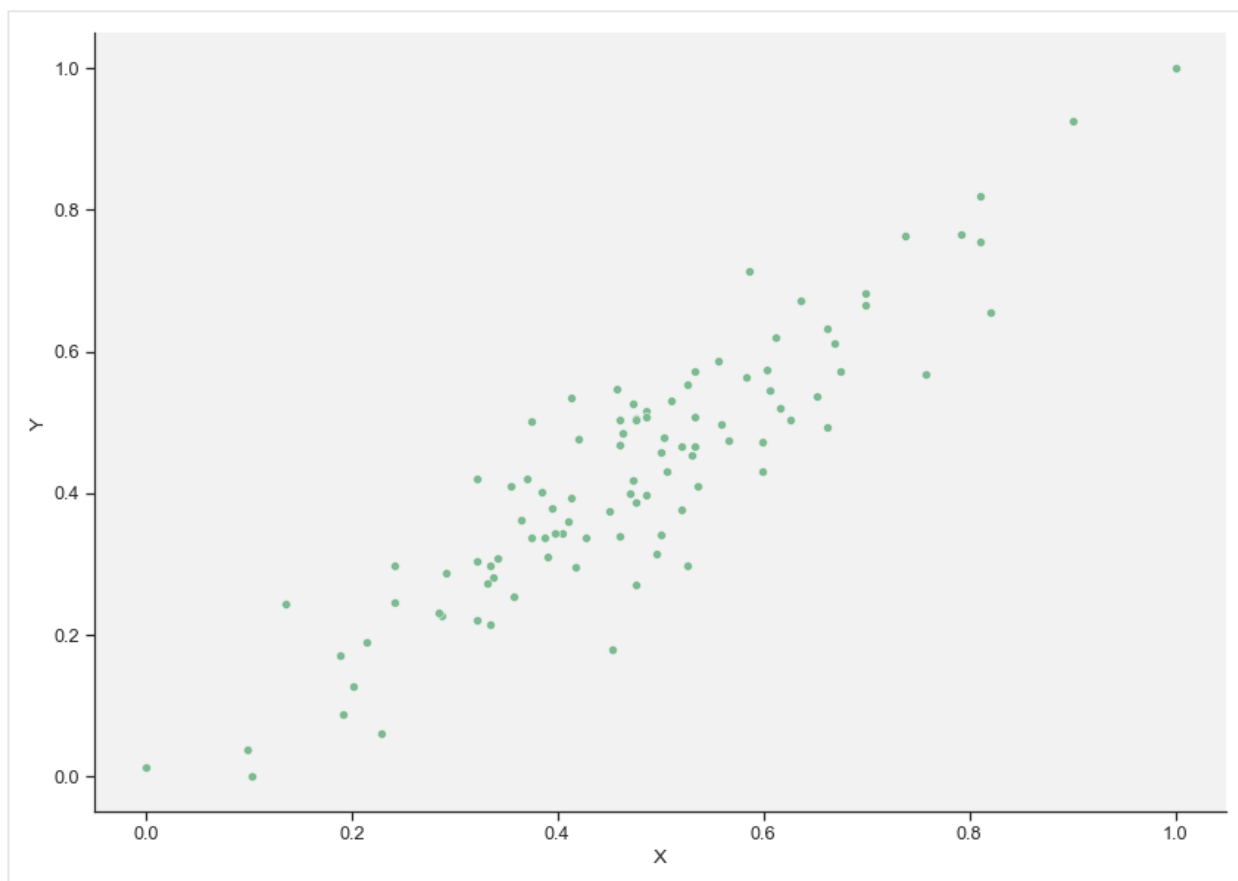
$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

После нормализации минимальное и максимальное масштабируемые значения равны 0 и 1 соответственно.

Нормализованная выборка представлена в табл. 2 и отображена на рис. 1.

