

Assignment 5 [OPTIONAL] Dealing with Data in an Unsupervised Fashion

Homeworks Guidelines and Policies

- **What you must hand in.** It is expected that the students submit an assignment report (.pdf) as well as – if necessary – required source codes (.m or .py) into an archive file named according to the following template: HW4_XXXXXXX.zip where XXXXXXXX must be replaced by their student ID.
 - **Pay attention to problem types.** Some problems are required to be solved *by hand* (shown by the ✍ icon), and some need to be implemented (shown by the 🚀 icon). Please don't use implementation tools when it is asked to solve the problem by hand, otherwise you'll be penalized and lose some points.
 - **Don't bother typing!** You are free to solve by-hand problems on a paper and include picture of them in your report. Here, cleanness and readability are of high importance. Images should also have appropriate quality.
 - **Reports are critical.** Your work will be evaluated mostly by the quality of your report. Don't forget to explain what you have done, and provide enough discussions when it's needed.
 - **Appearance matters!** In each homework, 5 points (out of a possible 100) belongs to compactness, expressiveness and neatness of your report and codes.
 - **Python is also allowable.** By default, we assume you implement your codes in MATLAB. If you're using Python, you have to use equivalent functions when it is asked to use specific MATLAB functions.
 - **Be neat and tidy!** Your codes must be separated for each question, and for each part. For example, you have to create a separate .m file for part b. of question 3. Please name it like p3b.m.
 - **Use bonus points to improve your score.** Problems with bonus points are marked by the ★ icon. These problems usually include uncovered related topics or those that are only mentioned briefly in the class.
 - **Moodle access is essential.** Make sure you have access to Moodle because that's where all assignments as well as course announcements are posted on. Homework submissions are also done through Moodle.
-
- **Assignment Deadline.** Please submit your work **before TBD**.
 - **Delay policy.** During the semester, students are given 7 free late days which they can use them in their own ways. Afterwards there will be a 25% penalty for every late day, and no more than three late days will be accepted.
 - **Collaboration policy.** We encourage students to work together, share their findings and utilize all the resources available. However you are not allowed to share codes/answers or use works from the past semesters. Violators will receive a zero for that particular problem.
 - **Any questions?** If there is any question, please don't hesitate to contact us through the following email addresses: ali.the.special@gmail.com and dalirani@aut.ac.ir. You may also find us in pattern recognition and image processing lab, 3rd floor, CEIT building.

1. Understanding the Behavior of Clustering Techniques

(15 Pts.)



Keywords: *Unsupervised Learning, Clustering Problem, K-Means Clustering*

Another type of machine learning algorithms lie under the concept of **Unsupervised Learning**. These methods make inferences from data using only inputs without referring to known or labelled outputs. **Clustering** – known as an unsupervised method – is the attempt of assigning objects to different groups, or **Clusters**, so that those in the same group are more similar to each other than those in other groups. One of the most popular clustering algorithm is **K-Means**. It keeps k **Centroids** that it uses to define clusters. In K-Means, a point is considered to be in a certain cluster if it is closer to that cluster's centroid than any other centroid.

This problem consists of several parts which aim to evaluate your basic understanding of clustering, mainly K-Means method.

First, you are given five different sets of 2-D points in Figure 1, and you are asked to provide a sketch of how K-Means would split them into clusters considering the given number of clusters. You must also indicate approximately where the final centroids would be.

