

# Домашнее задание

1. К алгоритму kNN, представленному на уроке, реализовать добавление весов для соседей по любому из показанных на уроке принципов.

In [1]:

```
1 import numpy as np
2 from sklearn import model_selection
3 from sklearn.datasets import load_iris
4 import matplotlib.pyplot as plt
5 from matplotlib.colors import ListedColormap
```

Загрузим один из "игрушечных" датасетов из sklearn.

In [2]:

```
1 X, y = load_iris(return_X_y=True)
2
3 # Для наглядности возьмем только первые два признака (всего в датасете их 4)
4 X = X[:, :2]
```

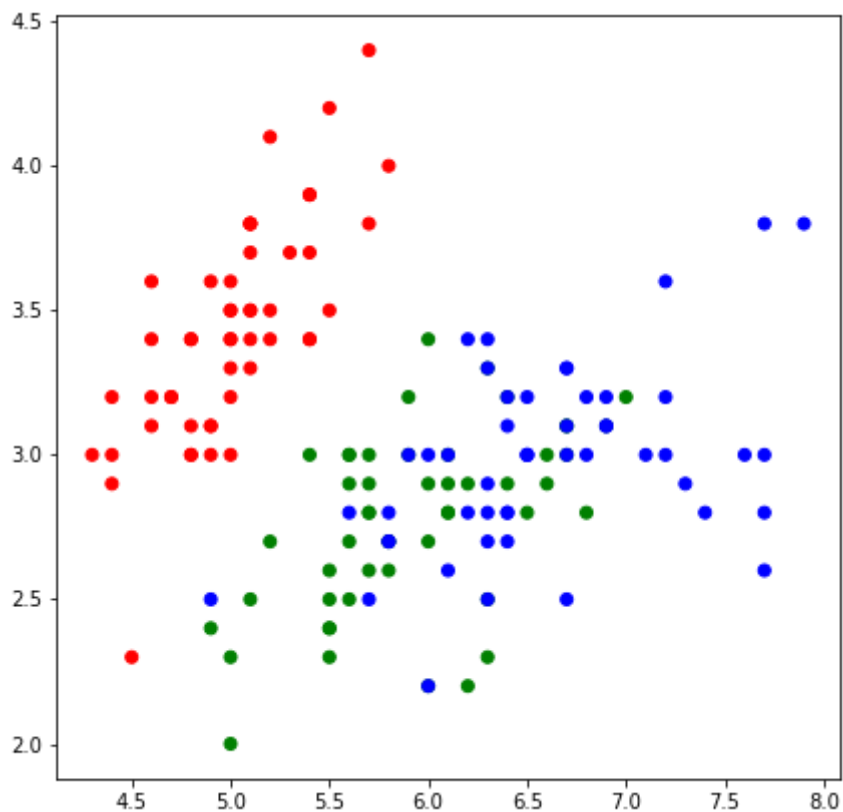
Разделим выборку на обучающую и тестовую

In [3]:

```
1 X_train, X_test, y_train, y_test = model_selection.train_test_split(X, y, test_size=0.2)
```

In [4]:

```
1 cmap = ListedColormap(['red', 'green', 'blue'])
2 plt.figure(figsize=(7, 7))
3 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap);
```



Используем евклидову метрику. Реализуем функцию для ее подсчета.

In [5]:

```
1 def e_metrics(x1, x2):
2
3     distance = 0
4     for i in range(len(x1)):
5         distance += np.square(x1[i] - x2[i])
6
7     return np.sqrt(distance)
```

Реализуем алгоритм поиска k ближайших соседей.

In [6]:

```
1 def knn(x_train, y_train, x_test, k):
2
3     answers = []
4     for x in x_test:
5         test_distances = []
6
7         for i in range(len(x_train)):
8
9             # расчет расстояния от классифицируемого объекта до
10            # объекта обучающей выборки
11            distance = e_metrics(x, x_train[i])
12
13            # Записываем в список значение расстояния и ответа на объекте обучающей вы
14            test_distances.append((distance, y_train[i]))
15
16
17            # создаем словарь со всеми возможными классами
18            classes = {class_item: 0 for class_item in set(y_train)}
19
20            # Сортируем список и среди первых k элементов подсчитаем частоту появления разн
21            for d in sorted(test_distances)[0:k]:
22                classes[d[1]] += 1
23            # print(classes)
24            # Записываем в список ответов наиболее часто встречающийся класс
25            answers.append(sorted(classes, key=classes.get)[-1])
26
27     return answers
```

