

Домашняя работа по кластеризации Борисова Алексея из 201 группы.

```
import numpy as np
from scipy.spatial import distance
import matplotlib.pyplot as plt

In [153].
# Функция выбора центроиды для инициализации методом k-means++
def furthest_sample(centers, X):
    dist = np.min(distance.cdist(centers, X)**2, axis=0)
    Sum = dist.sum()
    ind = np.random.choice(n * Sum
        cur_sum = 0
        i = 1
        while cur_sum <= Rnd:
            cur_sum += Dist[i]
            i += 1
        return X[i]

# Функция вычисления суммы средних внутрикластерных расстояний
def inner_dist(X, centers, cluster):
    sum = 0
    for i in range(centers.shape[0]):
        Dist = distance.cdist(X[cluster == i], centers[i][np.newaxis,:])**2
        if Dist.shape[0] == 0:
            return np.inf
        sum += Dist.mean()
    return sum

# Функция вычисления среднего межкластерного расстояния
def between_dist(X, centers):
    all_center = X.mean(axis=0)
    return 1 / centers.shape[0] * (distance.cdist(centers, all_center[np.newaxis,:])**2).sum()

# Функционал качества, который будем минимизировать.
# Он зависит от оптимального числа кластеров
def Q(X, centers, cluster):
    between = between_dist(X, centers)
    if between == 0:
        return np.inf
    return inner_dist(X, centers, cluster) / between

class MyKMeans:
    # Если число кластеров не задано или задано равным 0,
    # то алгоритм сам его подбирает.
    def __init__(self, n_clusters=0, init="random"):
        self.n_clusters = n_clusters
        self.init = init

    def fit_predict(self, X):
        # случай заданное число кластеров
        if self.n_clusters != 0:
            return self.predict(X, self.n_clusters)
        # Если подобрать оптимальное число кластеров самостоятельно
        n_clusters = 1
        prev_Q = np.inf
        # Даже используя способ инициализации k-means++
        # первый центр выбирается случайно.
        # чтобы минимизировать эту случайность будем брать
        # среднее значение функционала качества
        # за некоторое количество испытаний, например, 50.
        Q = np.zeros(50)
        while True:
            for j in range(50):
                pred = self.predict(X, n_clusters)
                Q[j] = pred[2]
                cur_Q = Q.mean()
            # Если функционал качества при n кластерах оказался больше,
            # чем при n-1, то останавливаемся
            if cur_Q > prev_Q:
                return (cluster, it, prev_Q, diff, centers)
                cluster = pred[0]
                it = pred[1]
                prev_Q = cur_Q
                diff = pred[3]
                centers = pred[4]
                n_clusters += 1

# Метод разбиения выборки на заранее заданное число кластеров
def predict(self, X, n_clusters):
    n_samples = X.shape[0]
    all_features = X.shape[2]
    cluster = np.zeros(n_samples)
    if self.init == "random":
        centers = X[np.random.choice(n_samples, n_clusters, replace=False)].astype(float)
    elif self.init == "k-means++":
        arg = np.random.choice(n_samples, 1)
        centers = X[arg]
        for i in range(n_clusters - 1):
            centers = np.vstack((centers, furthest_sample(centers, X)))
        differences = [1] # в массив значений на сколько на каждой итерации сдвинулись центры
        dif = 1 # переменная для подсчета на сколько сдвинулись центры на итерации
        while dif != 0:
            it = 0
            for i in range(n_samples):
                cluster[i] = np.argmax(distance.cdist(X[i][np.newaxis,:], centers).ravel())
                for j in range(n_clusters):
                    prev_center = centers[j].copy()
                    if (X[cluster == j].shape[0] != 0):
                        centers[j] = (X[cluster == j].sum(axis=0) / X[cluster == j].shape[0])
                        dif += distance.minkowski(prev_center, centers[j])
            differences.append(dif)
            # возвращаем разбиение, число итераций, функционал качества,
            # изменение координат центров, итоговые центры кластеров.
            return (cluster, it, Q(X, centers, cluster), differences, centers)

In [145].
# Реализация генератора k кластеров с Гауссовым распределением для тестирования алгоритма
def generate_gauss(k, X):
    size = int(N/k)
    X = np.zeros(2 * N)
    for i in range(k):
        c = (np.random.uniform(-1, 1), np.random.uniform(-1, 1))
        s = np.random.uniform(0.05, 0.2)
        for j in range(size):
            X[id1] = np.random.normal(c[0], s)
            X[id1 + 1] = np.random.normal(c[1], s)
        id1 += 2
    return X.reshape(-1, 2)

In [152].
# Другой генератор более плотных скопления точек
def make_blob(n, k):
    size = int(N/k)
    X = np.zeros(2 * N)
    for i in range(k):
        cx = 2 * np.random.random() - 1
        cy = 2 * np.random.random() - 1
        for j in range(size):
            X[id1] = cx + 0.4 * np.random.random() - 0.2
            X[id1 + 1] = cy + 0.4 * np.random.random() - 0.2
        id1 += 2
    return X.reshape(-1, 2)
```

