

Домашняя работа по кластеризации Борисова Алексея из 201 группы.

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

In [151]. # 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.argmax(dist)
    cur_sum = 0
    i = 0
    while cur_sum <= ind:
        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
# при поиске 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 кластера не задано или задано равным 0,
    # то автоматически 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 заданное число кластеров
        if self.n_clusters == 0:
            return self.predict(X, self.n_clusters)
        # будем подбирать оптимальное число кластеров самостоятельно
        n_clusters = 1
        cur_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-2, 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

# Итерация разбиений выборку на заранее заданное число кластеров
def predict(self, X, n_clusters):
    n_samples = X.shape[0]
    n_clusters = n_clusters
    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 = []
        # Если храним в массиве на каждой итерации сходимости центры
        dif = 1
        # переменная для подсчета на сколько сдвинулись центры на итерации
        while dif != 0:
            dif = 1
            for i in range(n_clusters):
                cluster[i] = np.argmin(distance.cdist(X[i][np.newaxis,:], centers).ravel())
            for i in range(n_clusters):
                prev_center = centers[i].copy()
                if (X[cluster == i].shape[0] != 0):
                    centers[i] = X[cluster == i].mean(axis=0)
                    dif = distance.minkowski(prev_center, centers[i])
            differences.append(dif)
        # возвращаем разбиение, число итераций, функционал качества, центры кластеров.
        return (cluster, it, Q(X, centers, cluster), differences, centers)

In [145. # Небольшой генератор k кластеров с гауссовым распределением для тестирования алгоритма
def generate_gauss(N, k):
    size = int(N/k)
    X = np.zeros(2 * N)
    idx = 0
    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[idx] = np.random.normal(c[0], s)
            X[idx+1] = np.random.normal(c[1], s)
            idx += 2
    return X.reshape(-1, 2)

In [152. # Другой генератор более плотных скопления точек
def make_blobs(N, k):
    size = int(N/k)
    X = np.zeros(2 * N)
    idx = 0
    for i in range(k):
        cx = 2 * np.random.random() - 1
        cy = 2 * np.random.random() - 1
        for j in range(size):
            X[idx] = cx + 0.4 * np.random.random() - 0.2
            X[idx+1] = cy + 0.4 * np.random.random() - 0.2
            idx += 2
    return X.reshape(-1, 2)
```

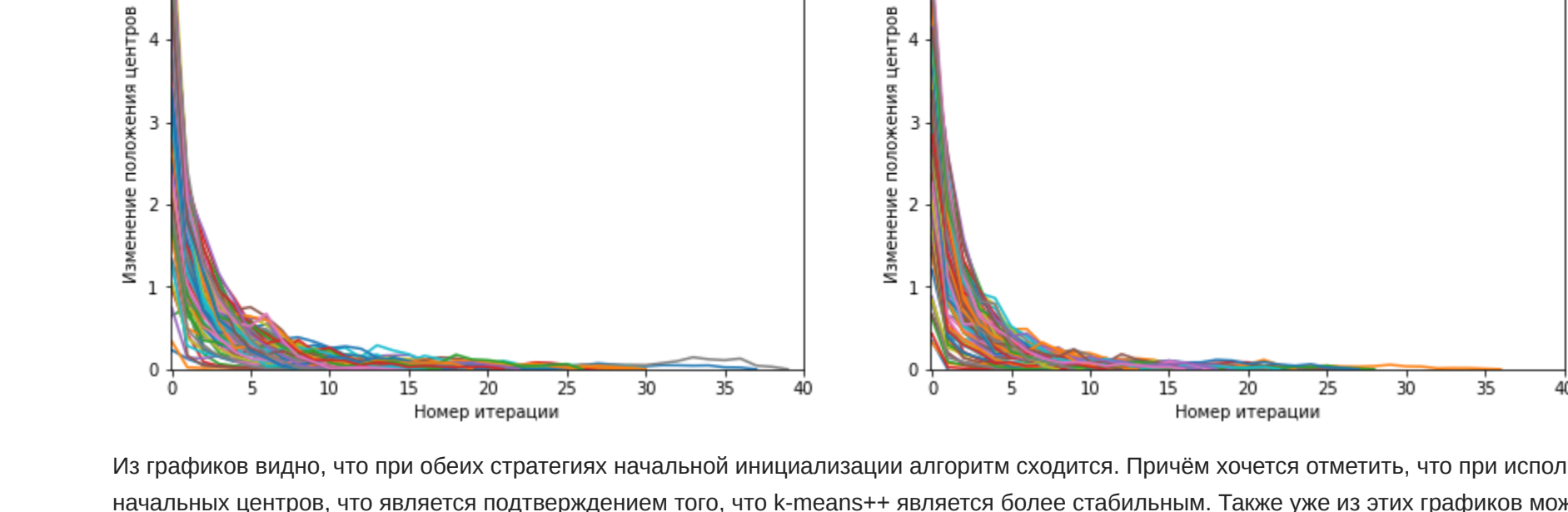
Исследование сходимости метода и зависимости сходимости от стратегии начальной инициализации

