## Lab 08

## Basic Image Processing Fall 2019

## **K-means Algorithm**

- *K-means* algorithm (MacQueen'67): a heuristic method
  - Each cluster is represented by the centre of the cluster and the algorithm converges to stable centroids of clusters.
  - K-means algorithm is the simplest partitioning method for clustering analysis and widely used in data mining applications.
  - Each sample will belong to the cluster with the nearest mean.
- The objective is to minimize the within-cluster sum of squares:

$$\underset{S}{\operatorname{argmin}} \sum_{i=1}^{K} \sum_{x \in S_{i}} d^{2}(x, \mu_{i}), \qquad d^{2}(x, \mu_{i}) = ||x - \mu_{i}||^{2} = \sum_{n=1}^{N} (x_{n} - \mu_{in})^{2}$$

• where  $x \in \mathbb{R}^N$  are the data samples,  $\mu_i$  is the mean (prototype) of the points in the cluster  $S_i$  ( $i = 1 \dots K$ ).

## **K-means Algorithm**

- Given the cluster number *K*, the *K*-means algorithm is carried out in three steps after initialization:
  - **1)** Initialisation: set the *K* cluster seed points (randomly)
  - 2) Assignment step: Assign each object to the cluster of the nearest seed point measured with a specific distance metric
  - Update step: Compute new seed points as the centroids of the clusters of the current partition (the centroid is the centre, i.e., *mean point*, of the cluster)

$$\mu_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$

4) Go back to Step 1), stop when no more new assignment (i.e., membership in each cluster no longer changes)

## Spaces & dimensions (in pure math)

It is important to understand all the spaces and their dimensionalities in this task.

Consider an  $S \underline{space}$ :  $S \subset \mathbb{R}^n$ 

Every element of S is a <u>vector</u> with n coordinates:

$$\mathbf{x}_a = [x_a^1, x_a^2, \dots, x_a^n] \in S$$

The  $S_i$  subsets of S are the <u>clusters</u>:  $S_i \subset S$   $\bigcup_{i=1}^k S_i = S$  and  $S_i \cap S_j = \emptyset$   $\forall i \neq j$ These are the 'k'-s in the term and where i = 1, 2, ..., k These are the 'k'-s in the term 'k-means'. They tell you the number of clusters. Also, there are  $\mu$  vectors representing the mean values aka the <u>centroids</u> of every cluster:  $\mu_i = [\mu_i^1, \mu_i^2, ..., \mu_i^n] \in S$  i = 1, 2, ..., k



## Spaces & dimensions (in MATLAB)

Let us translate the terms of the previous slide into MATLAB.

Consider an S <u>space</u>, represented by a *matrix*.

Every *row* of **S** is a <u>vector</u> with *n* items:

 $S(1,:) = x_1 = [x_{11} x_{12} \dots x_{1n}]$ 

The S\_i subsets of S are the <u>clusters</u>. Their representations are stored in a look-up-table (LUT). The index represents the index, the value represents the cluster # of a row vector x\_a of S.

Also, the  $\mu$  mean vectors are stored in a matrix similar to S, denoted by M. Its elements are

M(j,:) = mu\_j = [mu\_j1 ... mu\_jn]



## So what is stored in the LUT?



The LUT is a vector. It has as many elements as the number of vectors in the space *S*.

Every element of the LUT has an *index* and a *value*.

The value at position *j* tells us which cluster does vector  $\mathbf{x}_j$  belong to. Now please

#### download the 'Lab 08' code package

from the

submission system

#### Implement the function step1\_initialization in which:

• The function has 2 inputs and 2 outputs:

Inputs:

- S set of points to be clustered
- k number of clusters

Outputs:

- LUT the assignment vector
- M the matrix of centroids

The function should initialize LUT and M as described on the next slide!

After implementation, test your function with the **script** test1\_initialization!

## Step 1: Initialization

- initialize the LUT (as an  $1 \times m$  vector, filled with zeros),
- initialize the M matrix (as a  $k \times n$  matrix, filled with zeros),
- choose k-many vectors as the initial cluster center points and store them in M.

Now (to be able to reproduce the results) **we are NOT using random initialization** 

but an equidistant distribution!