Исследование сходимости метода и зависимости сходимости от стратегии начальной инициализации

```
In [144].
fig = plt.figure()
ax1 = fig.add_subplot(121)
ax2 = fig.add_subplot(122)
fig.set_figheight(5)
fig.set_figwidth(15)
n_clusters = np.arange(2, 100)

for n in n_clusters:
    X = make_blobs(1000, n)
    clf = MyKMeans(n_clusters=n, init="random")
    dif = clf.fit_predict(X)[3]
    ax1.set_title("random")
    ax1.set_xlabel("номер итерации")
    ax1.axis([-0.1, 40, 0, 5.5])
    ax1.plot(range(len(dif)), dif)

for n in n_clusters:
    X = make_blobs(1000, n)
    clf = MyKMeans(n_clusters=n, init="k-means++")
    dif = clf.fit_predict(X)[3]
    ax2.set_title("k-means++")
    ax2.set_xlabel("изменения положения центров")
    ax2.set_xlabel("номер итерации")
    ax2.axis([-0.1, 40, 0, 5.5])
    ax2.plot(range(len(dif)), dif)
```



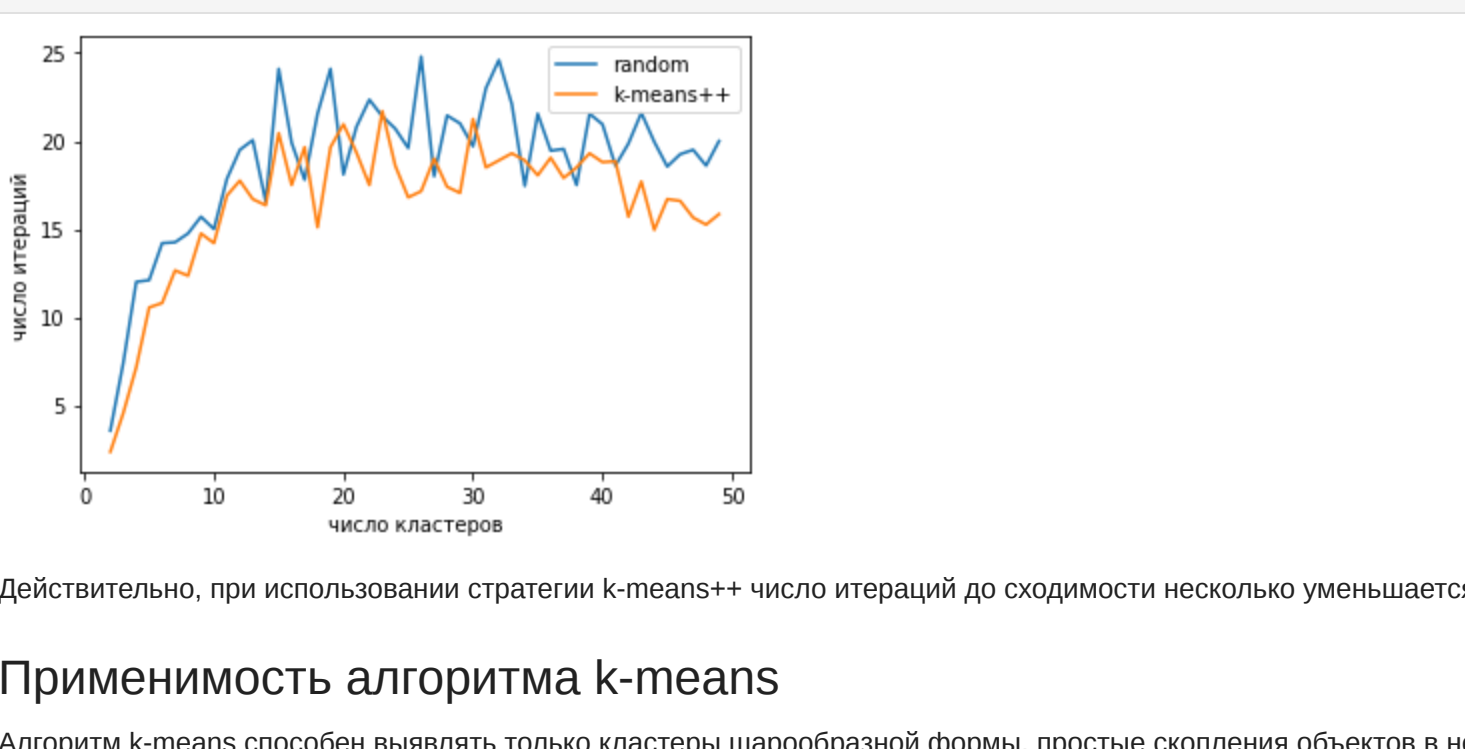
Из графиков видно, что при обеих стратегиях начальной инициализации алгоритм сходится. Причём хочется отметить, что при использовании стратегии k-means++ графики более плавные, меньше сильных скачков вверх, которые наблюдаются при случайном выборе начальных центров, что является подтверждением того, что k-means++ является более стабильным. Также уже из этих графиков можно заметить некоторое уменьшение числа итераций до сходимости.

```
In [144].
n_clusters = np.arange(2, 50)
mean = np.zeros(20)
arr_random = np.zeros(n_clusters.size)

for n in n_clusters:
    for i in range(20):
        X = make_blobs(1000, n)
        clf = MyKMeans(n_clusters=n, init="random")
        mean[i] = clf.fit_predict(X)[3]
        arr_random[n - 2] = mean.mean()