```
In [144. fig = plt.figure()
ax = fig.add_subplot(121)
ax2 = fig.add_subplot(122)
fig.set_figheight(15)
fig.set_figwidth(35)

n_clusters = np.arange(2, 100)

for n in n_clusters:
    X = make_blobs(1800, n)
    clf = MyKMeans(n_clusters=n, init="random")
    dif = clf.fit_predict(X)[3]
    ax1.set_title('random')
    ax1.set_ylabel('Изменение пополения центров')
    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(1800, n)
    clf = MyKMeans(n_clusters=n, init="k-means++")
    dif = clf.fit_predict(X)[3]
    ax2.set_title('k-means++')
    ax2.set_ylabel('Изменение пополения 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.

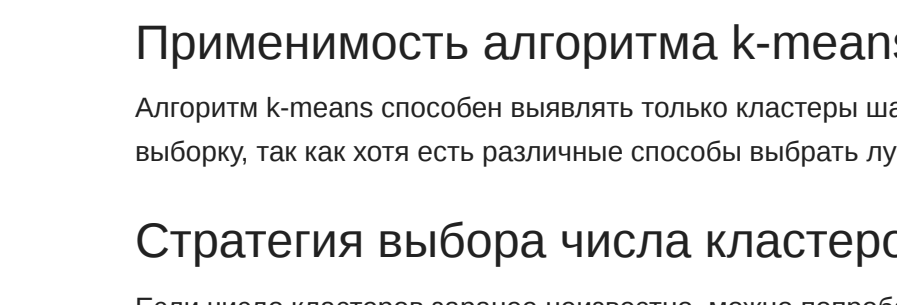
```
In [144. n_clusters = np.arange(2, 50)
mean = np.zeros(n_clusters)
arr_kmeans = np.zeros(n_clusters)

for n in n_clusters:
    for j in range(20):
        X = make_blobs(1800, n)
        clf = MyKMeans(n_clusters=n, init="random")
        mean[j] = clf.fit_predict(X)[1]
        arr_kmeans[n - 2] = mean.mean()

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

for n in n_clusters:
    for j in range(20):
        X = make_blobs(1800, n)
        clf = MyKMeans(n_clusters=n, init="k-means++")
        mean[j] = clf.fit_predict(X)[1]
        arr_kmeans[n - 2] = mean.mean()
```

```
In [144. plt.xlabel("число кластеров")
plt.ylabel("число итераций")
plt.plot(n_clusters, arr_kmeans, label="random")
plt.plot(n_clusters, arr_kmeans, label="k-means++")
plt.legend()
fig.set_figheight(30)
fig.set_figwidth(30)
```



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 которое мы хотим разбить выборку, так как хотя есть разные способы выбрать лучшее разбиение, они все не идеальны и допускают ошибки.

Стратегия выбора числа кластеров

Если число кластеров задано неизвестно, можно попробовать разбить выборку на разное 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.

Выбрав такой функционал в своем алгоритме и руководствуясь 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 разными, находиться на большом расстоянии.

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

Для минимизации внутрикластерного расстояния буду 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 буду максимизировать среднее межкластерное расстояние:

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

μ – центр всей выборки.

В качестве итогового функционала качества возьму следующую функцию:

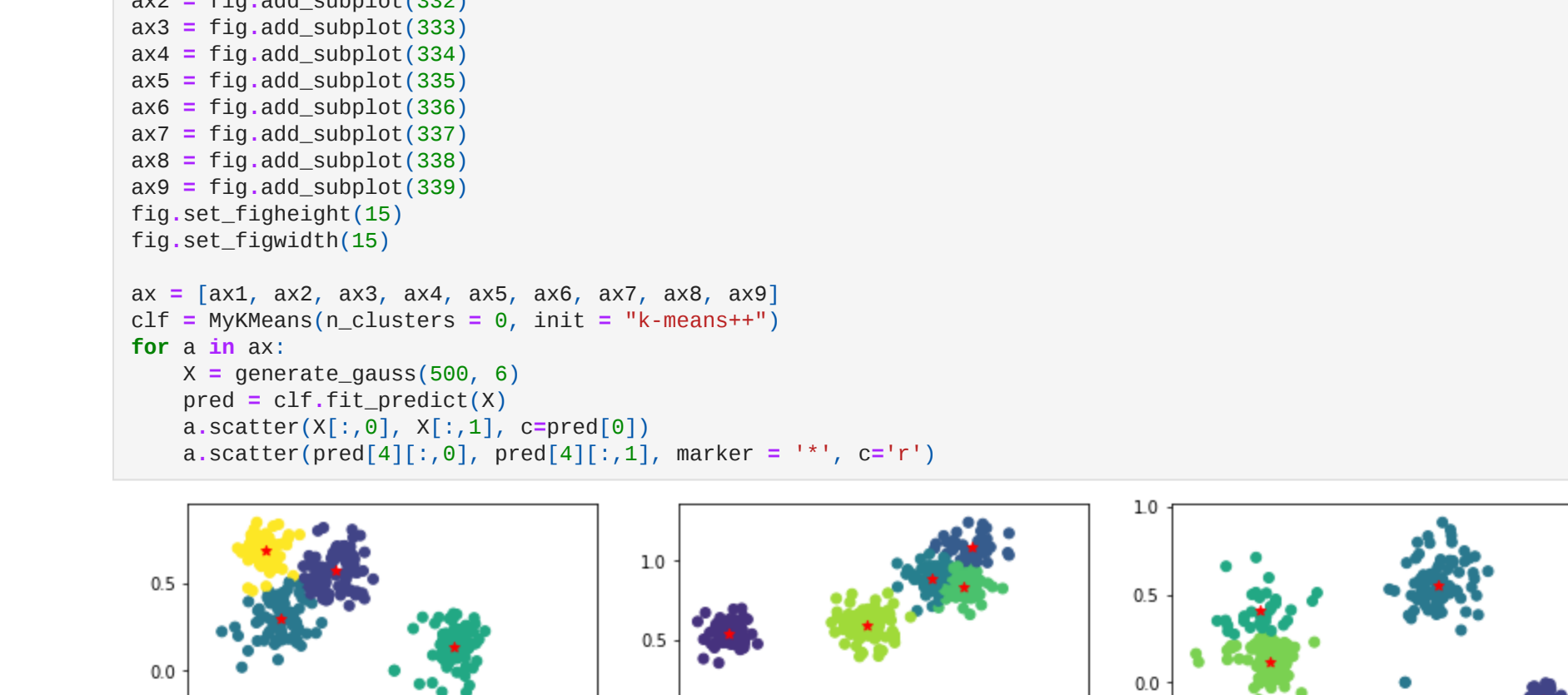
$$Q = \frac{Q_1}{Q_2}$$

Этой функцией и буду минимизировать при поиске оптимального числа кластеров.

Протестирую данную стратегию поиска числа кластеро на обоях генераторах.