- a. $K = 2$
- b. $K = 3$
- c. $K = 2$
- d. $K = 3$
- e. $K = 2$

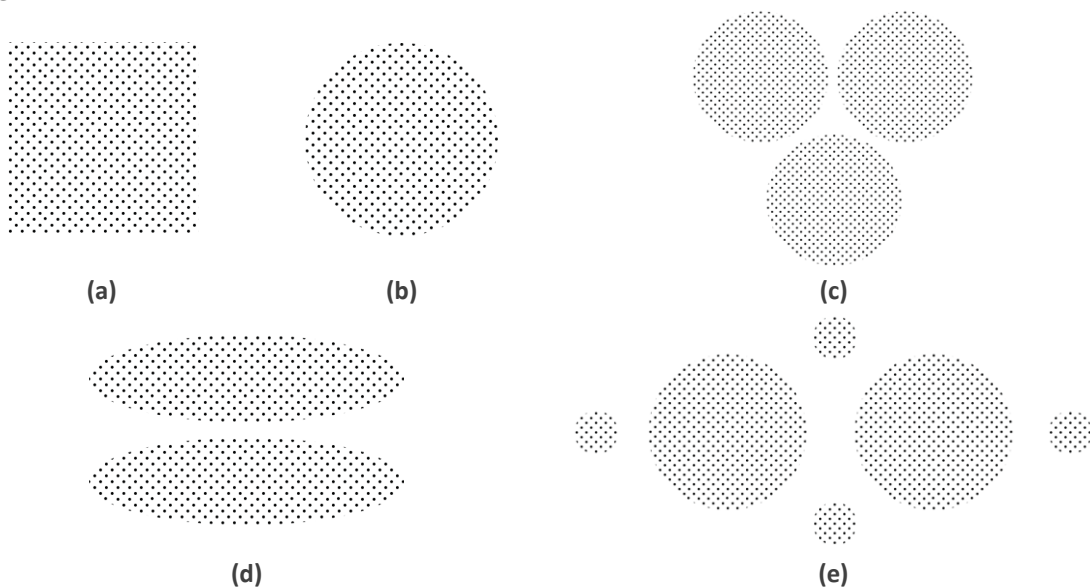


Figure 1 Sets of 2-D points provided for the first part of the problem 1

Note 1: Assume that the squared error objective function is being used.

Note 2: If there is more than one possible solution, then please specify for each solution whether it is a global or local minimum.

Note 3: Images in Figure 1 are given to you in “P1” directory attached to this assignment. You can use them in your report.

Next, consider the diagrams in Figure 2.

- f. In which one of the two diagrams do the classic clustering techniques, like single linkage, find the patterns represented by door and windows?
- g. Specify the limitations that clustering has in detecting the patterns formed by points.

Now, imagine there are five datasets, noted as a, b, c, d and e. The datasets are clustered using two different methods which one of them is K-Means. The distance measure used here is the Euclidean distance. Results are shown in Figure 3.

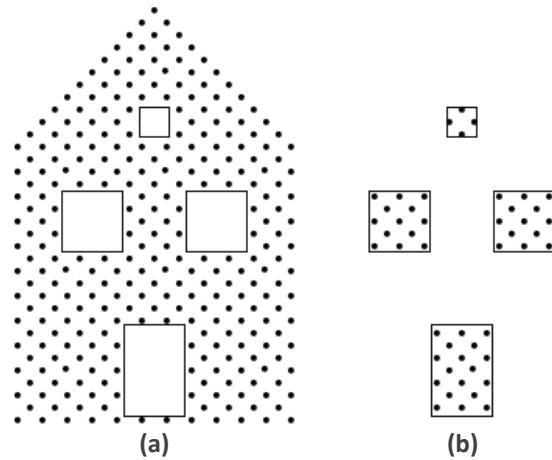


Figure 2 The goal of this clustering problem is to distinguish the main parts, i.e. door and windows

- h. Considering dataset (a) of the Figure 3, determine which result (1 or 2) is more likely to be generated by K-Means method.
- i. Repeat the part h for dataset (b).
- j. Repeat the part h for dataset (c).
- k. Repeat the part h for dataset (d).
- l. Repeat the part h for dataset (e).

Hint: You must check out the state when K-Means converges. Centres for each cluster are shown by red circles.