Таблица 2

№	<i>nu</i>	<i>E</i>	№	<i>nu</i>	<i>E</i>	№	<i>nu</i>	<i>E</i>	№	<i>nu</i>	<i>E</i>	№	<i>nu</i>	<i>E</i>
<b>1</b>	0.53	0.454	<b>21</b>	0.321	0.42	<b>41</b>	0.636	0.671	<b>61</b>	0.427	0.338	<b>81</b>	0.51	0.53
<b>2</b>	0.411	0.359	<b>22</b>	0.189	0.171	<b>42</b>	0.556	0.586	<b>62</b>	0.487	0.397	<b>82</b>	0.652	0.537
<b>3</b>	0.758	0.569	<b>23</b>	0.662	0.632	<b>43</b>	0.507	0.43	<b>63</b>	0.397	0.342	<b>83</b>	0.811	0.819
<b>4</b>	0.477	0.505	<b>24</b>	0.242	0.296	<b>44</b>	0.192	0.088	<b>64</b>	0.46	0.504	<b>84</b>	0.474	0.419
<b>5</b>	0.811	0.755	<b>25</b>	0.606	0.545	<b>45</b>	0.5	0.458	<b>65</b>	0.526	0.296	<b>85</b>	0.242	0.245
<b>6</b>	0.583	0.563	<b>26</b>	0.394	0.377	<b>46</b>	0.738	0.764	<b>66</b>	0.358	0.253	<b>86</b>	0.526	0.552
<b>7</b>	0.52	0.467	<b>27</b>	0.477	0.27	<b>47</b>	0.616	0.519	<b>67</b>	0.45	0.375	<b>87</b>	0.0	0.011
<b>8</b>	0.662	0.494	<b>28</b>	0.321	0.22	<b>48</b>	0.291	0.287	<b>68</b>	0.474	0.526	<b>88</b>	0.599	0.43
<b>9</b>	0.103	0.0	<b>29</b>	0.321	0.304	<b>49</b>	0.586	0.714	<b>69</b>	0.364	0.362	<b>89</b>	0.46	0.338
<b>10</b>	0.334	0.298	<b>30</b>	0.477	0.386	<b>50</b>	0.487	0.515	<b>70</b>	0.341	0.308	<b>90</b>	0.454	0.179
<b>11</b>	0.613	0.619	<b>31</b>	0.417	0.294	<b>51</b>	0.901	0.925	<b>71</b>	0.599	0.472	<b>91</b>	0.136	0.242
<b>12</b>	0.404	0.343	<b>32</b>	0.371	0.421	<b>52</b>	0.669	0.611	<b>72</b>	0.477	0.504	<b>92</b>	0.603	0.574
<b>13</b>	0.374	0.501	<b>33</b>	0.46	0.468	<b>53</b>	0.52	0.376	<b>73</b>	0.566	0.475	<b>93</b>	0.414	0.533
<b>14</b>	0.334	0.214	<b>34</b>	0.202	0.125	<b>54</b>	0.387	0.337	<b>74</b>	0.414	0.393	<b>94</b>	0.331	0.273
<b>15</b>	0.821	0.654	<b>35</b>	0.699	0.683	<b>55</b>	0.338	0.28	<b>75</b>	0.533	0.465	<b>95</b>	0.285	0.231
<b>16</b>	0.391	0.31	<b>36</b>	0.533	0.571	<b>56</b>	0.288	0.227	<b>76</b>	0.626	0.503	<b>96</b>	0.421	0.477
<b>17</b>	0.384	0.402	<b>37</b>	0.5	0.34	<b>57</b>	0.215	0.189	<b>77</b>	0.374	0.338	<b>97</b>	0.56	0.496
<b>18</b>	0.464	0.485	<b>38</b>	0.699	0.666	<b>58</b>	0.354	0.41	<b>78</b>	1.0	1.0	<b>98</b>	0.533	0.508
<b>19</b>	0.099	0.037	<b>39</b>	0.503	0.478	<b>59</b>	0.791	0.766	<b>79</b>	0.487	0.508	<b>99</b>	0.47	0.4
<b>20</b>	0.228	0.059	<b>40</b>	0.675	0.572	<b>60</b>	0.536	0.41	<b>80</b>	0.497	0.314	<b>100</b>	0.457	0.546



*Рисунок 1 – Нормализованная выборка*

○ Оценка количества кластеров

Была найдена верхняя оценка количества кластеров:

$$\bar{k} = \left\lfloor \sqrt{\frac{N}{2}} \right\rfloor = \left\lfloor \sqrt{\frac{104}{2}} \right\rfloor = 7$$

○ Метод k-средних

Реализован метод k-средних для количества кластеров [2; 7]. Полученные кластеры были отображены, выделены разным цветом, были отмечены центроиды, вычислены функционалы качества разбиения. В таблицах представлены количество элементов в кластерах и их центры.

