

МОСКОВСКИЙ ИНСТИТУТ ЭЛЕКТРОННОЙ ТЕХНИКИ
Институт системной и программной инженерии
и информационных технологий (Институт СПИНТех)

Лабораторный практикум по курсу
"Интеллектуальные системы"
(09/23 – 01/24)

Лабораторная работа 4
Кластеризация данных методом k -средних ¹.

На этом занятии компьютерного практикума предстоит изучить кластеризацию методом k -средних и применить данный метод для сжатия изображения. Прежде чем приступить собственно к программированию, настоятельно рекомендуется ознакомиться с материалом лекций, а также с дополнительными материалами, имеющими отношение к задаче кластеризации методом k -средних.

Файлы, включенные в задание:

lab_k_means.ipnb - ноутбук, реализующий пошаговое выполнение задания по разделу «Кластеризация методом k -средних»;

ex4data.mat – набор данных для метода k -средних;

bird_small.png – пример изображения;

displayData – отображает 2-мерные данные, сохраненные в матрице;

drawLine – рисует линию на существующем графике;

plotDataPoints – выводит в графическом виде центры для метода k -средних;

plotProgresskMeans – выводит в графическом виде каждый шаг метода k -средних во время расчета;

runkMeans – запуск метода k -средних.

* *findClosestCentroids* – находит ближайшие центры (для метода k -средних);

* *computeCentroids* – рассчитывает среднее значение центров (для метода k -средних)

* *kMeansInitCentroids* – задание центров (для метода k -средних)

Функции, отмеченные знаком *, следует написать самостоятельно.

В этой Лабораторной работе следует использовать ноутбук *lab_k_means*. В данном ноутбуке подготовлены обращения к исходным данным. Далее производится вызов функций, написанных Вами, и отображаются результаты вычислений. Необходимо дописать функции в файлах по инструкциям упражнения.

1 Метод k -средних

В этом упражнении предстоит изучить как, собственно, сам метод k -средних, так и применить его для сжатия изображений. В первой части Вы применяете метод k -средних для обработки 2-мерного набора данных, что способствует лучшему пониманию алгоритма. Далее,

¹ <https://github.com/alexeykuznetsov/k-means>

метод k -средних применяется для сжатия изображения, посредством уменьшения количества цветов и сохранения только тех цветов, которые встречаются в данном изображении чаще всего.

Указание: Вы располагаете также написанными ранее (Лабораторные работы 1-2) программами, которые здесь можно использовать.

1.1 Использование метода k -средних

Метод k -средних автоматически объединяет в группы схожие данные. Например, пусть задан набор данных для обучения $\{x^{(1)}, \dots, x^{(m)}\}$, где $x^{(i)} \in \mathbb{R}^n$, и Вы хотите его объединить в несколько сплоченных «кластеров». Принцип работы метода k -средних, заключается в том, что он является итерационной процедурой, в которой начальные центры задаются некоторым образом, часто случайным, а далее повторно уточняются. Уточнение сводится к перегруппировке примеров (точек) к ближайшим центрам и последующему пересчету местоположения центров.

Указание: код для метода k -средних выглядит следующим образом:

```
% Задание начального положение центров
centroids = kMeansInitCentroids(X, K);
for iter = 1:iterations
% Присвоение каждой точки ближайшему центру
% idx(i) соответствует  $\hat{c}^{(i)}$  - индексу центра  $i$ -го примера
idx = findClosestCentroids(X, centroids);
% Шаг перемещения центра: Вычисление среднего местоположения
centroids = computeMeans(X, idx, K);
end
```

Внутренний цикл осуществляет две операции: 1) группировка каждого примера $x^{(i)}$ к ближайшему центру кластеризации и 2) перерасчет среднего значения каждого центра, используя примеры, которые отнесены к его группе. Метод k -средних всегда сходится и находит конечное среднее для центров. Обратите внимание, что решение не всегда будет оптимальным и зависит от начальных значений центров. Соответственно, метод k -средних обычно запускают несколько раз с различными начальными параметрами. Один из способов выбрать между различными решениями - это найти то решение, для которого функция стоимости (искажение) имеет наименьшее значение.

В следующей части требуется выполнить два этапа метода k -средних по отдельности.

1.1.1 Нахождение ближайших центров

На первом этапе алгоритм метода k -средних, группирует каждый пример $x^{(i)}$ к его ближайшему центру для выбранного расположения центров. Для каждого примера i :

$$c^{(i)} := j \text{ минимизирует } \|x^{(i)} - \mu_j\|^2,$$

где $c^{(i)}$ - индекс ближайшего к $x^{(i)}$ центра, а μ_j - положение (численное значение) j -го центра. Обратите внимание на то, что в коде $c^{(i)}$ обозначен как $idx(i)$.

