

LECTURER: TAI LE QUY

MACHINE LEARNING

UNSUPERVISED LEARNING AND FEATURE ENGINEERING

INTRODUCTION TO UNSUPERVISED MACHINE LEARNING AND FEATURE
ENGINEERING

1

CLUSTERING

2

DIMENSIONALITY REDUCTION

3

FEATURE ENGINEERING

4

FEATURE SELECTION

5

AUTOMATED FEATURE GENERATION

6

UNIT 2

CLUSTERING



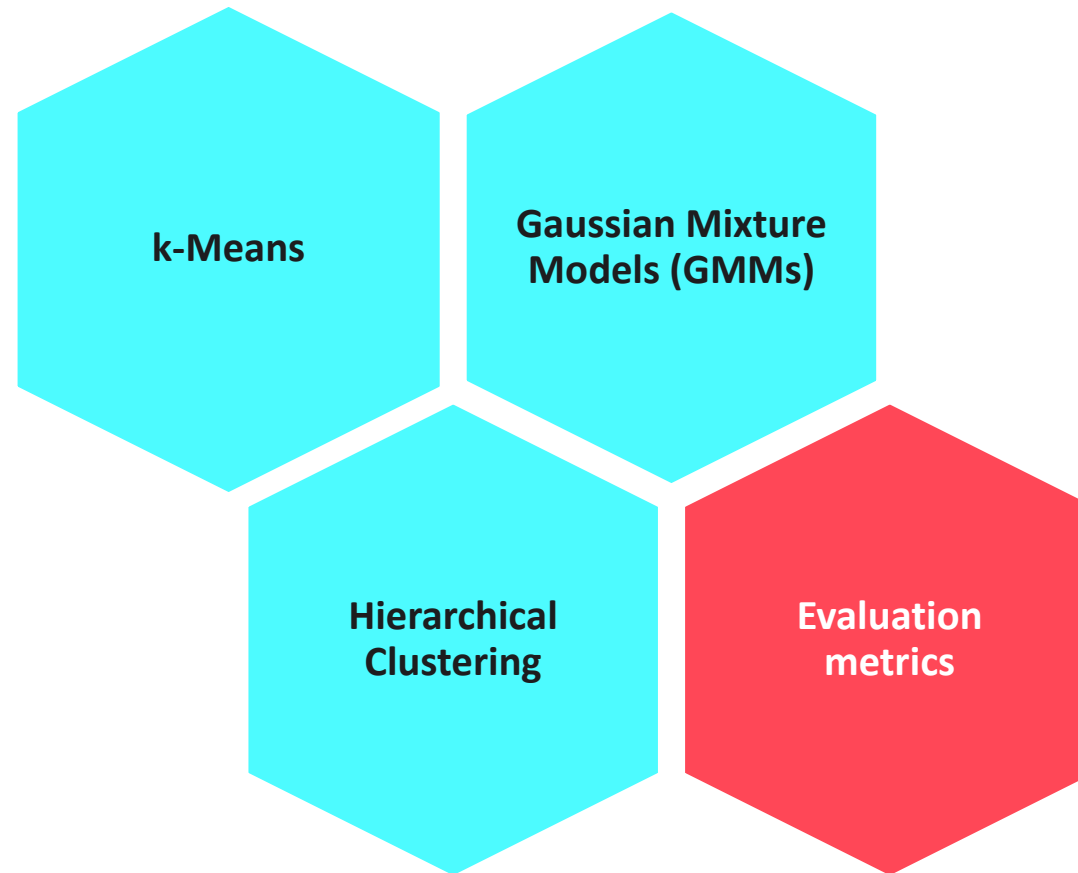
- Explain the **functioning principal of clustering** approaches and how they work.
- **Implement** a clustering approach.
- **Test and evaluate the quality** of the obtained clusters.
- **Choose the clustering approach** with respect to the challenges and constraints of the dataset.



1. Explain how it is possible to obtain **two different clustering results** for the **same dataset** using **k-Means** clustering.
2. In k-Means, the **centroids** are **updated** in each iteration. Explain the equivalent in **Gaussian Mixture Models** that is updated in each iteration.
3. For a **100-sample dataset**, explain **how many samples** will be in **each leaf** and how many will be in the **stem** of the dendrogram when applying hierarchical clustering.