Реализуем алгоритм поиска k ближайших соседей с весами

Добавим соседям веса, вычисленные по формуле:

$$w(d) = q^d, q \in (0, 1),$$

где d - расстояние

In [7]:

```
1 def knn_weight(x_train, y_train, x_test, k):
2
3     answers = []
4     for x in x_test:
5         test_distances = []
6
7         for i in range(len(x_train)):
8
9             # расчет расстояния от классифицируемого объекта до
10            # объекта обучающей выборки
11            distance = e_metrics(x, x_train[i])
12
13            # Записываем в список значение расстояния и ответа на объекте обучающей вы
14            test_distances.append((distance, y_train[i]))
15
16
17            # создаем словарь со всеми возможными классами
18            classes = {class_item: 0 for class_item in set(y_train)}
19
20            # Сортируем список и среди первых k элементов подсчитаем частоту появления разн
21            for d, i in zip(sorted(test_distances)[0:k], range(1, k+1)):
22                classes[d[1]] += 0.8**d[0]
23
24            # Записываем в список ответов наиболее часто встречающийся класс
25            answers.append(sorted(classes, key=classes.get)[-1])
26        # print(classes)
27
28    return answers
```

Напишем функцию для вычисления точности

In [8]:

```
1 def accuracy(pred, y):
2     return (sum(pred == y) / len(y))
```

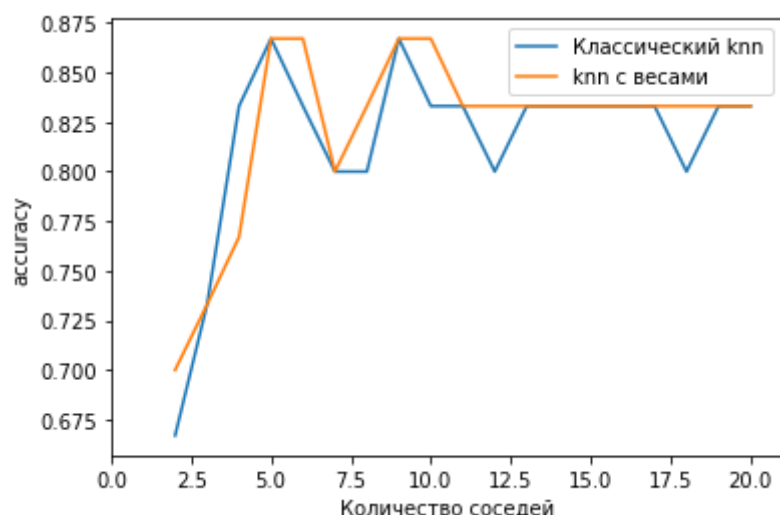
Проверим работу алгоритма при различных k

In [9]:

```
1 accuracies = []
2 for k in range(2, 21):
3     y_pred = knn(X_train, y_train, X_test, k)
4     y_pred_weight = knn_weight(X_train, y_train, X_test, k)
5     accuracies.append([k, float(f'{accuracy(y_pred, y_test):.3f}'), float(f'{accuracy(y
```

In [10]:

```
1 # Сравним ассурасу алгоритмов
2 plt.xlabel('Количество соседей')
3 plt.ylabel('accuracy')
4 plt.xlim(0, list(zip(*accuracies))[0][-1] + 1)
5 plt.plot(list(zip(*accuracies))[0], list(zip(*accuracies))[1], label='Классический knn')
6 plt.plot(list(zip(*accuracies))[0], list(zip(*accuracies))[2], label='knn с весами')
7 plt.legend(loc='best')
8 plt.show()
```



**Вывод:** Из графика видно, что алгоритм с весами работает лучше. Только при количестве соседей 3 классический knn работает лучше, во всех остальных случаях knn с весами работает или одинаково или лучше

2. (\*) Написать функцию подсчета метрики качества кластеризации как среднее квадратичное внутриклассовое расстояние и построить график ее зависимости от количества кластеров k (взять от 1 до 10) для выборки данных из этого урока (создать датасет, как в методичке).