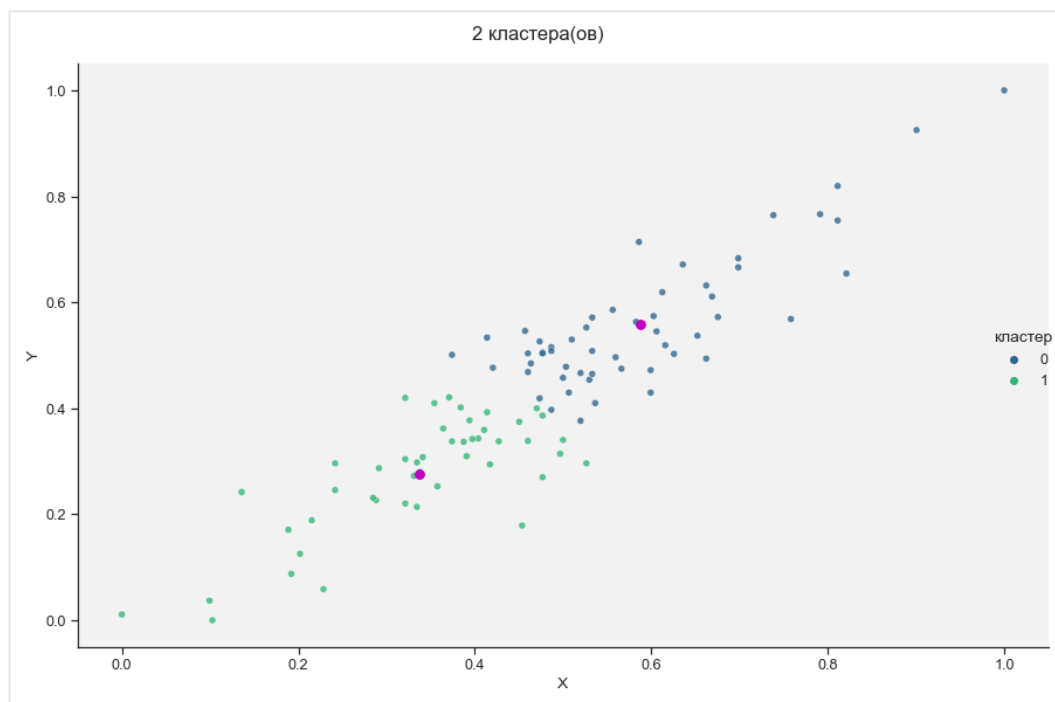


Рисунок 2 – 2 кластера (3 шага)

