

```
In [1]: import warnings
import pandas as pd
import numpy as np
import random
import seaborn as sns
import matplotlib.pyplot as plt

from os.path import join
from IPython import display
from sklearn.datasets import load_digits
from sklearn.datasets import make_blobs
from sklearn.decomposition import PCA
from sklearn.manifold import TSNE
from sklearn.metrics import silhouette_score # и другие метрики
from sklearn.cluster import KMeans # а также другие алгоритмы
```

```
In [2]: DATA_PATH = "data"
plt.rcParams["figure.figsize"] = 12, 9
sns.set_style("whitegrid")
warnings.filterwarnings("ignore")

SEED = 111
random.seed(SEED)
np.random.seed(SEED)
```

Задание 1. Реализация Kmeans

5 баллов

В данном задании вам предстоит дописать код класса `MyKMeans`. Мы на простом примере увидим, как подбираются центры кластеров и научимся их визуализировать.

Сгенерируем простой набор данных, 400 объектов и 2 признака (чтобы все быстро работало и можно было легко нарисовать):

```
In [3]: X, true_labels = make_blobs(400, 2, centers=[[0, 0], [-4, 0], [3.5, 3.5], [3.5, -2.0]])
```

Напишем функцию `visualize_clusters`, которая по данным и меткам кластеров будет рисовать их и разукрашивать:

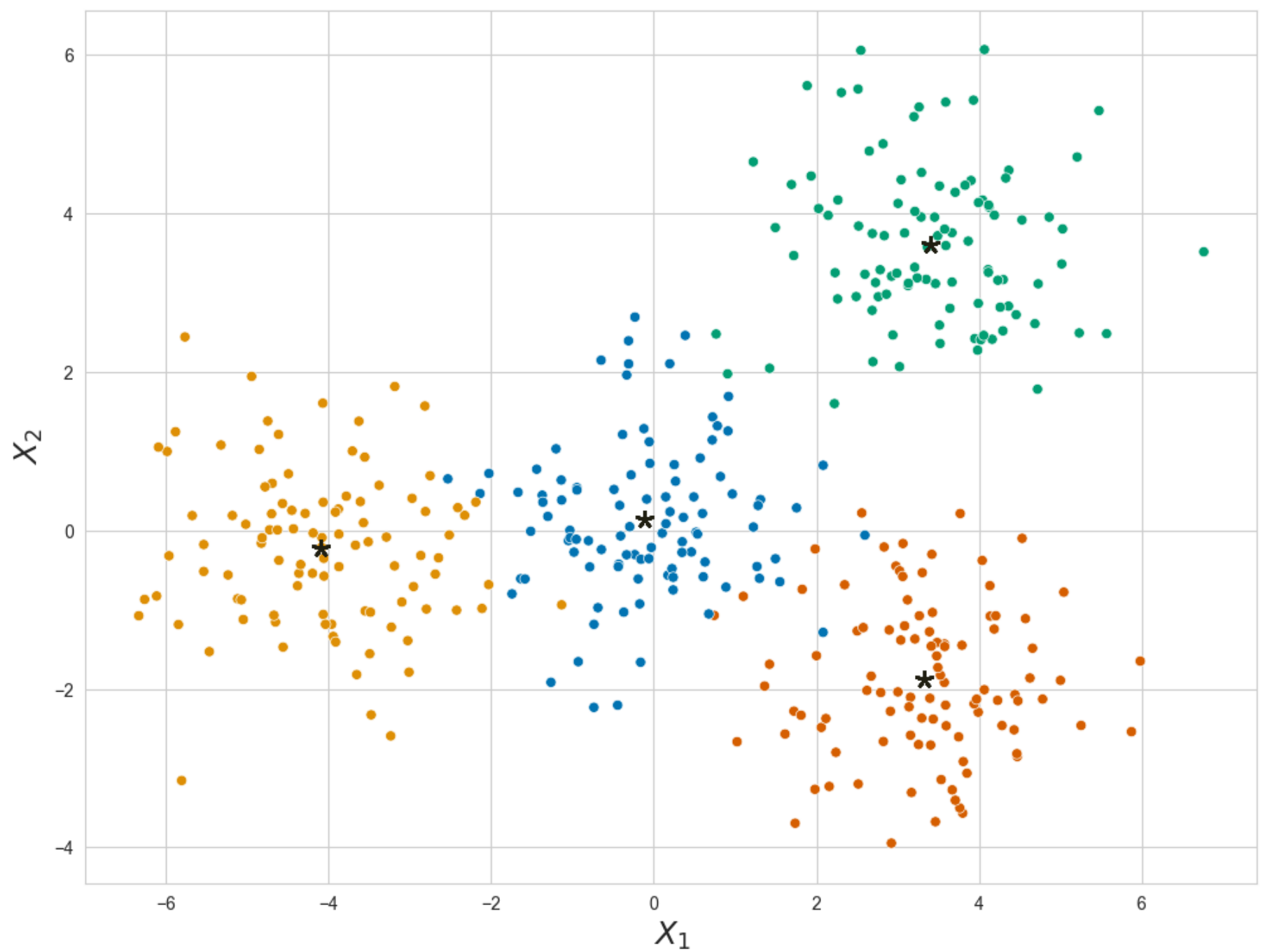
```
In [4]: def visualize_clusters(X, labels):
        """
        Функция для визуализации кластеров
        :param X: таблица объекты x признаки
        :param labels: np.array[n_samples] – номера кластеров
        """

        #     print(X[:, 0], X[:, 1])

        unique_labels = np.sort(np.unique(labels))
        sns.scatterplot(x=X[:, 0], y=X[:, 1], hue=labels,
                        palette="colorblind", legend=False,
                        hue_order=unique_labels)
        plt.xlabel("$X_1$", fontsize=18)
        plt.ylabel("$X_2$", fontsize=18)

        for label in labels:
            center = X[(labels == label)].mean(axis=0)
            plt.scatter(x=center[0], y=center[1], s=80, c="#201F12", marker=(5, 2))
```

```
In [5]: visualize_clusters(X, true_labels)
```