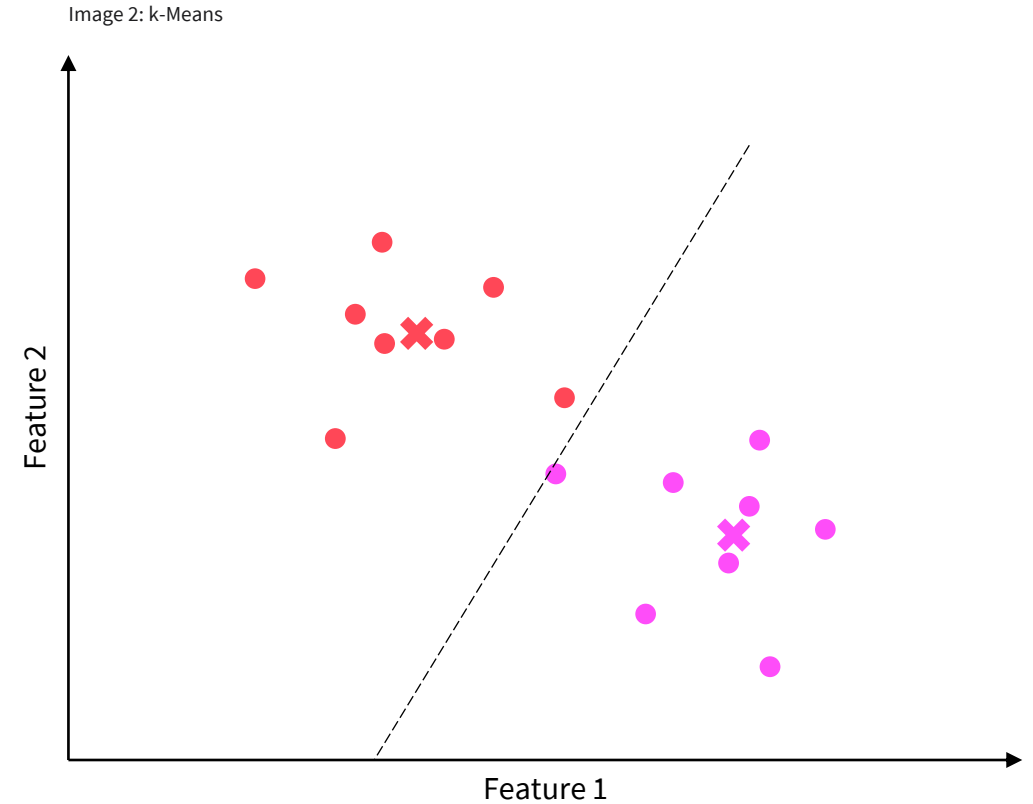
UNIT CONTENT

Image 1: Unit content - Clustering



K-MEANS

1. Choose a **number of clusters, k**
2. **Randomly select** a data point for each cluster (**seed = interim centroid**).
3. Calculate the **distance** between **each data point** and the **centroids**.
4. Assign each data point **to the nearest centroid**.
5. Select new centroids as the **mid-point** of each cluster.
6. Repeat steps 3 to 5 until the **stop criterion is fulfilled**.



