

## ENGR 421 / DASC 521: Introduction to Machine Learning

### Homework 06: Spectral Clustering

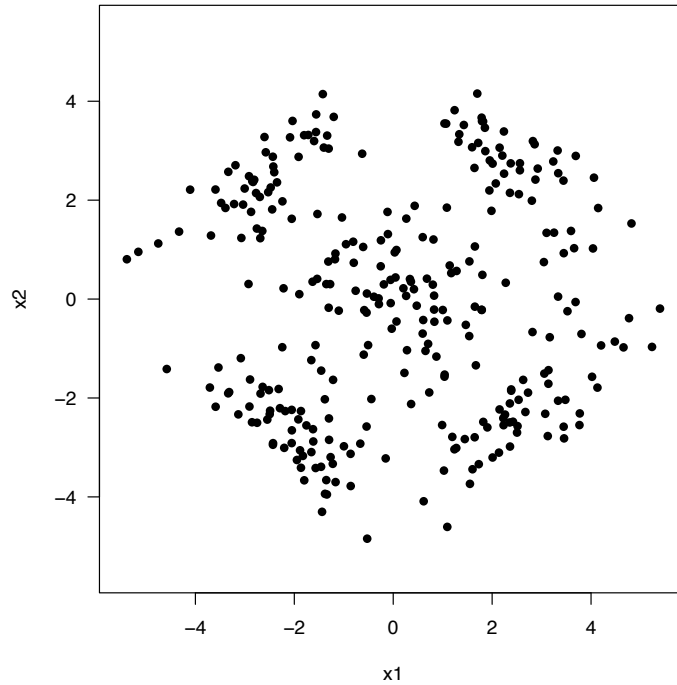
Deadline: June 11, 2021, 11:59 PM

In this homework, you will implement a spectral clustering algorithm in Python. Here are the steps you need to follow:

1. You are given a two-dimensional data set in the file named `hw06_data_set.csv`, which contains 300 data points generated randomly from five bivariate Gaussian densities with the following parameters.

$$\begin{aligned}\mu_1 &= \begin{bmatrix} +2.5 \\ +2.5 \end{bmatrix}, & \Sigma_1 &= \begin{bmatrix} +0.8 & -0.6 \\ -0.6 & +0.8 \end{bmatrix}, & N_1 &= 50 \\ \mu_2 &= \begin{bmatrix} -2.5 \\ +2.5 \end{bmatrix}, & \Sigma_2 &= \begin{bmatrix} +0.8 & +0.6 \\ +0.6 & +0.8 \end{bmatrix}, & N_2 &= 50 \\ \mu_3 &= \begin{bmatrix} -2.5 \\ -2.5 \end{bmatrix}, & \Sigma_3 &= \begin{bmatrix} +0.8 & -0.6 \\ -0.6 & +0.8 \end{bmatrix}, & N_3 &= 50 \\ \mu_4 &= \begin{bmatrix} +2.5 \\ -2.5 \end{bmatrix}, & \Sigma_4 &= \begin{bmatrix} +0.8 & +0.6 \\ +0.6 & +0.8 \end{bmatrix}, & N_4 &= 50 \\ \mu_5 &= \begin{bmatrix} +0.0 \\ +0.0 \end{bmatrix}, & \Sigma_5 &= \begin{bmatrix} +1.6 & +0.0 \\ +0.0 & +1.6 \end{bmatrix}, & N_5 &= 100\end{aligned}$$

The given data points are shown in the following figure.