Напишем свой класс `MyKMeans`, который будет реализовывать алгоритм кластеризации К-средних. Напомним сам алгоритм:

1. Выбираем число кластеров (K)
2. Случайно инициализируем K точек (или выбираем из данных), это будут начальные центры наших кластеров
3. Далее для каждого объекта считаем расстояние до всех кластеров и присваиваем ему метку ближайшего
4. Далее для каждого кластера считаем "центр масс" (среднее значение для каждого признака по всем объектам кластера)
5. Этот "центр масс" становится новым центром кластера
6. Повторяем п.3, 4, 5 заданное число итераций или до сходимости

Во время предсказания алгоритм просто находит ближайший центроид (центр кластера) для тестового объекта и возвращает его номер.

Реализуйте методы:

- `_calculate_distance(X, centroid)` - вычисляет Евклидово расстояние от всех объектов в `X` до заданного центра кластера (`centroid`)
- `predict(X)` - для каждого элемента из `X` возвращает номер кластера, к которому относится данный элемент

```
In [6]: class MyKMeans:
def __init__(self, n_clusters, init="random", max_iter=300, visualize=False):
    """
    Конструктор класса MyKMeans
    :param n_clusters: число кластеров
    :param init: способ инициализации центров кластеров
        'random' – генерирует координаты случайно из нормального распределения
        'sample' – выбирает центроиды случайно из объектов выборки
    :param max_iter: заданное число итераций
        (мы не будем реализовывать другой критерий остановки)
    :param visualize: рисовать ли кластеры и их центроиды в процессе работы
        код будет работать сильно дольше, но красиво...
    """

    assert init in ["random", "sample"], f"Неизвестный метод инициализации {init}"
    self.n_clusters = n_clusters
    self.init = init
    self.max_iter = max_iter
    self.centroids = None
    self.visualize = visualize

def fit(self, X):
    """
    Подбирает оптимальные центры кластеров
    :param X: наши данные (n_samples, n_features)
    :return self: все как в sklearn
    """

    n_samples, n_features = X.shape

    # Инициализация центров кластеров
    if self.init == "random":
        centroids = np.random.randn(self.n_clusters, n_features)
    elif self.init == "sample":
        centroids_idx = np.random.choice(np.arange(n_samples),
                                         size=self.n_clusters,
                                         replace=False)

        centroids = X[centroids_idx]

    # Итеративно двигаем центры
    for _ in range(self.max_iter):
        # Посчитаем расстояния для всех объектов до каждого центроида
        dists = []
        for centroid in centroids:
            dists.append(self._calculate_distance(X, centroid))
        dists = np.concatenate(dists, axis=1)
        # Для каждого объекта найдем, к какому центроиду он ближе
        cluster_labels = np.argmin(dists, axis=1)

        # Пересчитаем центр масс для каждого кластера
        centroids = []
        for label in np.sort(np.unique(cluster_labels)):
            center = X[(cluster_labels == label)].mean(axis=0)
            centroids.append(center)

        # Отрисовываем точки, покрасим по меткам кластера, а также изобразим центроиды
        if self.visualize:
            visualize_clusters(X, cluster_labels)
            display.clear_output(wait=True)
            display.display(plt.gcf())
            plt.close()

    self.centroids = np.array(centroids)

    return self

def predict(self, X):
    """
    Для каждого X возвращает номер кластера, к которому он относится
    :param X: наши данные (n_samples, n_features)
    :return cluster_labels: метки кластеров
    """

    dists = []
    for centroid in self.centroids:
        dists.append(self._calculate_distance(X, centroid))

    dists = np.concatenate(dists, axis=1)

    cluster_labels = np.argmin(dists, axis=1)

    return cluster_labels

def _calculate_distance(self, X, centroid):
```