Таблица 3

Центр кластера		Количество элементов
0.5882	0.5593	54.0
0.3372	0.276	46.0

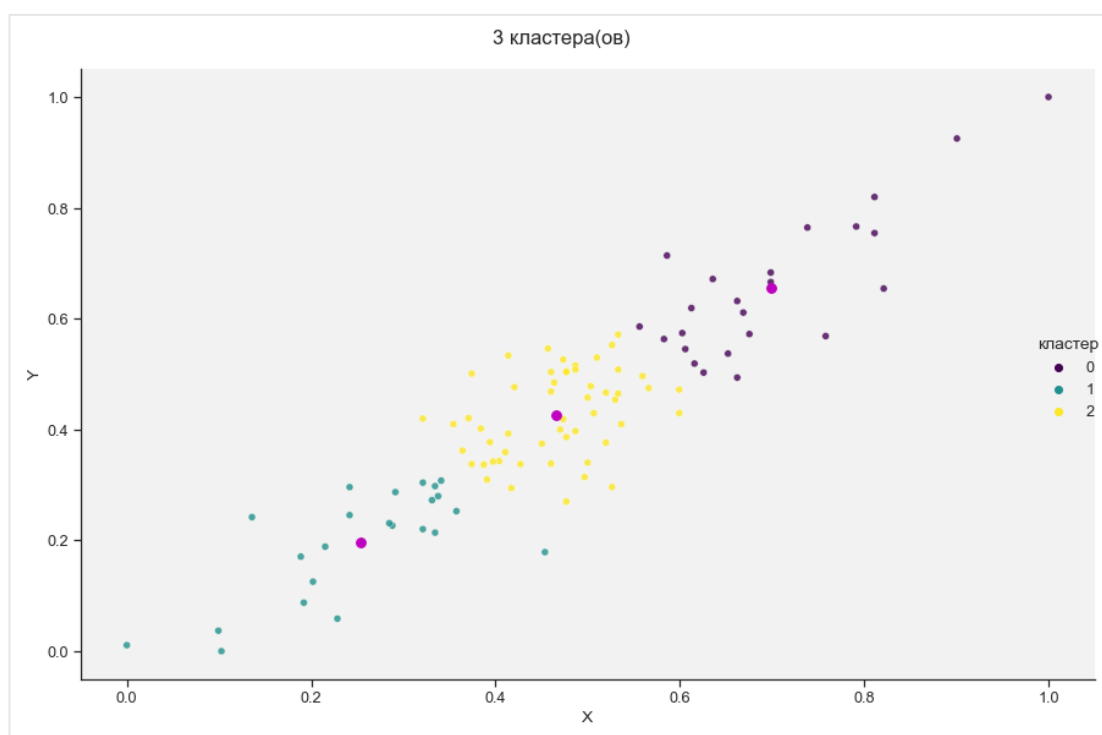


Рисунок 3 – 3 кластера (5 шагов)



Таблица 4

Центр кластера		Количество элементов
0.699	0.6559	24.0
0.2541	0.1971	23.0
0.4652	0.4268	53.0

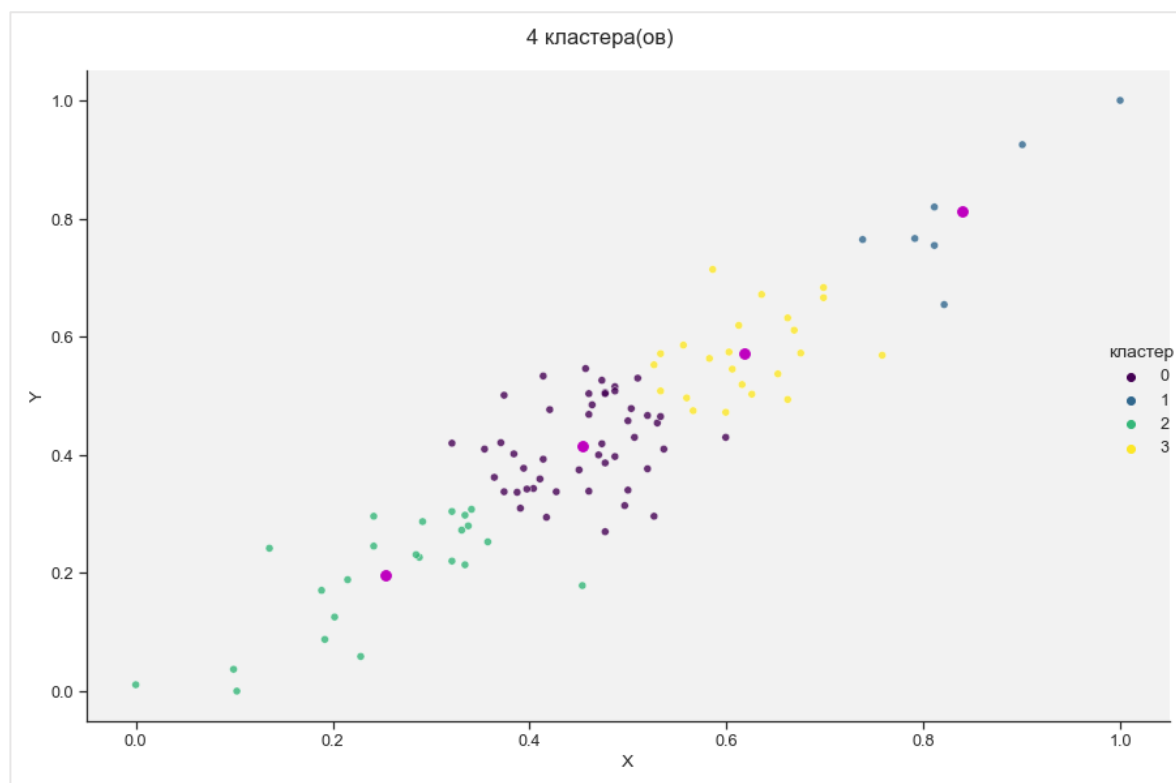


Рисунок 4 – 4 кластера (8 шагов)