 $\begin{bmatrix} \cdot \end{bmatrix} \text{ means floor ()} \\ \text{We choose every } \begin{bmatrix} \frac{m}{k} \end{bmatrix} \text{-th element as an initial center.} \\ \text{E.g: } m = 16 \\ k = 5 \\ \Rightarrow \\ \begin{bmatrix} \frac{16}{5} \end{bmatrix} = 3 \\ \text{The first and then every 3^{rd} element of S} \\ \text{will be selected to be an element of M.} \\ \end{bmatrix}$ 



#### Implement the function step2\_assignment in which:

• The function has 4 inputs and 1 output:

Inputs:

- $\circ$  S set of points to be clustered  $\circ$  LUT the assignment vector
- k number of clusters M the matrix of centroids

Outputs:

• LUT the <u>updated</u> assignment vector

The function should update LUT as described on the next slide!

After implementation, test your function with the **script** test2\_assignment!

#### Step 2: Assignment

In this step:

For every x\_i vector in S (i=1..m) For every mu\_j vector in M (j=1..k) Calculate the distance between x\_i and mu\_j:

$$d_{ij} = d^2(\mathbf{x}_i, \mu_j) = \|\mathbf{x}_i - \mu_j\|^2 = \sum_{p=1}^n (x_i^p - \mu_j^p)^2$$

From the calculated **d\_ij** distances choose the smallest one, and store the index of the minimum in the LUT at position **i** :

$$LUT_i = \arg\min_{j=1..k} d_{ij}$$

#### Implement the function step3\_update in which:

• The function has 4 inputs and 1 output:

Inputs:

- $\circ$  S set of points to be clustered  $\circ$  LUT the assignment vector
- k number of clusters M the matrix of centroids

Outputs:

• M the <u>updated</u> matrix of centroids

The function should update M as described on the next slide!

After implementation, test your function with the **script** test3\_update!

## Step 3: Update

In this step:

#### For every mu\_j vector in M (j=1..k) Select every x vector of S that is assigned to the j-th cluster:

**MATLAB hint:** You can index a vector logically! If the LUT is a vector and you write LUT == 1 then this expression will return a logical vector: 1 if the element == 1, 0 otherwise.

If A = [1 2 3 1 1 2 1] then A == 1 returns [1 0 0 1 1 0 1]

The other trick is that if you index a vector or matrix with a logical vector, the result will be the set of those elements that has the same indices where the logical vector contained 1-s.

 If B = [1 2 3 4 5 6 7 then B(:, [1 0 0 1 1 0 1]) returns [1 4 5 7 .

 2 0 2 0 1 0 7]

Update mu\_j: the new value is the mean of the vectors of this cluster:

$$\mu_j^{(t+1)} = \frac{1}{|S_j|} \sum_{\mathbf{x}_i \in S_j} \mathbf{x}_i$$

#### Implement the function mykmeans in which:

• The function has 2 inputs and 2 outputs:

Inputs:

- S set of points to be clustered
- k number of clusters

Outputs:

- LUT the final assignment vector
- M the <u>final</u> matrix of centroids

The function should realize the iterative procedure described on the next slide. Please print the number of iterations after the execution of the iterative procedure.

After implementation, test your function with the **script** test4\_mykmeans!

#### Pseudo-code of the k-means algorithm

# function mykmeans(S, k) Initialization step while not converged and number of iterations is less than 100 Assignment step Update step

The algorithm *converged* if in the update step the sum of the distances between the old and new cluster center points is less than a threshold:

$$\sum_{j=1}^{k} \left\| \mu_{j}^{(t+1)} - \mu_{j}^{(t)} \right\|^{2} < \varepsilon \qquad \varepsilon = 0.02$$

#### Exercise 4 – result



#### Let's work with image data

An RGB color image is represented by a 3D matrix. It has *h* rows, *w* columns and 3 layers along the 3rd dimension. These layers are the R, G and B color layers.

To be able to cluster an image with the the k-means function we *somehow* has to transform the 3D matrix into a 2D array of row vectors.



#### How to create a segmented image?

![](_page_17_Figure_1.jpeg)

#### Let's work with image data

After the segmentation is complete we should rearrange the resulted A matrix into a matrix which has the sizes of the original image. This operation is the reverse of the transformation seen on the previous slide.

![](_page_18_Figure_2.jpeg)

#### Implement the function image\_segmenter in which:

• The function has 2 inputs and 1 output:

Inputs:

- I the RGB image to be clustered (segmented) in double format
- k number of clusters

Outputs:

• **OUT** the **final** segmented image

The function should convert the image to an S representation, call the function mykmeans on this array, and using the returned LUT and M create the out image.

After implementation, test your function with the **script** test5\_segmenter!

#### Exercise 5 – result

Input image

![](_page_20_Picture_2.jpeg)

RGB segmented image (k=2)

# THE END