In [144].
arr_kmeans = np.zeros(n_clusters.size)

for n in n_clusters:
    for i in range(20):
        X = make_blobs(1000, n)
        clf = MyKMeans(n_clusters=n, init="k-means++")
        mean[i] = clf.fit_predict(X)[3]
        arr_kmeans[n - 2] = mean.mean()
```



Действительно, при использовании стратегии k-means++ число итераций до сходимости несколько уменьшается.

Применимость алгоритма k-means

Алгоритм k-means способен выявлять только кластеры шарообразной формы, простые скопления объектов в некоторой области, для поиска кластеров более сложной формы он не приспособлен. Желательно знать количество кластеров, на которые мы хотим разбить выборку, так как хотя есть различные способы выбрать лучшее разбиение, они все не идеальны и допускают ошибки, а неверный выбор числа кластеров может привести к плохим результатам.

Стратегия выбора числа кластеров

Если число кластеров заранее неизвестно, можно попробовать разбить выборку на различное число кластеров и посмотреть какое из разбиений лучше. Для этого нужно иметь некоторую функцию, по значению которой можно будет понять какое из разбиений лучше, так называемый функционал качества.

Выбирая такой функционал в своем алгоритме я руководствовался несколькими разумными требованиями к желаемому разбиению. 1) Элементы внутри одного кластера должны быть "похожи", находиться на небольшом расстоянии. 2) Элементы относящиеся к разным кластерам должны быть различны, находиться на большом расстоянии.

Переведем на математический язык, я хочу минимизировать внутрикластерное расстояние для каждого кластера и максимизировать межкластерное расстояние.