Таблица 5

Центр кластера		Количество элементов
0.454	0.4159	47.0
0.8392	0.812	7.0
0.2541	0.1971	23.0
0.6182	0.571	23.0

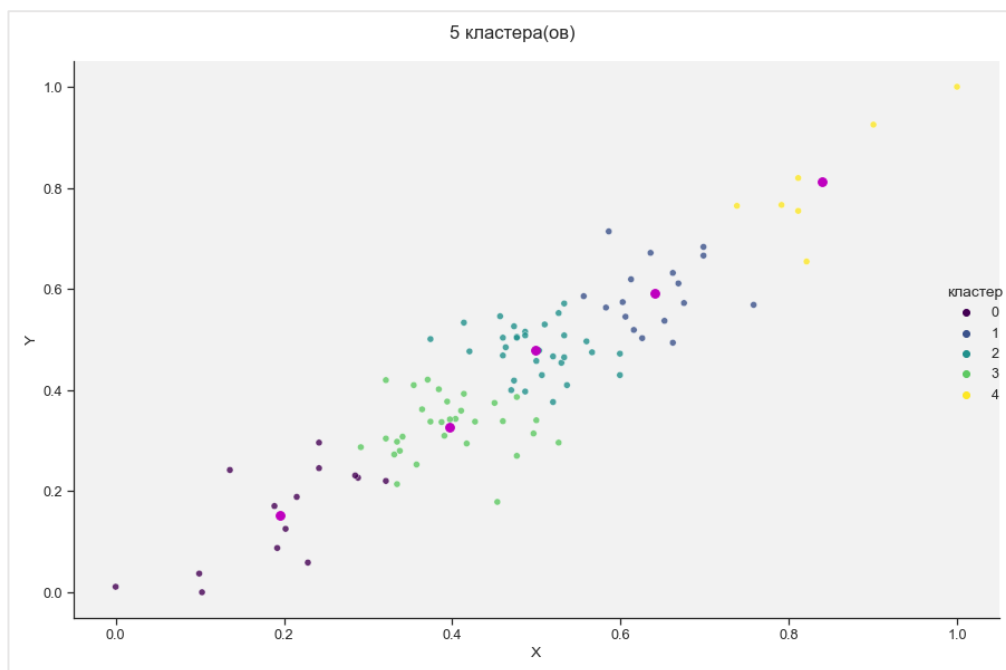


Рисунок 5 – 5 кластеров (10 шагов)

Таблица 6

Центр кластера		Количество элементов
0.1958	0.1528	14.0
0.6412	0.5916	17.0
0.4986	0.4793	31.0
0.3968	0.3276	31.0
0.8392	0.812	7.0