Создадим датасет

In [11]:

```
1 from sklearn.datasets import make_blobs
2 import random
3 X, y = make_blobs(n_samples=100, random_state=1)
```

In [12]:

```
1 def kmeans(data, k, max_iterations, min_distance):
2     # Создадим словарь для классификации
3     classes = {i: [] for i in range(k)}
4
5     # инициализируем центроиды как первые k элементов датасета
6     centroids = [data[i] for i in range(k)]
7
8     for _ in range(max_iterations):
9         # классифицируем объекты по центроидам
10        old_classes = classes.copy()
11        classes = {i: [] for i in range(k)}
12        # print(f'На {_} итерации классы = {old_classes}')
13        for x in data:
14
15            # определим расстояния от объекта до каждого центроида
16            distances = [e_metrics(x, centroid) for centroid in centroids]
17            # отнесем объект к кластеру, до центроида которого наименьшее расстояние
18            classification = distances.index(min(distances))
19            classes[classification].append(x)
20            # print(f'class = {classes}')
21
22            # сохраним предыдущие центроиды в отдельный список для последующего сравнения с
23            old_centroids = centroids.copy()
24            # print(f'После цикла алгоритма класс = {classes}')
25            # пересчитаем центроиды как среднее по кластерам
26            for classification in classes:
27                centroids[classification] = np.average(classes[classification], axis=0)
28
29            # print(wss(centroids, classes))
30            # сравним величину смещения центроидов с минимальной
31            optimal = True
32            for centroid in range(len(centroids)):
33                if np.sum(abs((centroids[centroid] - old_centroids[centroid]) / old_centroids[centroid])) > min_distance:
34                    optimal = False
35
36            # если все смещения меньше минимального, останавливаем алгоритм
37            if optimal:
38                break
39
40        return old_centroids, classes
```

Напишем функцию для расчета среднего квадратичного внутриклассового расстояния по формуле

$$\sum_{k=1}^K \frac{1}{|k|} \sum_{i=1}^l [a(x_i) = k] \rho^2(x_i, c_k),$$

где  $|k|$  - количество элементов в кластере под номером  $k$ .

In [13]:

```
1 def average_square_class_distance(centroids, clusters):
2     a_sq_cla_dist = 0
3     for i in clusters:
4         distance = np.sum([e_metrics(x, centroids[i])**2 for x in clusters[i]]) / len(clusters[i])
5         a_sq_cla_dist += distance
6     return a_sq_cla_dist
```

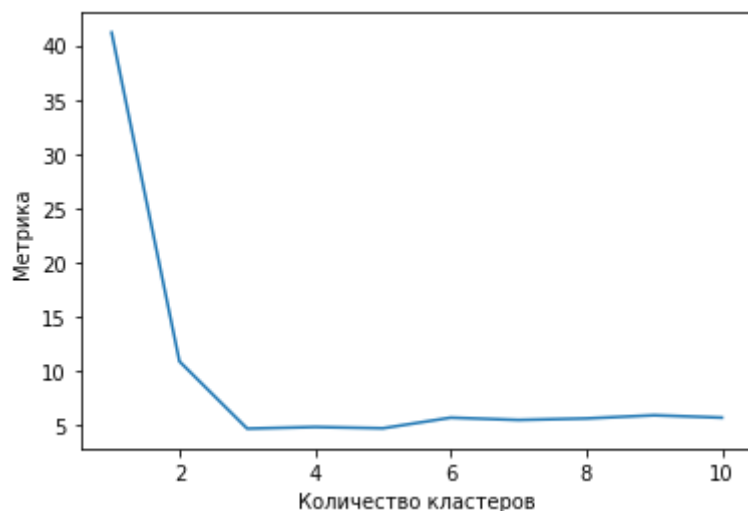
Рассчитаем средние квадратичные внутриклассовые расстояния для различного количества центроидов

In [14]:

```
1 a_sq_cla_dists = []
2 for k in range(1, 11):
3     centroids, clusters = kmeans(X, k, max_iterations=10, min_distance=1e-4)
4     a_sq_cla_dists.append(average_square_class_distance(centroids, clusters))
```

In [15]:

```
1 # построим график зависимости среднего квадратичного внутриклассового расстояния от кол
2 plt.xlabel('Количество кластеров')
3 plt.ylabel('Метрика')
4 plt.plot([i for i in range(1, 11)], a_sq_cla_dists);
```



Видим, что наименьшая метрика достигается при количестве центроидов 3

Сделаю проверку правильности расчетов. В библиотеке sklearn есть модель KMeans. У нее есть метод `inertia_`, который показывает значение целевой функции - суммы квадратов внутриклассовых расстояний, рассчитанной по формуле:

$$\sum_{k=1}^K \sum_{i=1}^l [a(x_i) = k] \rho^2(x_i, c_k)$$

Рассчитаем ее с помощью методов библиотеки sklearn для разного количества центроидов

In [16]:

```
1 from sklearn.cluster import KMeans
2 inercias = []
3 for k in range(1, 11):
4     model = KMeans(n_clusters=k, random_state=42)
5     model.fit(X)
6     inercias.append(model.inertia_)
```

In [17]:

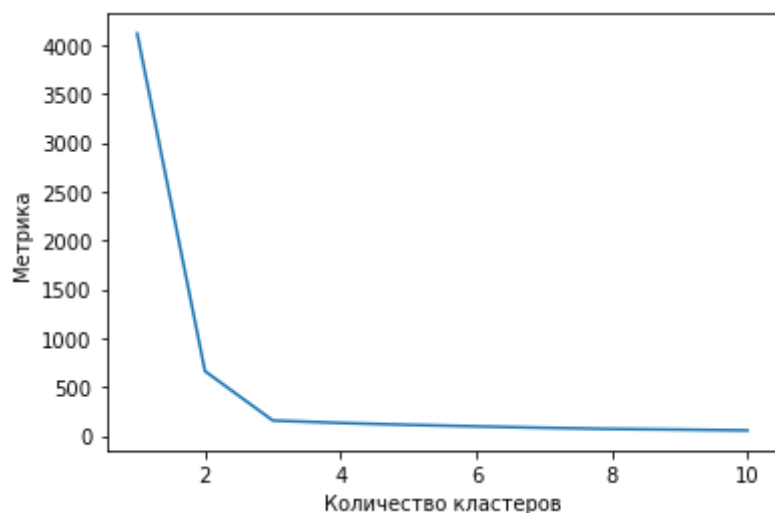
```
1 inercias
```

Out[17]:

```
[4118.153777704471,  
661.5698490972003,  
156.28289251170003,  
130.96121900774804,  
112.05653036061504,  
96.26315742705445,  
81.49760027212565,  
69.70669902233496,  
62.82959029286148,  
54.408866472636056]
```

In [18]:

```
1 # построим график зависимости среднего квадратичного внутриклассового расстояния от кол  
2 plt.xlabel('Количество кластеров')  
3 plt.ylabel('Метрика')  
4 plt.plot([i for i in range(1, 11)], inercias);
```



Изменим нашу функцию и посмотрим какие расчеты покажет она

In [19]:

```
1 def average_square_class_distance(centroids, clusters):  
2     a_sq_cla_dist = 0  
3     for i in clusters:  
4         distance = np.sum([e_metrics(x, centroids[i])**2 for x in clusters[i]])  
5         a_sq_cla_dist += distance  
6     return a_sq_cla_dist
```

In [20]:

```
1 a_sq_cla_dists = []  
2 for k in range(1, 11):  
3     centroids, clusters = kmeans(X, k, max_iterations=10, min_distance=1e-4)  
4     a_sq_cla_dists.append(average_square_class_distance(centroids, clusters))
```



In [21]:

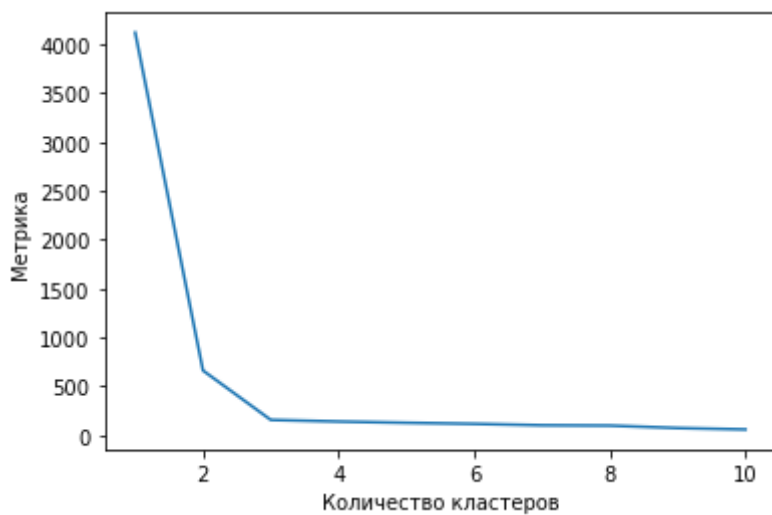
```
1 a_sq_cla_dists
```

Out[21]:

```
[4118.153777704471,  
661.5698490972001,  
156.28289251170003,  
139.3741146136838,  
126.57537841891627,  
115.8130108490437,  
101.33706078363487,  
98.00829163101346,  
72.88769594879396,  
58.46251834933362]
```

In [22]:

```
1 # построим график зависимости среднего квадратичного внутриклассового расстояния от кол  
2 plt.xlabel('Количество кластеров')  
3 plt.ylabel('Метрика')  
4 plt.plot([i for i in range(1, 11)], a_sq_cla_dists);
```



Видим, что цифры и графики очень похожи, значит наши расчеты были верные