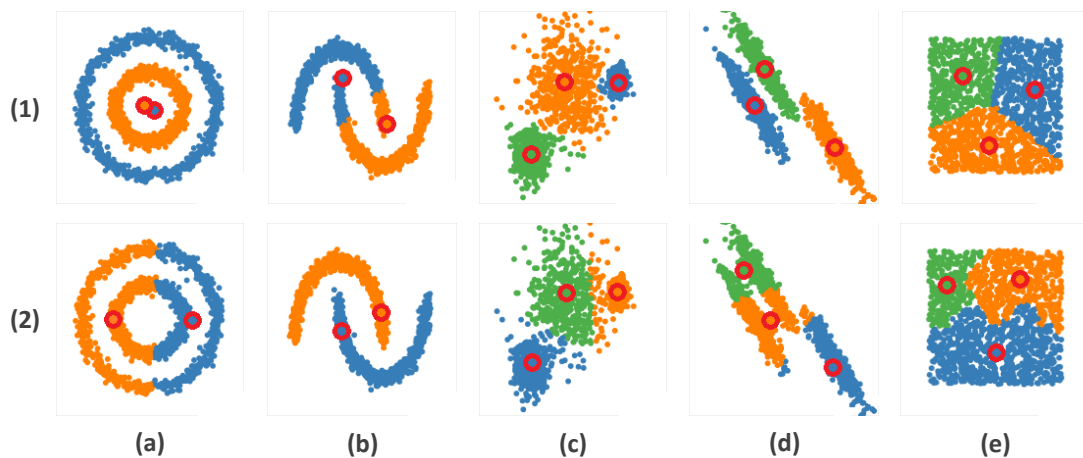


Figure 3 The clustering results of five datasets, each with different data distributions. Final centres are shown by red circles.

Finally, consider the points in Figure 4.

- m. Determine all well-separated clusters in the given set of points.

Note: You may use the equivalent image of Figure 4 placed in "P1" directory.

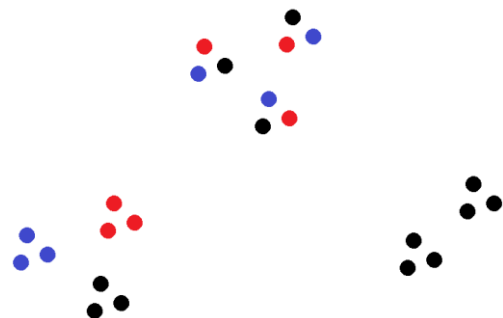


Figure 4 Clustering these points may yield to different results depending on the method and settings considered

2. K-Means and Beyond: Solving Clustering Problems by Hand

(10 Pts.)



Keywords: *Unsupervised Learning, Clustering Problem, K-Means, K-Median, K-Medoids, Hierarchical Clustering, Single-Linkage, Complete Linkage, Dendrogram*

K-Means clustering method has several variations, each structurally similar to the main algorithm but with slight differences. **K-Medians**, as an example, calculates the median in order to determine the centroids. However, there is another group of clustering strategies called **Hierarchical Clustering**, which also attempt to partition similar objects into same groups, but by building a hierarchy of clusters, known as **Dendrogram**.

In this part, you are to solve some basic clustering problems. First, consider the following toy dataset consisting of six points:

	x	y
1	1.00	3.00
2	2.00	1.00
3	2.00	5.00
4	4.00	1.00
5	4.00	5.00
6	5.00	3.00

- Assuming points 1 and 6 as the initial centres of clusters A and B respectively, use K-Means method to determine the final clustering result.
- Now consider points 3 and 4 as the initial centres of clusters A and B respectively, and repeat the previous part.
- Compare these two clustering results based on SSE measure and determine which one is better.

Now for hierarchical clustering, assume the following set of points:

	x	y
1	0.2	0.19
2	0.23	0.68
3	0.33	0.61
4	0.56	0.05
5	0.9	0.2
6	0.9	0.7
7	0.9	0.95

- Draw a sketch of the hierarchical clustering tree (dendrogram) we would obtain for single linkage method, considering Euclidean distance measure.
- Repeat the previous part for complete linkage method.

3. K-Means Maneuver in Image Processing Territory

(15 Pts.)