Для минимизации внутрикластерного расстояния будем минимизировать сумму средних внутрикластерных расстояний по всем кластерам:

$$Q_1 = \sum_{c \in C} \frac{1}{|S_c|} \sum_{x \in S_c} \rho^2(x, \mu_c)$$

C – множество кластеров, $|S_c|$ – число элементов в одном кластере, μ_c – центр кластера c .

Для максимизации межкластерного расстояния будем максимизировать среднее межкластерное расстояние:

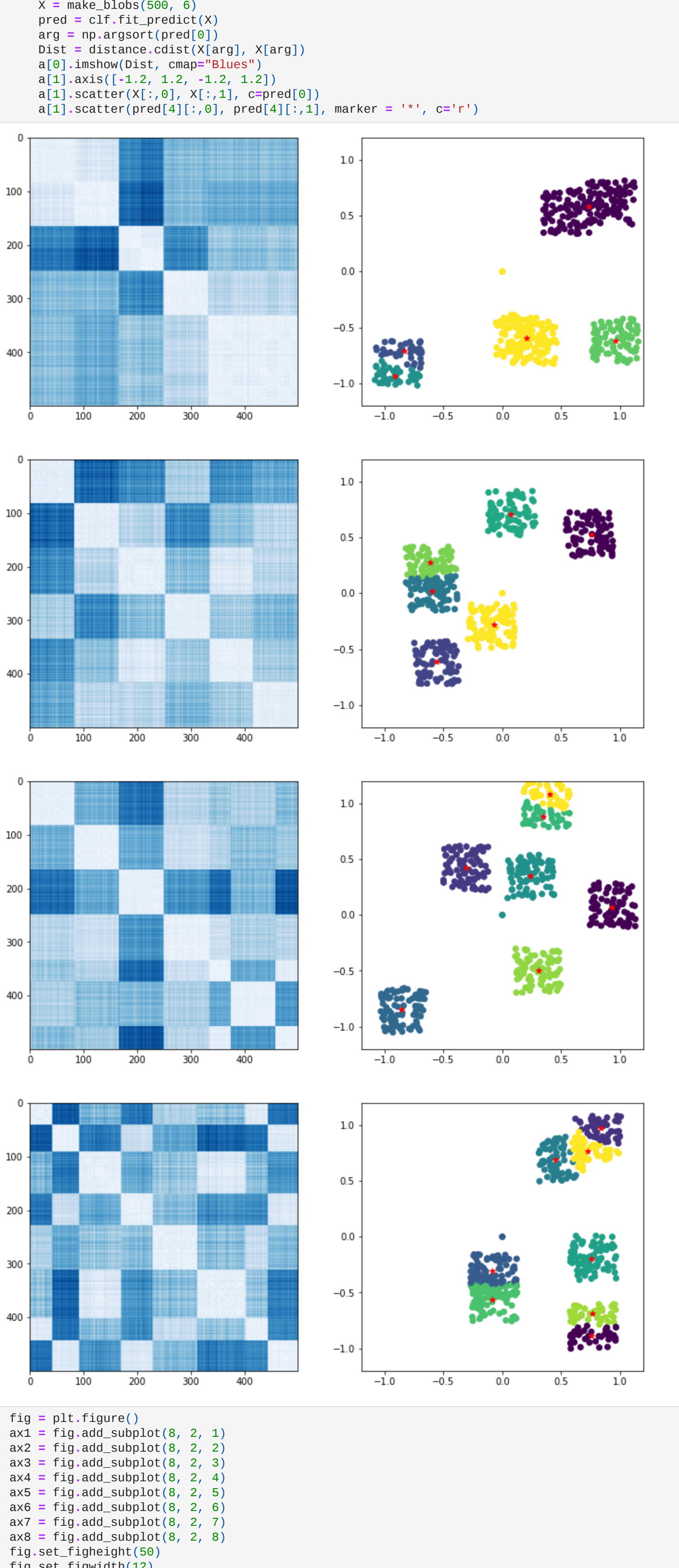
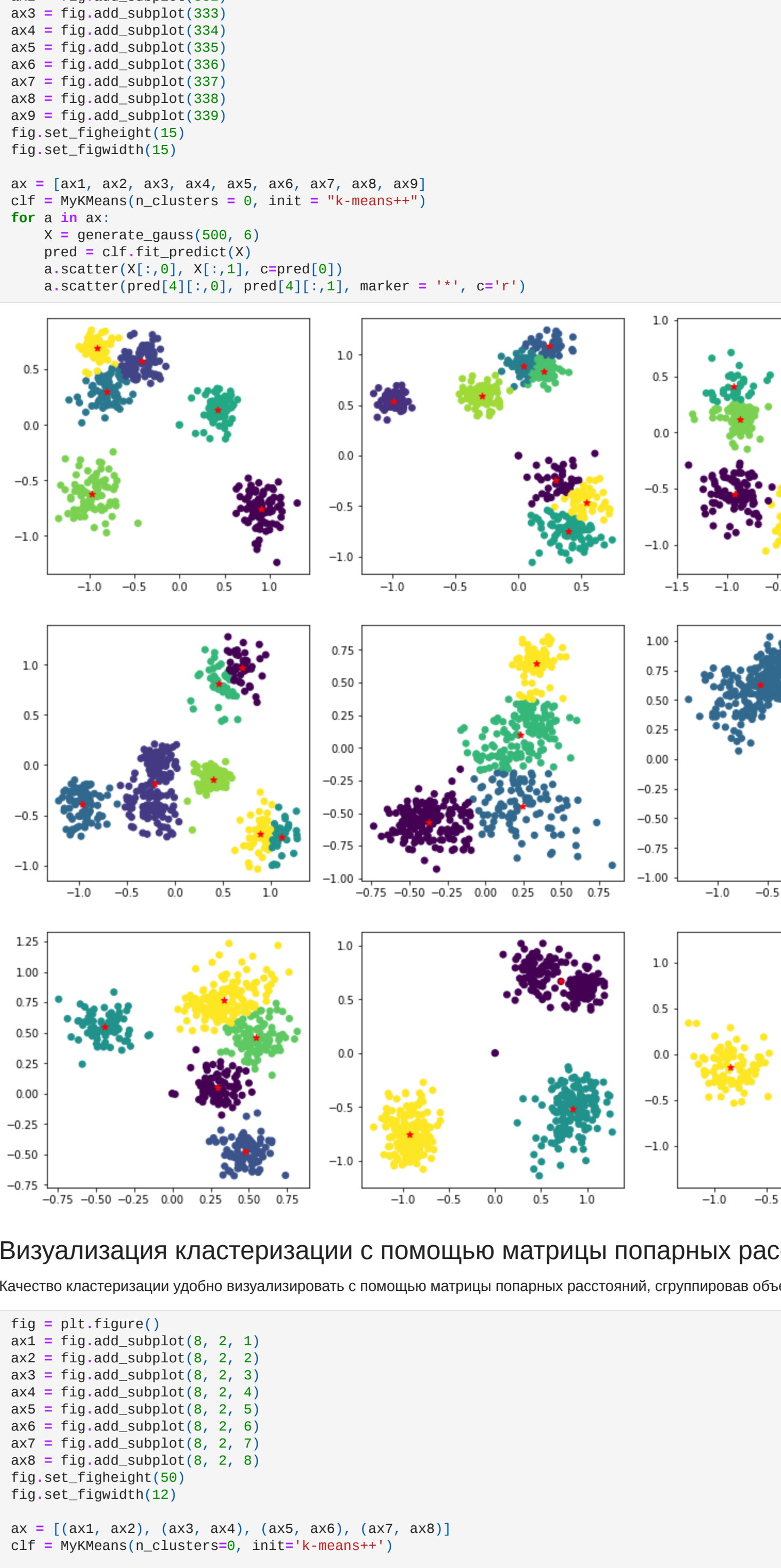
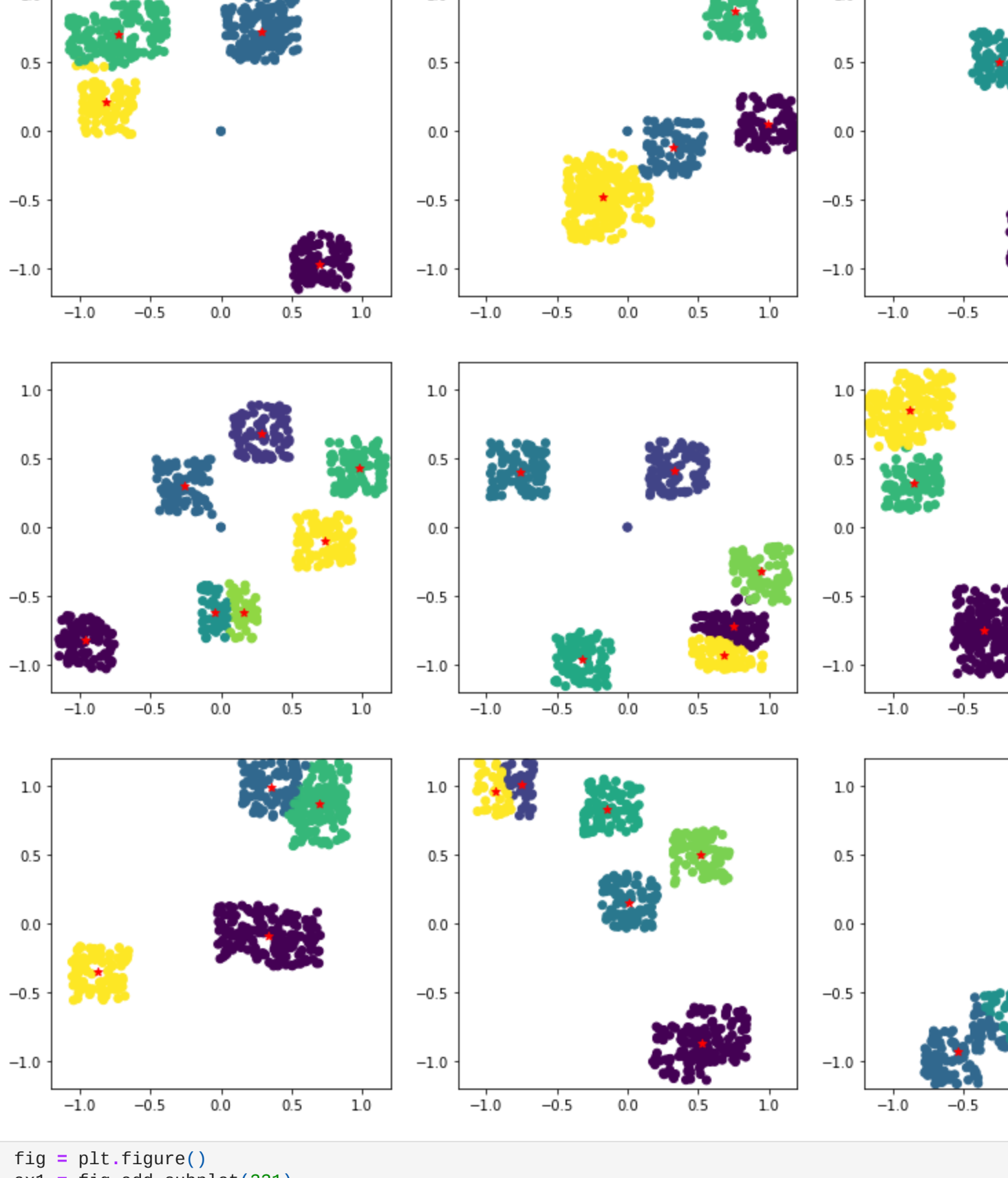
$$Q_2 = \frac{1}{|C|} \sum_{c \in C} \rho^2(\mu_c, \mu)$$

μ – центр всей выборки.

$$Q = \frac{Q_1}{Q_2}$$

Этот функционал я буду минимизировать при поиске оптимального числа кластеров.

Протестируем данную стратегию поиска числа кластеро на обоих генераторах.



Визуализация кластеризации с помощью матрицы попарных расстояний

Каждому кластеризации удобно визуализировать с помощью матрицы попарных расстояний, структурированных объекты по кластерам. Ниже это реализовано и протестировано.

