

# Project 2: K-means-Strategy

## Purpose

In this project, you are required to implement the K-means algorithm and apply your implementation on the given dataset (AllSamples.npy), which contains a set of 2-D points. You are required to implement two different strategies for choosing the initial cluster centers.

## Technology Requirements

### Algorithms:

- k-Means Clustering

### Resources:

- A 2-D dataset to be provided

### Workspace:

- Any Python programming environment

### Software:

- Python environment

### Language(s):

- Python

## Directions

Download Mat File: “CSE 575\_Project 2\_AllSamples” (attached in the project description in the course)

### Lab: Project 2: K-mean-Strategy, Part 1

#### K-means\_algorithm\_Strategy Part 1 (K-mean) Overview:

You are required to implement the following strategy for choosing the initial cluster centers.

Part 1 is to randomly pick the initial centers from the given samples.

You need to test your implementation on the given data, with the number  $k$  of clusters ranging from 2-10, output the final coordinate of the centroids and compute the loss based on the objective function.

(Referring to the course notes: When clustering the samples into  $k$  clusters/sets  $D_i$ , with respective center/mean vectors  $\mu_1, \mu_2, \dots, \mu_k$ , the objective function is defined as  $\sum_{i=1}^k \sum_{x \in D_i} \|x - \mu_i\|^2$ )

You are highly suggested to use the built-in Jupiter Notebook to implement your algorithm. Another python environment is okay but you must be responsible for any numerical error caused by a different programming environment.

### Lab: Project 2: K-mean-Strategy, Part 2

#### K-means\_algorithm\_Strategy2 (K-mean ++) Overview:

You are required to implement the following strategy for choosing the initial cluster centers.

Part 2 is to pick the first center randomly; for the  $i$ -th center ( $i > 1$ ), choose a sample (among all possible samples) such that the average distance of this chosen one to all previous ( $i - 1$ ) centers is maximal.

You need to test your implementation on the given data, with the number  $k$  of clusters ranging from 2-10, output the final coordinate of the centroids and compute the loss based on the objective function.

(Referring to the course notes: When clustering the samples into  $k$  clusters/sets  $D_i$ , with respective center/mean vectors  $\mu_1, \mu_2, \dots, \mu_k$ , the objective function is defined as  $\sum_{i=1}^k \sum_{x \in D_i} \|x - \mu_i\|^2$

You are highly suggested to use the built-in Jupiter Notebook to implement your algorithm. Another python environment is okay but you must be responsible for any numerical error caused by a different programming environment.

## Required Tasks:

1. Write code to implement the k-means algorithm with Strategy 1.
2. Use your code to do clustering on the given data; compute the objective function as a function of  $k$  ( $k = 2, 3, \dots, 10$ ).
3. Repeat the above step with another initialization.
4. Write code to implement the k-means algorithm with Strategy 2.
5. Use your code to do clustering on the given data; compute the objective function as a function of  $k$  ( $k = 2, 3, \dots, 10$ ).
6. Repeat the above step with another initialization.
7. Submit a short report summarizing the results, including the plots for the objective function values under different settings described above.

## Submission Directions for Project Deliverables

### What to Submit:

1. **Result Submissions:** Code file with comments explaining what you do for each part as directed in their respective result submission quiz.
2. **Report Submission:** A report that summarizes the results and includes the plots for each of the objective function values.

## Result Submission

### Quiz: K-means-Strategy, Part 1 Result Submission

Once you launch the lab, you will be able to find two sets of initial  $k$  and points, which are associated with your ID. Please test your algorithm with this initialization.

The 1) final coordinate of the centroids and 2) the loss computed by the objective function should be submitted to the quiz titled "**Quiz: K-means-Strategy, Part 1 Result Submission**".

**Note:**

1. You should implement your own K-means algorithm.
2. You will not get any points for the project by simply programming here. Please remember you need to submit the results in the result submission quiz.

## **Quiz: K-means-Strategy, Part 2 Result Submission**

Once you launch the lab, you will be able to find two sets of initial  $k$  and points, which are associated with your ID. Please test your algorithm with this initialization.

The 1) final coordinate of the centroids and 2) the loss computed by the objective function should be submitted to the quiz titled "**Quiz: K-means-Strategy, Part 2 Result Submission**".

**Note:**

1. You should implement your own K-means algorithm.
2. You will not get any points for the project by simply programming in the Lab. Please remember you need to submit the results in the result submission quiz.

## **Report Submission**

### **Graded Assignment: K-means-Strategy, Report Submission**

Please submit your report regarding Project 2 to the item titled "**Graded Assignment: K-means-Strategy, Report Submission**".

- Acceptable file types: .pdf or .doc/docx.
- Length of the report: no more than 2 A4 pages.
- Content: (The following must be included)
  - Please include the  $k$  and initial points assigned to you as well as the final clustering centroid and loss in the report.
  - Your observation and analysis about the two strategies.

## Evaluation

### Report Submission

- 1 point for the final clustering centroid and the loss
- 1 point for the analysis.