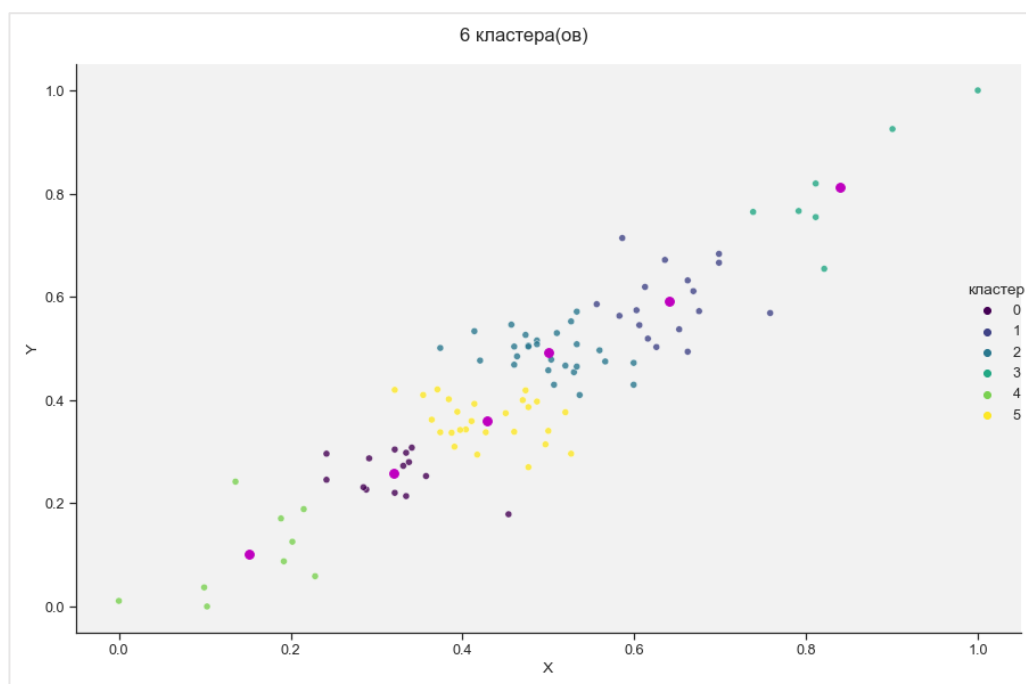


Рисунок 6 – 6 кластеров (11 шагов)

Таблица 7

Центр кластера		Количество элементов
0.32	0.2581	14.0
0.6412	0.5916	17.0
0.5002	0.4914	27.0
0.8392	0.812	7.0
0.1516	0.1023	9.0
0.4288	0.3598	26.0

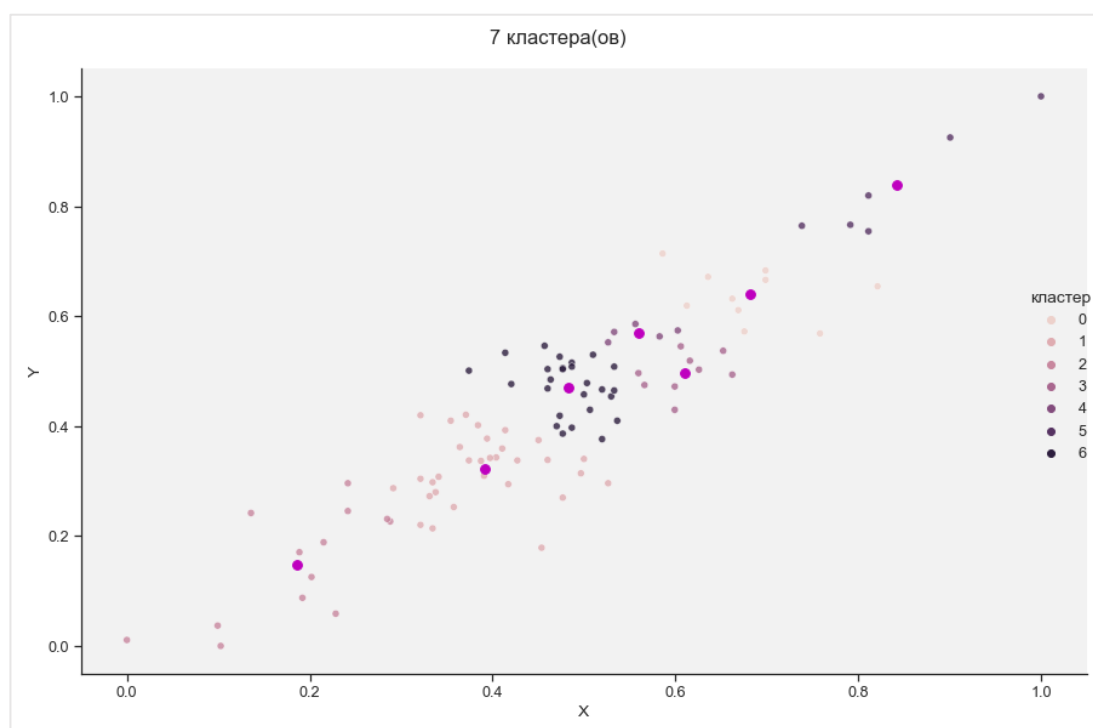


Рисунок 7 – 7 кластеров (16 шагов)

Таблица 8

Центр кластера		Количество элементов
0.6818	0.6392	10.0
0.3918	0.3223	31.0
0.1862	0.1477	13.0
0.6096	0.4967	9.0
0.5603	0.5693	5.0
0.8422	0.8383	6.0
0.4827	0.4711	26.0

○ Оценка качества разбиения

Для каждого разбиения были вычислены функционалы качества  $F_1, F_2, F_3$ , которые были определены в начале. В таблице 9 приведены значения функционалов на первой и последней итерациях алгоритма для разного количества кластеров разбиения.

