

CSE 4309 - Assignments - Assignment 6

List of assignment due dates.

The assignment should be submitted via [Canvas](#). Submit a file called assignment6.zip, containing the following two files:

- answers.pdf, for your answers to the written tasks, and for the output that the programming task asks you to include. Only PDF files will be accepted. All text should be typed, and if any figures are present they should be computer-generated. Scans of handwritten answers will NOT be accepted.
- k_means.m, or k_means.py, containing your Matlab or Python code for the programming part. Your Matlab or Python code should run on the versions specified in the syllabus, unless permission is obtained via e-mail from the instructor or the teaching assistant. Also submit any other files that are needed in order to document or run your code (for example, additional source code files).

These naming conventions are mandatory, non-adherence to these specifications can incur a penalty of up to 20 points.

Your name and UTA ID number should appear on the top line of both documents.

Task 1 (70 points, programming)

In this task you will implement k-means clustering.

Arguments

You must implement a Matlab function or a Python executable file called `k_means`. Your function should be invoked as follows:

```
k_means(<data_file>, <K>, <initialization>)
```

If you use Python, just convert the Matlab function arguments shown above to command-line arguments. The arguments provide to the function the following information:

The arguments provide the following information:

- The first argument, `<data_file>`, is the path name of a file where the data is stored. The path name can specify any file stored on the local computer.
- The second argument, `<K>`, specifies the number of clusters.
- The third argument, `<initialization>` specifies how the initial assignment of data points to clusters is done. This argument can have two possible values: `random`, or `round_robin`.

Example 1D and 2D datasets that you can test your code with can be found in the [toy_data](#) directory. In each file, each row corresponds to a data point, which can be a single number or a 2D vector. The two values of each 2D vector are separated by space.

Your code should work with any other files that follow the same format. You can assume that the input data will be either 1D or 2D, you do not need to worry about higher dimensions.

Implementation Guidelines

- To make the result deterministic, if the third command-line argument is `round_robin`, the initial cluster assignments should be done in round-robin fashion. More specifically, the first object gets assigned to cluster 1, the second object gets assigned to 2, ..., the k -th object is assigned to cluster K , and so on. For example, for file [set1a.txt](#), and $K=3$, this is how the initial assignments of points to clusters should be:
 - 7 --> cluster 1
 - 29 --> cluster 2
 - 11 --> cluster 3
 - 2 --> cluster 1
 - 16 --> cluster 2
 - 4 --> cluster 3
 - 37 --> cluster 1
 - 22 --> cluster 2

Note that cluster IDs range from 1 to K , and there is NO cluster ID 0.

- Your algorithm should terminate when the assignment of objects to clusters stops changing. In other words, the algorithm should stop if the assignment of objects to clusters at the end of the current iteration is identical to the assignment at the end of the previous iteration.

Output

At the end, your program should print the final cluster assignments. If the dataset is one-dimensional (like file [set1a.txt](#)), the output should contain one line per data point, printed with the following formatting specification:

```
fprintf('%10.4f --> cluster %d\n', data_point, cluster_id);
```

If the dataset is two-dimensional (like file [set2a.txt](#)), the output should contain again one line per data point, printed with the following formatting specification:

```
fprintf('(%10.4f, %10.4f) --> cluster %d\n', data_point[0], data_point[1], cluster_id);
```

You can see examples of results printed using this format right below.

Example Results

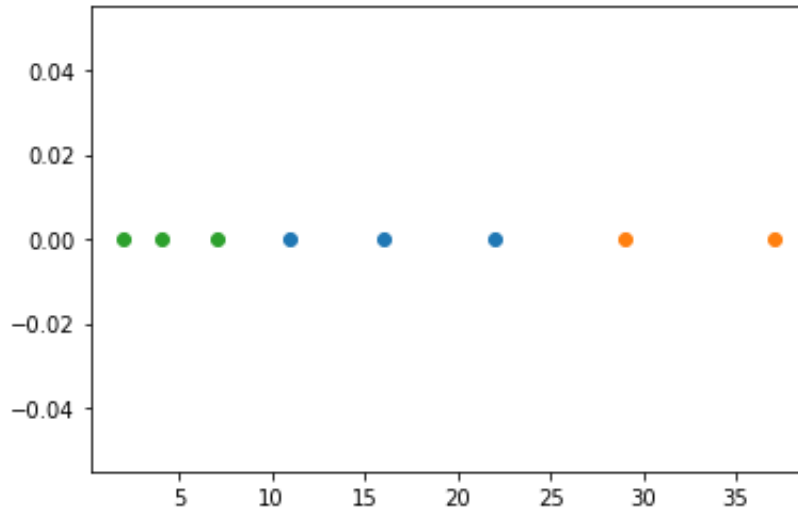
The following files show some results that my code gets. In all cases, the runtime is under one second.