Задание: Ваша задача завершить код в *findClosestCentroids.m*. Это функция, использует матрицу X и расположение центров для вывода одномерного массива idx . Массив idx содержит индекс (значение принадлежащее $\{1, \dots, K\}$ где K – число центров) ближайшего центра к каждому примеру. Это может быть реализовано циклом для каждого примера для каждого центра. После завершения кода в *findClosestCentroids*, ноутбук запустит код и Вы увидите результат $\begin{bmatrix} 1 & 3 & 2 \end{bmatrix}$, соответствующий ближайшим центрам для первых 3-х примеров. Вы также можете увидеть результат $\begin{bmatrix} 0 & 2 & 1 \end{bmatrix}$.

1.1.2 Расчет среднего расположения центров

После распределения всех точек к ближайшему центру, на втором этапе алгоритм находит новые положения центров: среднее положение точек которые к нему принадлежат к конкретному центру. Для каждого центра k имеем

$$\mu_k := \frac{1}{|C_k|} \sum_{i \in C_k} x^{(i)},$$

где C_k набор примеров, которые относятся к центру k . Например, если два примера $x^{(3)}$ и $x^{(5)}$ относятся к центру $k = 2$, тогда следует обновить

$$\mu_2 = 1/2(x^{(3)} + x^{(5)}).$$

Задание: Завершить код в программе *computeCentroids*. Это можно сделать, например, используя цикл для каждого центра или цикл для каждого примера, но если расчеты будут производиться в векторном виде, программа будет работать значительно быстрее. После завершения кода в *computeCentroids*, ноутбук запустит вашу работу и отобразит центры после первого этапа метода К-среднего.

1.2 Метод k -средних с набором данных из примера

После завершения работы двух функций (*findClosestCentroids* и *computeCentroids*) ноутбук применит метод k -средних с 2-мерным набором данных. Функции, которые Вы написали, вызываются из *runKMeans*. Для лучшего понимания материала рекомендуется ознакомиться с этим кодом. Обратите внимание на то, что функции вызываются в цикле. Следующий шаг программы - это последовательное отображение каждой итерации алгоритма. Нажмите клавишу «Enter» несколько раз, чтобы увидеть изменение расположения центров. Окончательный вариант должен выглядеть так, как представлено на рис. 1.

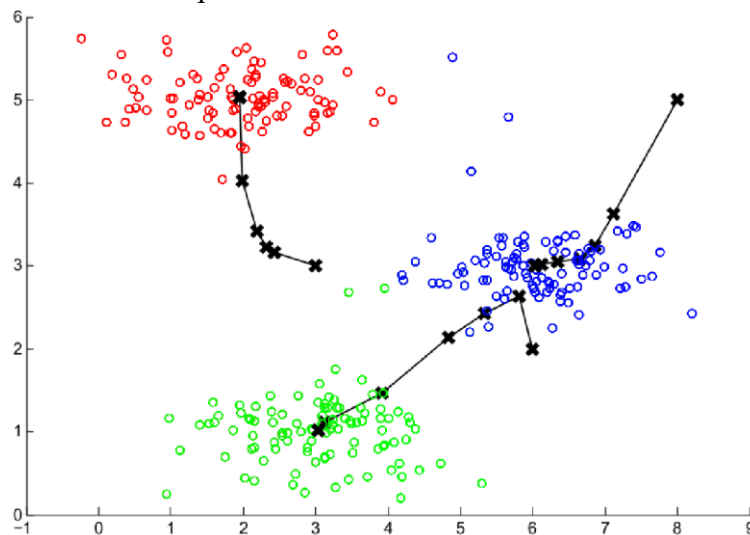


Рис. 1. Ожидаемый результат (10 итераций)

1.3 Случайный запуск

Начальное расположение центров для предыдущего набора данных было выбрано таким образом, чтобы Вы смогли увидеть такую же картинку, как на рис. 1. В действительности, лучшим подходом будет случайный выбор начального расположения центров. В этой части упражнения Вам следует завершить код в *kMeansInitCentroids*, добавив следующие строчки:

```
% Случайное задание центров из примеров
% Случайно перегруппировать индексы примеров
centroids = np.zeros((K, X.shape[1]))
randidx = np.random.permutation(X.shape[0])
centroids = X[randidx[:K], :]
```

Здесь случайным образом перемешиваются индексы примеров (используя *randperm*) и выбираются первые K примеров в качестве центров.

1.4 Сжатие изображений, используя метод k -средних



Рис. 2. Исходное изображение 128 x 128

Используйте метод k -средних для сжатия изображения. В обычном, 24-х битовом цветовом представлении каждый пиксель соответствует трем 8-ми битным целым числам (в интервале от 0 до 255) - красный, зеленый, синий. Данную цветовую модель обычно обозначают как *RGB*. Приведенная картинка птицы содержит тысячи цветов и в этой части упражнения предстоит уменьшить их количество до 16, что потребует только 4 бита. Используйте метод k -средних для выбора 16-ти цветов, которые будут отображать рассматриваемое изображение. Каждый пиксель представляется как набор данных и метод k -средних находит 16 цветов, которые лучше всего группируют пиксели в 3-х мерном *RGB*-пространстве. После нахождения центров Вы замените цвет каждого пикселя исходного изображения на один из 16-ти цветов.

