

ENGR 421/DASC 521: Introduction to Machine Learning

Homework 5: Expectation-Maximization Clustering

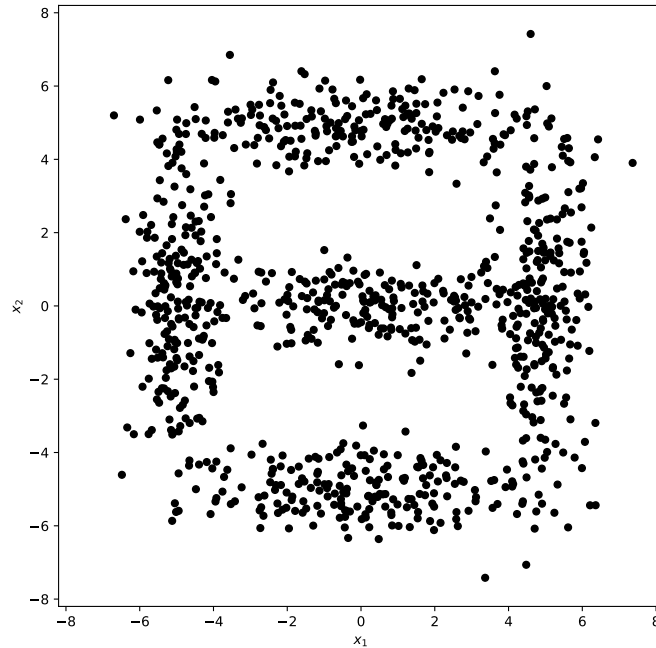
Deadline: May 28, 2025, 11:59 PM

In this homework, you will implement an expectation-maximization (EM) 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 `hw05_data_set.csv`, which contains 1000 data points randomly generated from five bivariate Gaussian densities with the following parameters.

$$\begin{array}{lll} \boldsymbol{\mu}_1 = \begin{bmatrix} -5.0 \\ +0.0 \end{bmatrix} & \boldsymbol{\Sigma}_1 = \begin{bmatrix} +0.4 & +0.0 \\ +0.0 & +6.0 \end{bmatrix} & N_1 = 200 \\ \boldsymbol{\mu}_2 = \begin{bmatrix} +0.0 \\ +5.0 \end{bmatrix} & \boldsymbol{\Sigma}_2 = \begin{bmatrix} +6.0 & +0.0 \\ +0.0 & +0.4 \end{bmatrix} & N_2 = 200 \\ \boldsymbol{\mu}_3 = \begin{bmatrix} +5.0 \\ +0.0 \end{bmatrix} & \boldsymbol{\Sigma}_3 = \begin{bmatrix} +0.4 & +0.0 \\ +0.0 & +6.0 \end{bmatrix} & N_3 = 200 \\ \boldsymbol{\mu}_4 = \begin{bmatrix} +0.0 \\ -5.0 \end{bmatrix} & \boldsymbol{\Sigma}_4 = \begin{bmatrix} +6.0 & +0.0 \\ +0.0 & +0.4 \end{bmatrix} & N_4 = 200 \\ \boldsymbol{\mu}_5 = \begin{bmatrix} +0.0 \\ +0.0 \end{bmatrix} & \boldsymbol{\Sigma}_5 = \begin{bmatrix} +6.0 & +0.0 \\ +0.0 & +0.4 \end{bmatrix} & N_5 = 200 \end{array}$$

The given data points are shown in the following figure.



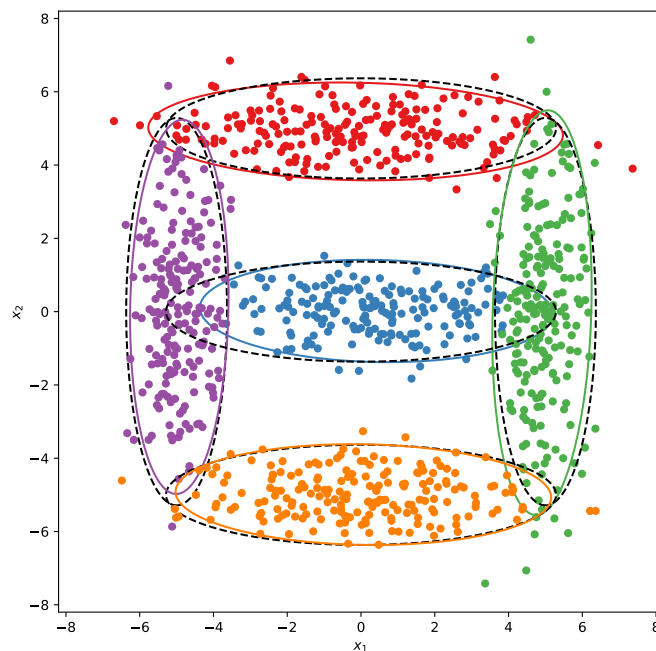
2. To initialize your EM algorithm, you should take the centroids given in the file named `hw05_initial_centroids.csv` as the initial values for the mean vectors. By assigning the data points to the nearest center, estimate the initial covariance matrices and prior probabilities in your EM algorithm. (20 points)
3. After the initialization step, run your EM algorithm for 100 iterations. Report the mean vectors your EM algorithm finds. Your results should be like the following matrix. (50 points)

```

print(means)
[[-0.14662825  4.91787672]
 [ 0.41688612  0.01706596]
 [ 4.90712573 -0.02874697]
 [-4.92579882  0.1300626 ]
 [ 0.06413512 -4.99302976]]
print(priors)
[0.20056729 0.19159891 0.20987053 0.20081455 0.19714872]

```

4. Draw the clustering result obtained by your EM algorithm by coloring each cluster with a different color. You should also draw the original Gaussian densities you use to generate data points and the Gaussian densities your EM algorithm finds with dashed and solid lines, respectively. Draw these Gaussian densities where their values are equal to 0.01. Your figure should be like the following figure. (30 points)



What to submit: You need to submit your source code in a single file (.py file). You are provided with a template file named as 0099999.py, where 99999 should be replaced with your 5-digit student number. You are allowed to change the template file between the following lines.

```
# your implementation starts below
```

```
# your implementation ends above
```

How to submit: Submit the file you edited to LearnHub by following the exact style mentioned. 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.