- File [set1a_3.txt](#) shows the output for [set1a.txt](#), with $K=3$ and round-robin initialization.
- File [set2a_3.txt](#) shows the output for [set2a.txt](#), with $K=3$ and round-robin initialization.
- File [set2b_3.txt](#) shows the output for [set2b.txt](#), with $K=3$ and round-robin initialization.
- File [set2c_2.txt](#) shows the output for [set2c.txt](#), with $K=2$ and round-robin initialization.

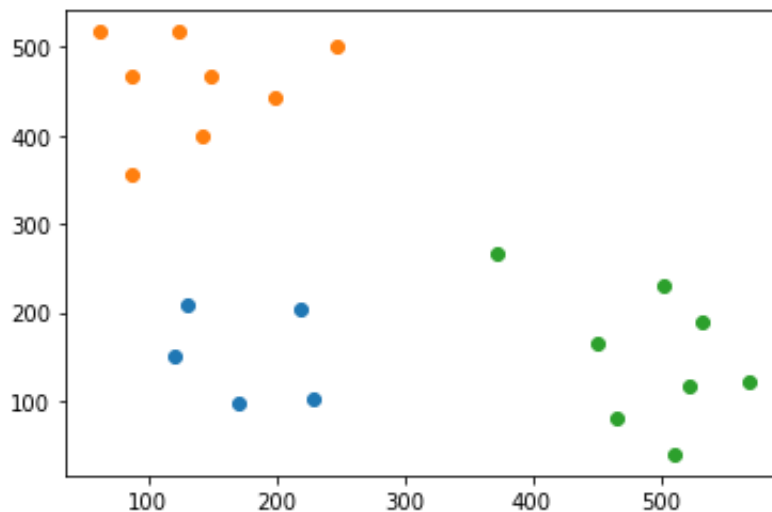
Plotting Datasets and Results

This is entirely optional, but if you are interested, Python file [drawing.py](#) contains code that you can use to visualize the datasets and results. Function `draw_points (data)` plots a 1D or 2D dataset. Function `draw_assignments(data, assignments)` plots each cluster with a different color. Here are some examples of such plots:

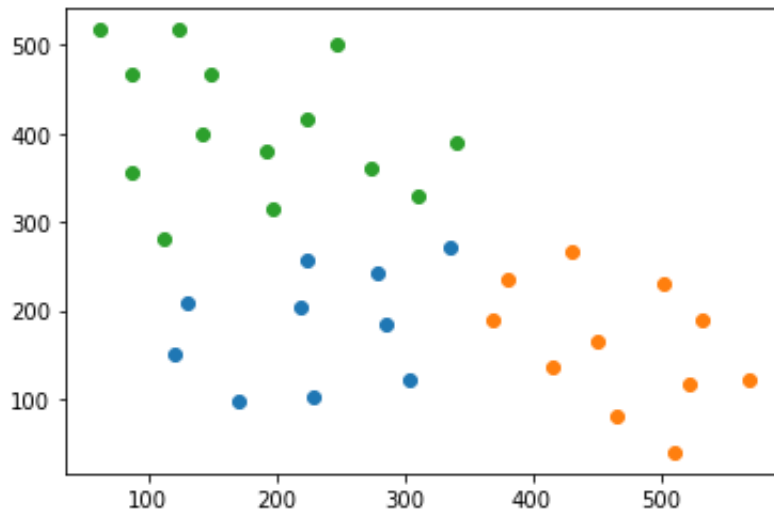
- Plot of the result for [set1a.txt](#), with $K=3$ and round-robin initialization:



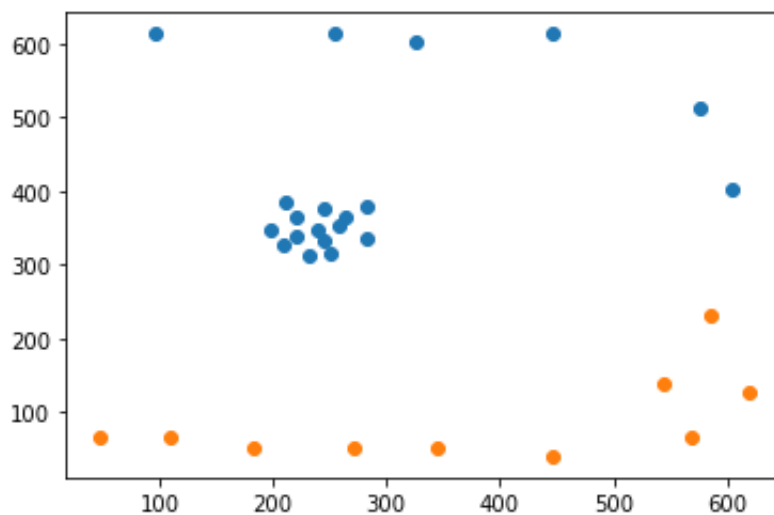
- Plot of the result for [set2a.txt](#), with $K=3$ and round-robin initialization:



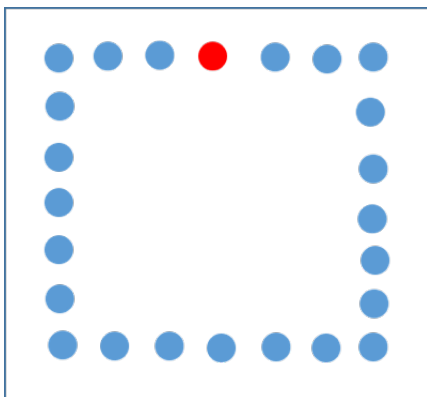
- Plot of the result for [set2b.txt](#), with $K=3$ and round-robin initialization:



- Plot of the result for [set2c.txt](#), with $K=2$ and round-robin initialization:



Task 2 (10 points)



Consider the set of points above. Each dot corresponds to a point. Consider a clustering consisting of two clusters, where the first cluster is the set of all the blue dots, and the second cluster has the red dot as its only element. Can this clustering be the final result of the k-means algorithm? Justify your answer.

Your answer should be based on your visual estimations of Euclidean distances between dots. For this particular question, no greater precision is needed.

Task 3 (10 points)

For both parts of this task, assume that the algorithm in question **never needs to break a tie** (i.e., it never needs to choose between two distances that are equal).

Part a: Will the EM algorithm always give the same results when applied to the same dataset with the same K ? To phrase the question in an alternative way, is there any dataset where the algorithm can produce different results if run multiple times, with the value of K kept the same? If your answer is that EM will always give the same results, justify why. If your answer is that EM can produce different results if run multiple times, provide an example.

Part b: Same question as in part a, but for agglomerative clustering with the d_{\min} distance. Here, as "result" we consider all intermediate clusterings, between the first step (with each object being its own cluster) and the last step (where all objects belong to a single cluster). Will this agglomerative algorithm always give the same results when applied to the same dataset? If your answer is "yes", justify why. If your answer is "no", provide an example.

Task 4 (10 points)

Consider a dataset consisting of these eight points: 2, 4, 7, 11, 16, 22, 29, 37.

Part a: Show the results (all intermediate clusterings) obtained by applying agglomerative clustering to this dataset, using the d_{\min} distance.

Part b: Show the results (all intermediate clusterings) obtained by applying agglomerative clustering to this dataset, using the d_{\max} distance.

Extra Credit Task (optional, 10 points extra credit)

Implement EM clustering. For more precise instructions, please refer to [this page](#).

Optional Tasks, No Credit

Here some tasks that are optional, and for which you will receive NO credit. These tasks will give you some additional hands-on experience with clustering.

- **Optional Task 1, 0 points:** If you have implemented EM clustering, implement a way to identify when the EM algorithm has converged, so as to stop the iterations of the main loop.

- **Optional Task 2, 0 points:** Implement k-medoid clustering, and apply it on the time series dataset of the [optional DTW assignment](#), with DTW as the underlying distance measure.
- **Optional Task 3, 0 points:** Implement agglomerative clustering, with d_{mean} as the distance measure between clusters. Apply your implementation on the time series dataset of the [optional DTW assignment](#), with DTW as the underlying distance measure.

[CSE 4309](#) - [Assignments](#) - Assignment 6