1.4.1 Метод k -средних и пиксели (“элизы”)

После считывания изображения создается 3-х мерную матрицу A , в которой первые два элемента обозначают местоположение пикселя, а третий обозначает интенсивность красного, зеленого или синего цветов. Например, $A(50, 33, 3)$ обозначает пиксель с координатами 50, 30 и интенсивностью синего цвета 3. Код в программе *ex4.m* загружает изображение и переделывает ее в матрицу размером $m \times 3$ (где $m = 16384 = 128 \times 128$), и запускает метод k -средних.

Задание: После определения первых $K = 16$ цветов для представления изображения, требуется распределить пиксели в зависимости от ближайшего к ним центра, используя функцию *findClosestCentroids*. Это преобразует исходное изображение, используя ближайшие для каждого пикселя центры. Обратите внимание на то, что количество битов, необходимых для хранения изображения значительно уменьшилось. Исходное изображение требовало 24 бита для каждого из

128×128 пикселей, занимая в итоге $128 \times 128 \times 24 = 393.216$ битов. Новое отображение требует «словаря» из 16 цветов, занимающих 24 бита каждый и 4 бита для каждого пикселя, получая $16 \times 24 + 128 \times 128 \times 4 = 65.920$ битов, что соответствует сжатию исходного изображения примерно в 6 раз.

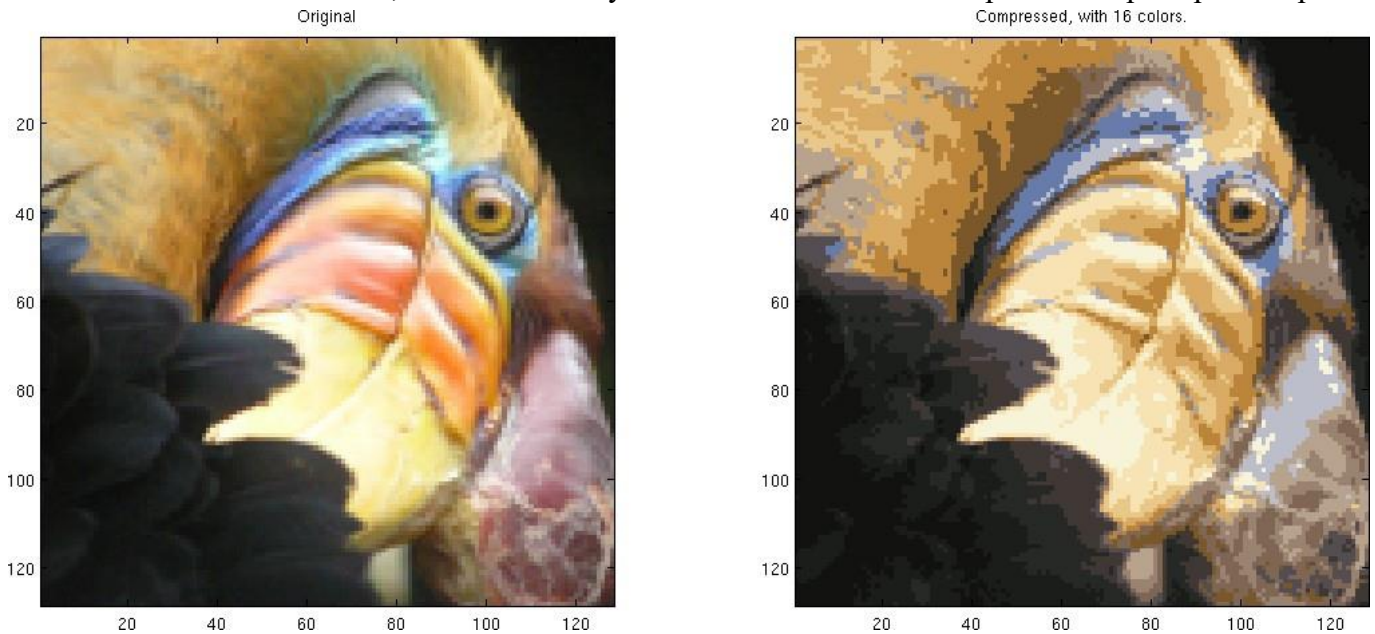


Рис. 3: Исходное и преобразованное изображение (метод k -средних).

Результат сжатия с помощью метода k -средних можно увидеть на рис. 3. Каждый пиксель можно заменить на среднее значение центра, к которому он относится. Несмотря на то, что изображение сохраняет большинство своих характеристик, последствия сжатия весьма заметны.

1.5 Сжатие собственного изображения

В этом упражнении измените предыдущую программу, применив её для Вашего собственного изображения. Обратите внимание на то, что если изображение имеет большой размер, то для обработки может потребоваться достаточно долгое время. Рекомендуется вначале его уменьшить. Вы также можете использовать различное количество центров и посмотреть на разницу при сжатии.