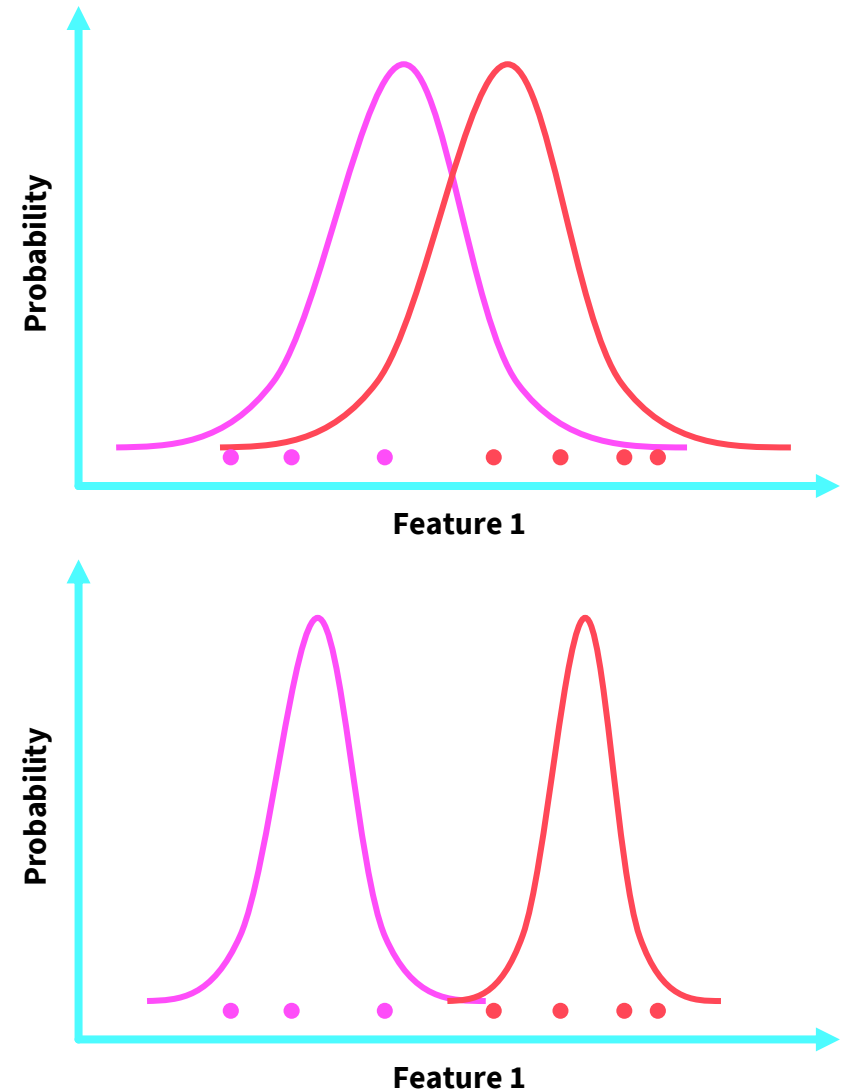
K-MEANS

- k-Means is **not deterministic**.
- There are several **variations** to k-Means.
- for large datasets
 - **Clustering for Large Applications (CLARA)**
 - Partitioning
 - Batch-processing

GAUSSIAN MIXTURE MODEL CLUSTERING (GMM)

1. Choose „**prior probabilities**“ at random.
2. Assign **each sample** to the **closest cluster centroid** based on the **Maximum Likelihood**.
3. **Re-calculate** the cluster **centroids** based on the **mean and variance** of the samples in this cluster.
4. Repeat steps 2 and 3 until the **stop criterion is fulfilled**.

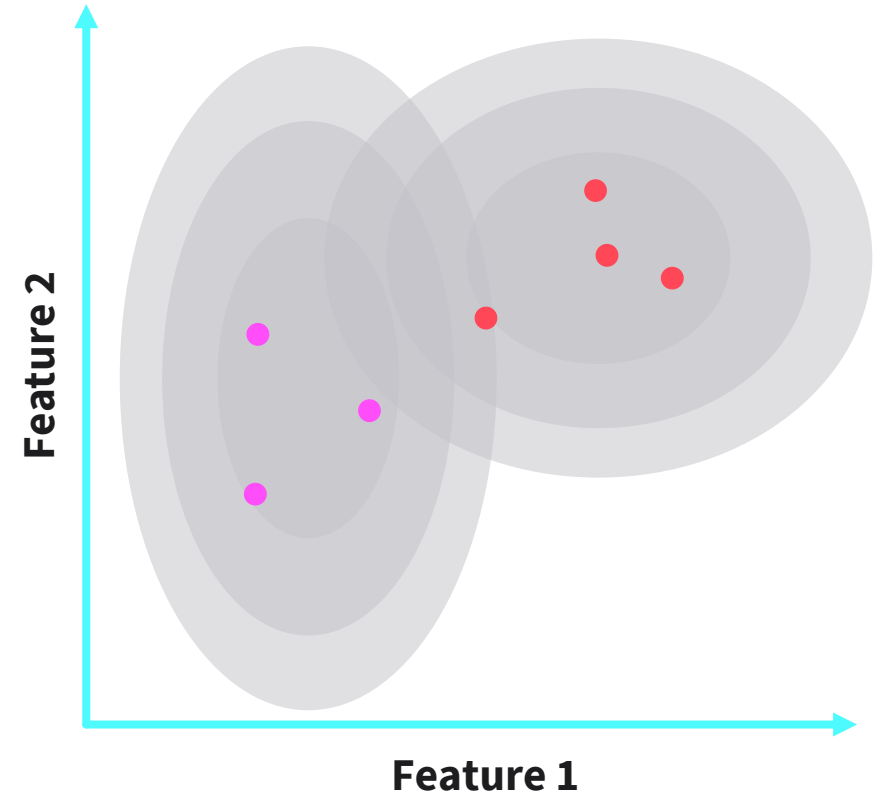
Image 3: GMM



Advantages

- **"Fuzzy" clusters** with probability zones.
- Different **"probability slopes"** for each feature.