Таблица 9

Количество кластеров	2	3	4	5	6	7
$F_1$	2.92	1.948	2.045	1.791	1.134	1.27
$F_1$ (конец)	2.891	1.645	1.171	0.812	0.625	0.56
$F_2$	151.736	103.754	95.325	75.691	22.247	29.21
$F_2$ (конец)	146.915	54.414	34.918	16.851	10.607	7.662
$F_3$	0.058	0.061	0.07	0.066	0.065	0.063
$F_3$ (конец)	0.058	0.06	0.058	0.057	0.054	0.055

Из таблицы можно увидеть, что при увеличении числа кластеров, минимизируются все функционалы, а также, то насколько сильно они меняются в сравнении с первой итерацией.

Алгоритм был реализован в двух вариантах: в первом, который был представлен выше, центр пересчитывается только по завершении шага процедуры, второй же вариант предполагает изменение центра кластера после обработки каждого объекта. Сравнение алгоритмов представлено в таблице 10.

Таблица 10

Количество кластеров	2	3	4	5	6	7
Количество итераций (первый алгоритм)	3	5	8	10	11	16
Количество итераций (второй алгоритм)	2	2	3	4	3	3

Из таблицы видно, что при увеличении количества кластеров увеличивается число итераций, а также, что количество итераций второго алгоритма

меньше, чем первого, что связано с тем, что центр меняется после обработки каждого объекта.

## **Выводы**

Освоены основные понятия кластерного анализа и метода k-means. Первоначальная двумерная выборка была нормализована методом минимакс, после которого минимальное и максимальное масштабируемые значения равны 0 и 1 соответственно. Нормализованная выборка была отображена на рисунке.

С помощью алгоритма k-средних было произведено разбиение выборки на 2, 3, ..., 7 кластеров. Для каждого разбиения были выведены центры кластеров и количество элементов в кластерах. Было оценено качество разбиения с помощью функционалов качества. В сравнительной таблице можно заметить, что при увеличении числа кластеров, минимизируются все функционалы качества, а также, то насколько сильно они меняются при сравнении первой и последней итераций.

Алгоритм k-средних был реализован в двух вариантах. В первом, центр пересчитывается только по завершении шага процедуры, второй же вариант предполагает изменение центра кластера после обработки каждого объекта. Из сравнительной таблицы можно заметить, что при увеличении количества кластеров увеличивается число итераций, а также, что количество итераций второго варианта алгоритма меньше, чем первого, это связано с тем, что центр меняется после обработки каждого элемента.

## ПРИЛОЖЕНИЕ А

### ИСХОДНЫЙ КОД

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import MinMaxScaler
from scipy.spatial import distance

df0 = pd.read_csv('data/main_data.csv')
X = df0['nu']
Y = df0['E']

X_norm = MinMaxScaler().fit_transform(df0)
df = pd.DataFrame(data=X_norm, columns=['nu', 'E'])
df.round(3).to_csv('data/df_norm.csv', index=False)

sns.set_theme(palette='crest', font_scale=1.15)
sns.set_style("ticks", {"axes.facecolor": ".95"})
ax = sns.relplot(data=df, x='nu', y='E', kind='scatter', height=8.27,
aspect=11.7/8.27)
ax.set_axis_labels('X', 'Y')
plt.tight_layout()
plt.savefig('pics/1.png')

