

- Submit Gxxx.ZIP in Fenix where xxx is your group number. The ZIP should contain two files: Gxxx_report.pdf with your report and Gxxx_notebook.ipynb with your notebook demo according to the suggested templates
- It is possible to submit several times on Fenix to prevent last-minute problems. Yet, only the last submission is kept
- Exchange of ideas is encouraged. Yet, if copy is detected after automatic or manual clearance, homework is nullified and IST guidelines apply for content sharers and consumers, irrespectively of the underlying intent
- Please consult the FAQ before posting questions to your faculty hosts

I. Pen-and-paper [9v]

Consider the bivariate observations $\{\mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \mathbf{x}_2 = \begin{bmatrix} 0 \\ 2 \end{bmatrix}, \mathbf{x}_3 = \begin{bmatrix} 3 \\ -1 \end{bmatrix}\}$ and the multivariate

Gaussian mixture given by

$$\mathbf{u}_1 = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \quad \mathbf{u}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \Sigma_1 = \begin{bmatrix} 4 & 1 \\ 1 & 4 \end{bmatrix}, \quad \Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad \pi_1 = 0.5, \quad \pi_2 = 0.5$$

Answer the following questions by presenting all intermediary steps, and use 3 decimal places in each.

1. [6v] Perform two epochs of the EM clustering algorithm and determine the new parameters.
2. Using the final parameters computed in previous question:
 - a. [1v] Perform a hard assignment of observations to clusters under a MAP assumption.
 - b. [2v] Compute the silhouette of the larger cluster (the one that has more observations assigned to it) using the Euclidean distance.

II. Programming and critical analysis [11v]

In the next exercise you will use the `accounts.csv` dataset. This dataset contains account details of bank clients, and the target variable y is binary ('has the client subscribed a term deposit?'). Select the first 8 features and remove duplicates and null values.

Hint: You can use `get_dummies()` to change the feature type (e.g. `pd.get_dummies(data, drop_first=True)`).

1. Normalize the data using MinMaxScaler:
 - a. [4v] Using *sklearn*, apply k -means clustering (without targets) on the normalized data with $k=\{2,3,4,5,6,7,8\}$, `max_iter=500` and `random_state=42`. Plot the different sum of squared errors (SSE) using the `_inertia` attribute of k -means according to the number of clusters.
 - b. [1.5v] According to the previous plot, how many underlying customer segments (clusters) should there be? Explain based on the trade-off between the clusters and inertia.
 - c. [1.5v] Would k -modes be a better clustering approach? Explain why based on the dataset features.
2. Normalize the data using StandardScaler:
 - a. [1v] Apply PCA to the data. How much variability is explained by the top 2 components?
 - b. [1v] Apply k -means clustering with $k=3$ and `random_state=42` (all other arguments as default) and use the original 8 features. Next, provide a scatterplot according to the first 2 principal components. Can we clearly separate the clusters? Justify.
 - c. [2v] Plot the cluster conditional features of the frequencies of "job" and "education" according to the clusters obtained in the previous question (2b.). Use `sns.displot` (see Data Exploration notebook), with `multiple="dodge"`, `stat='density'`, `shrink=0.8` and `common_norm=False`. Describe the main differences between the clusters in no more than half a page.

END