Keywords: *K-Means Clustering, Vector Quantisation, Color Extraction, Image Segmentation, Image Compression*

Despite its simplicity, **K-Means** can be applied to various machine learning tasks. From document classification and data analysis to fraud detection and collaborative filtering, this algorithm still challenges even newly introduced methods in different applications, which many of them are known to be state-of-the-art.

The goal of this problem is to get you more familiar with some of the thing you can do with K-Means in the area of image processing. Given below are three different, but structurally similar image processing tasks and you are required to propose a method to solve them using K-Means.

- a. **Color Extraction** is the process of identifying and extracting key colors in images. It gives a better visual understanding of images while providing significant features for other computer vision tasks.

Load the image “tiny_trump_1.jpg”. Use K-Means to extract its 3, 5, 7 and 9 main colors. Display these colors properly in separate square shapes.



(a)

Proportional palette	Hex color	Area	Closest color name
	#211b1b	31.7 %	Bokara Grey (Grey)
	#c6b293	23.5 %	Sour Dough (Brown)
	#784028	14.1 %	Copper Canyon (Brown)
	#d1c6c1	9.9 %	Swiss Coffee (Grey)
	#6a5a4d	8.9 %	Domino (Brown)
	#464c2b	4.4 %	Waiouru (Green)
	#4e301f	3.2 %	Indian Tan (Brown)
	#c3a36a	2.5 %	Putty (Yellow)
	#b57a67	1.2 %	Toast (Brown)
	#847e8b	0.7 %	Topaz (Grey)

(b)

Figure 5 An example image with its extracted main colors, called “palette”, sorted by area they occupy in the input image. The result is obtained using an online tool called TinEye [here](#) (a) Original image (b) Color extraction results.

- b. **Image Segmentation** is a common technique in image processing in which the goal is to divide an image into multiple parts or regions, often based on the characteristics of the pixels in the image.

Load the image “tiny_trump_2.jpg”. Use K-Means to divide the given image into 3, 5, 7 and 9 partitions.



(a)



(b)

Figure 6 Image segmentation applied to an example input image (a) Original image (b) The result of image segmentation

- c. **Image Compression** refers to techniques used for minimising the size of an image using the image data which are repeated in the image.

Load the image "tiny_trump_3.jpg". Use K-Means to reduce the size of the input image to %50, %75, %90 and %97 of the original image size (in KB).

Hint: You must find appropriate values for parameter k .

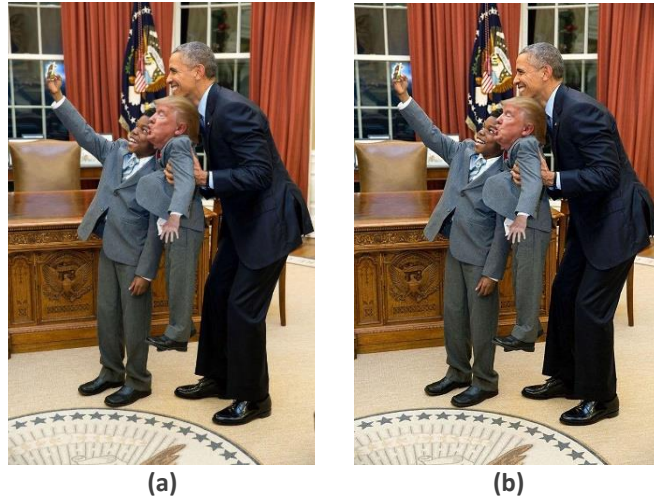


Figure 7 A compression technique has reduced the image size from 211KB to 83KB. Note that the difference is not properly noticeable here
(a) Original image (b) Compressed image.

4. Some Explanatory Questions

(5 Pts.)



Please answer the following questions as clear as possible:

- K-Means with typical settings can only return *circular* clusters. However, in many cases clusters are not circular and may appear as lines. Give a method or a distance function to capture both *circular* clusters and *appear-as-line* clusters.
- How can you relate K-Means and K-NN algorithms?
- If we run K-Means algorithm several times with the same number of clusters and starting points, does it always converge to the same solution? Why / Why not?

Good Luck!
Ali Abbasi, Farhad Dalirani