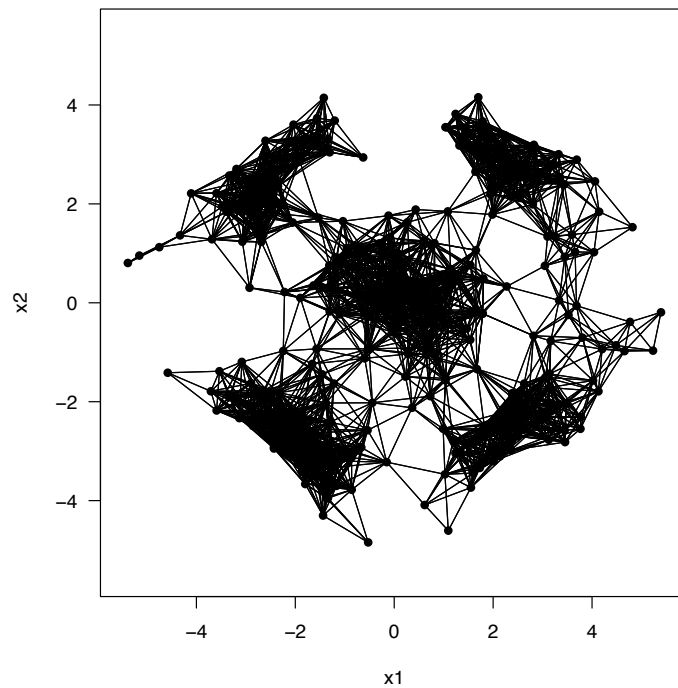


2. You should first calculate the Euclidean distances between the pairs of data points. The data point pairs with distance less than or equal to  $\delta = 1.25$  are considered as connected. Construct the matrix **B** as follows:

$$b_{ij} = \begin{cases} 0, & \|\mathbf{x}_i - \mathbf{x}_j\|_2 > \delta \\ 1, & \text{otherwise.} \end{cases}$$

$$b_{ii} = 0$$

You should also visualize this connectivity matrix by drawing a line between two data points if they are connected. Your figure should be similar to the following figure.

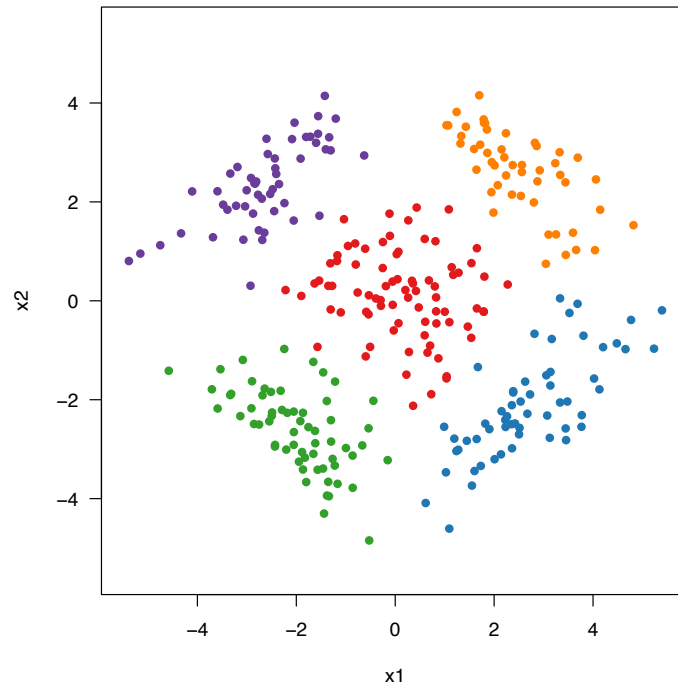


3. You should then calculate **D** and **L** matrices as described in the lecture notes. You should normalize the Laplacian matrix using the following formula:

$$\mathbf{L}_{symmetric} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{B} \mathbf{D}^{-1/2}$$

4. Find the eigenvectors of the normalized Laplacian matrix and pick  $R = 5$  eigenvectors that corresponds to  $R$  smallest eigenvectors. Using these eigenvectors construct the matrix **Z** as described in the lecture notes.
5. Run  $k$ -means clustering algorithm on **Z** matrix to find  $K = 5$  clusters. When initializing your algorithm, use the following rows of **Z** matrix for initial centroids: 85, 129, 167, 187, and 270.

6. Draw the clustering result obtained by your spectral clustering algorithm by coloring each cluster with a different color. Your figure should be similar to the following figure.



**What to submit:** You need to submit your source code in a single file (.py file) and a short report explaining your approach (.doc, .docx, or .pdf file). You will put these two files in a single zip file named as ***STUDENTID.zip***, where ***STUDENTID*** should be replaced with your 7-digit student number.

**How to submit:** Submit the zip file you created to Blackboard. Please follow the exact style mentioned and do not send a zip file named as ***STUDENTID.zip***. Submissions that do not follow these guidelines will not be graded.

**Late submission policy:** Late submissions will not be graded.

**Cheating policy:** Very similar submissions will not be graded.