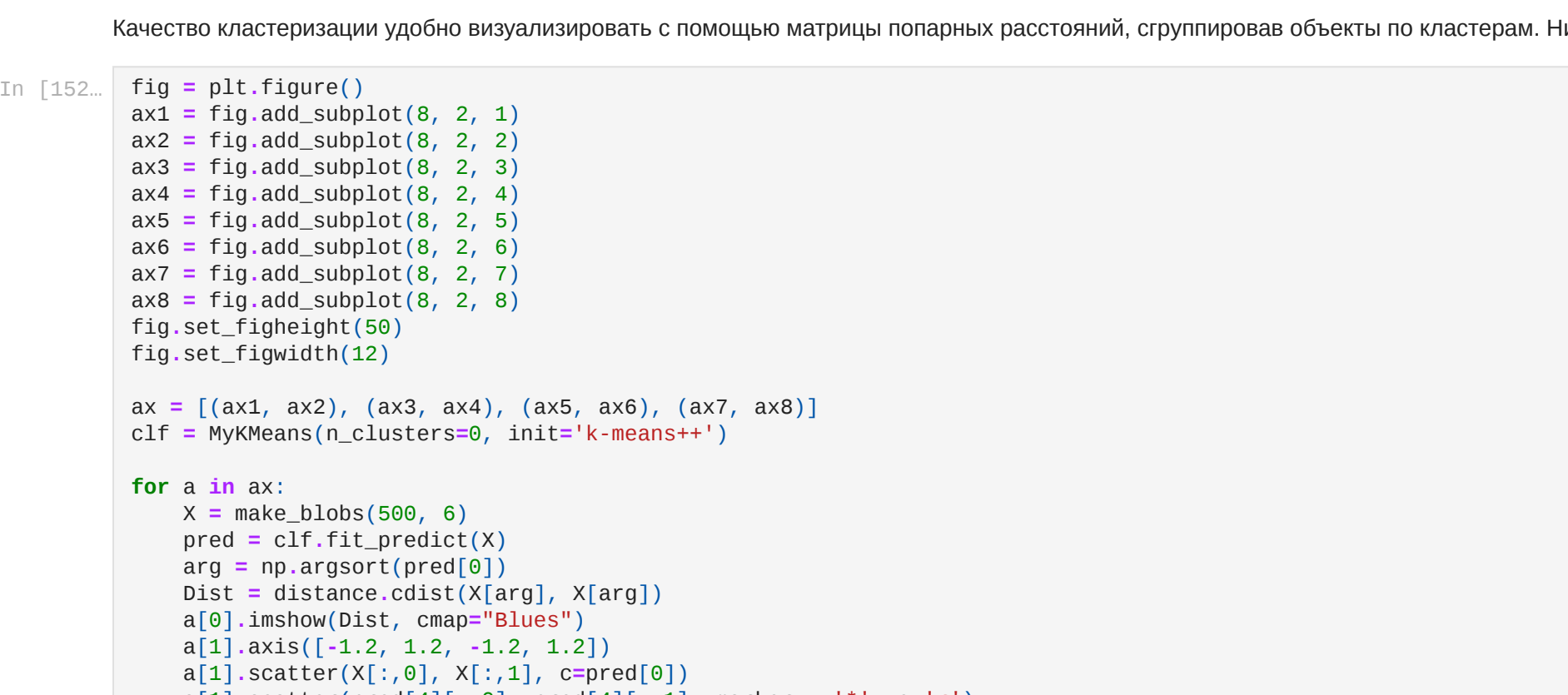
```
In [153. fig = plt.figure()
ax1 = fig.add_subplot(331)
ax2 = fig.add_subplot(332)
ax3 = fig.add_subplot(333)
ax4 = fig.add_subplot(334)
ax5 = fig.add_subplot(335)
ax6 = fig.add_subplot(336)
ax7 = fig.add_subplot(337)
ax8 = fig.add_subplot(338)
ax9 = fig.add_subplot(339)
fig.set_figheight(15)
fig.set_figwidth(15)

ax = [ax1, ax2, ax3, ax4, ax5, ax6, ax7, ax8, ax9]
clf = MyKMeans(n_clusters = 0, init = "k-means++")
for a in ax:
    X = make_blobs(500, 6)
    pred = clf.fit_predict(X)
    scatter(X[:,0], X[:,1], c=pred[0])
    a.axis([-1.2, 1.2, -1.2, 1.2])
    a.scatter(pred[4][:,0], pred[4][:,1], marker = '+', c='r')
```



```
In [145. fig = plt.figure()
ax1 = fig.add_subplot(331)
ax2 = fig.add_subplot(332)
ax3 = fig.add_subplot(333)
ax4 = fig.add_subplot(334)
ax5 = fig.add_subplot(335)
ax6 = fig.add_subplot(336)
ax7 = fig.add_subplot(337)
ax8 = fig.add_subplot(338)
ax9 = fig.add_subplot(339)
fig.set_figheight(15)
fig.set_figwidth(15)

ax = [ax1, ax2, ax3, ax4, ax5, ax6, ax7, ax8, ax9]
clf = MyKMeans(n_clusters = 0, init = "k-means++")
for a in ax:
    X = generate_gauss(500, 6)
    pred = clf.fit_predict(X)
    scatter(X[:,0], X[:,1], c=pred[0])
    a.scatter(pred[4][:,0], pred[4][:,1], marker = '+', c='r')
```



```
In [152. fig = plt.figure()
ax1 = fig.add_subplot(8, 2, 1)
ax2 = fig.add_subplot(8, 2, 2)
ax3 = fig.add_subplot(8, 2, 3)
ax4 = fig.add_subplot(8, 2, 4)
ax5 = fig.add_subplot(8, 2, 5)
ax6 = fig.add_subplot(8, 2, 6)
ax7 = fig.add_subplot(8, 2, 7)
ax8 = fig.add_subplot(8, 2, 8)
fig.set_figheight(30)
fig.set_figwidth(30)

ax = [[ax1, ax2], [ax3, ax4], [ax5, ax6], [ax7, ax8]]
clf = MyKMeans(n_clusters=0, init="k-means++")

for a in ax:
    X = make_blobs(500, 0)
    pred = clf.fit_predict(X)
    arg = np.argmax(pred[0])
    Dist = distance.cdist(X[arg], X[arg])
    a[0].imshow(Dist, cmap='blue')
    a[1].axis([-1.2, 1.2, -1.2, 1.2])
    a[1].scatter(X[:,0], X[:,1], c=pred[0])
    a[1].scatter(pred[4][:,0], pred[4][:,1], marker = '+', c='r')
```