N = len(df)
sup_k = np.floor(np.sqrt(N/2))
sup_k

def sc_plots(data, means, Ncl, step):
    if Ncl > 6:
        ax = sns.relplot(data=data, x='nu', y='E', hue='кластер',
kind='scatter', alpha=0.8,
                        height=8.27, aspect=11.7/8.27)
        plt.scatter(means[:,0],means[:,1], c='m', s=60)
    else:
        ax = sns.relplot(data=data, x='nu', y='E', hue='кластер',
kind='scatter', palette='viridis', alpha=0.8,
                        height=8.27, aspect=11.7/8.27)
        plt.scatter(means[:,0],means[:,1], c='m', s=60)

    print(f'кластеры = {Ncl}; шаги = {step}')
    ax.set_axis_labels('X', 'Y')
    ax.fig.suptitle(f'{Ncl} кластера(ов)')
    plt.tight_layout()
    plt.savefig(f'pics/2_{Ncl}.png')
    plt.close()
```

```

def nearest_center(data, cts):
    dist1 = np.array([], dtype=np.float64)
    for i in cts:
        dist1 = np.append(dist1, np.linalg.norm(i[:-1]-data)) # евклидово
расстояние
    min_dist = np.argmin(dist1)
    return min_dist

def Fs(data):
    curr_data = data.copy()
    cts = curr_data.groupby('кластер').mean()
    F1,F2,F3 = 0,0,0

    # F1 - сумма кв. расст. точек до центров соотв. кластеров
    for i in range(len(curr_data)):
        dist_F1 = np.linalg.norm(curr_data.iloc[i,:-1].values-
cts.values[curr_data.iloc[i,2]])
        F1 += dist_F1**2

    # F2 - сумма кв. расст. до всех точек соотв. кластеров
    for i in range(len(cts)):
        coords = curr_data[curr_data['кластер']==i].iloc[:,2].values
        dist_F2 = distance.cdist(coords, coords, 'euclidean')
        F2 += (np.triu(dist_F2,0)**2).sum()

    # F3 - сумма внутрикластерных дисперсий
    F3 = curr_data.groupby('кластер').var().values.sum()

    return F1,F2,F3

def custKM(dataf, n_clusters, chng_ctr=1, max_iter=30, tol=0.01):
    data = dataf.copy()
    centers = data.sample(n_clusters) # случайные центры
    data['кластер'] = -1 # нет принадлежности кластерам
    cts = np.array([], dtype=np.float64)
    F1,F2,F3 = 0,0,0
    df_Fs = pd.DataFrame(columns=['F1', 'F2', 'F3'])

    for i in range(n_clusters):
        data.loc[centers.index[i],'кластер'] = i # кластеры для центров
        cts = np.append(cts, [data.loc[centers.index[i]].values])
    centers = cts.reshape((n_clusters,3))

    for j in range(max_iter):
        for i in range(len(data)): # ближ. центр для каждой точки
            curr_clust = nearest_center(data.iloc[i,:-1].values, centers)
            data.loc[i,'кластер'] = curr_clust # соотносим кластер

```

```

        if chng_ctr: # пересчет центра при новой точке
            centers[curr_clust][:2] =
data[data['кластер']==curr_clust].iloc[:,2].mean()

    if chng_ctr == 0: # пересчет центра на каждой итерации
        for i in range(n_clusters):
            centers[i][:2] = data[data['кластер']==i].iloc[:,2].mean()

    cur_F1,cur_F2,cur_F3 = Fs(data) # функционалы
    df_Fs = df_Fs.append({'F1':cur_F1,'F2':cur_F2,'F3':cur_F3},
ignore_index=True)

    if np.abs(F1-cur_F1) < tol:
        data['кластер'].astype('int')
        sc_plots(data, centers, n_clusters, j+1)
        break
    F1,F2,F3 = cur_F1,cur_F2,cur_F3
    data['кластер'] = -1

    df_ctrs = pd.DataFrame(np.concatenate((centers[:,2],
data.groupby('кластер')['nu'].count().values.reshape(-1,1)), axis=1),
        columns=['nu_mean', 'E_mean', 'num'])
    silhouette_avg = silhouette_score(data.values[:,2], data.values[:,2])

    return df_Fs, df_ctrs, silhouette_avg

for i in range(2,8):
    F, ctrs, sil = custKM(df, n_clusters=i, chng_ctr=1)
    F.round(3).to_csv(f'data/Fs_{i}c.csv', index=False)
    ctrs.round(4).to_csv(f'data/centers_{i}c.csv', index=False)

```