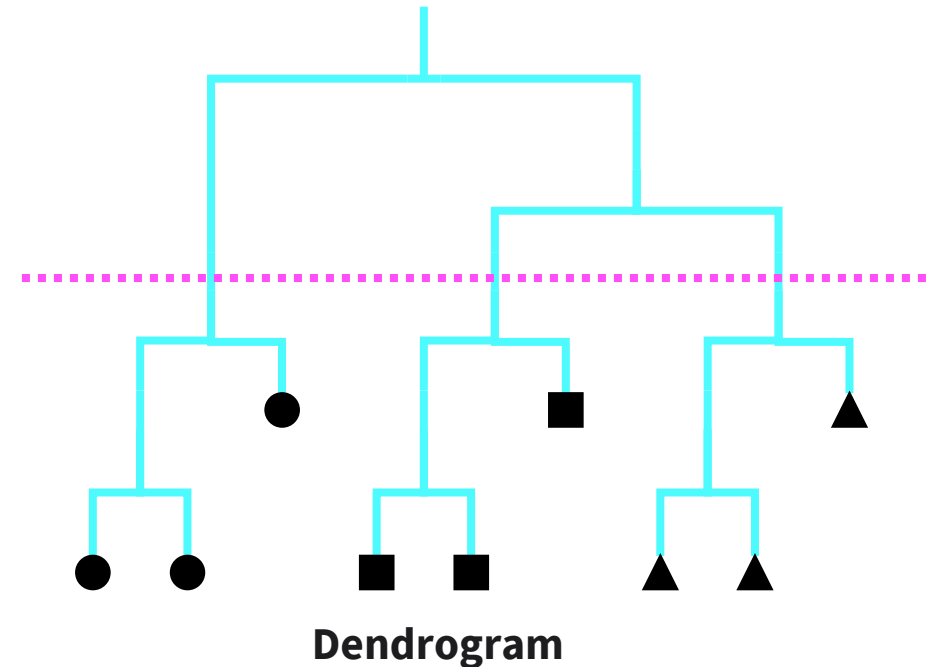
Image 4: GMM clusters with probability zones



HIERARCHICAL CLUSTERING

1. **Each sample** in one **cluster** (leaves).
2. Calculate **distances** between all samples.
3. Group the **two closest samples**, respectively.
4. Continue **grouping samples and groups**.
5. Ultimately, all sample end up in **one cluster** (stem).
6. Choose the **number of clusters** by horizontally drawing the **decision boundary** through the dendrogram.

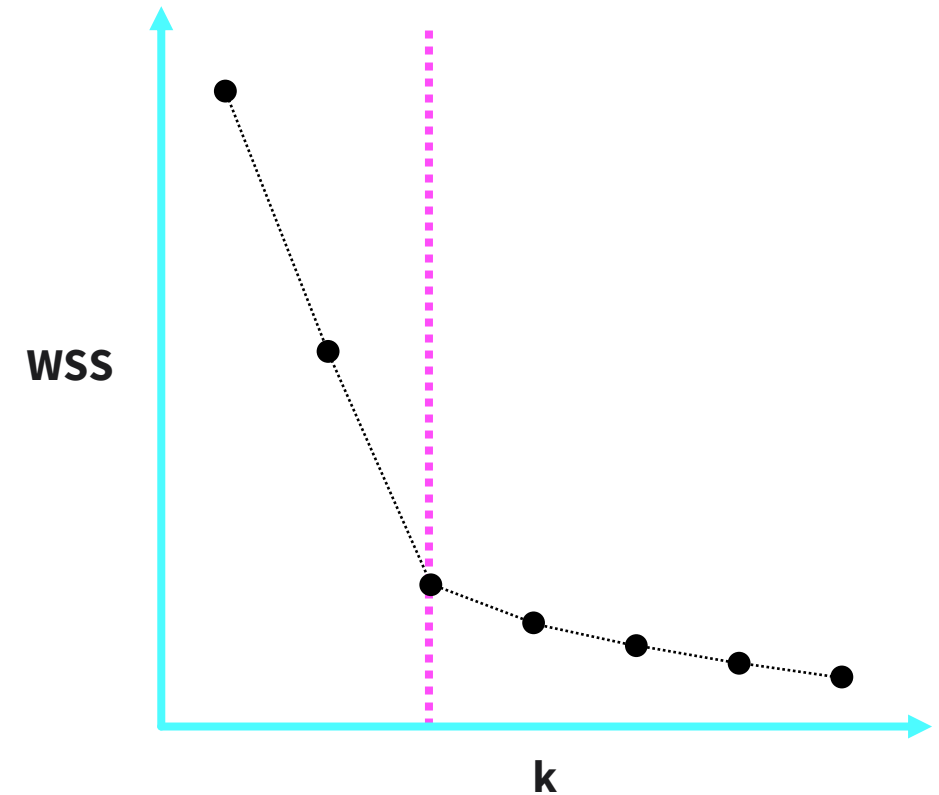
Image 5: Hierarchical clustering



Elbow method with WSS

- **Within-Cluster Sum of Squares (WSS)**
 - Squared **distance** between **data points** and respective cluster **centroids**.
 - $$WSS = \sum_{j=1}^k \sum_{x_i \in c_j} (x_i - c_j)^2$$

Image 6: Elbow method

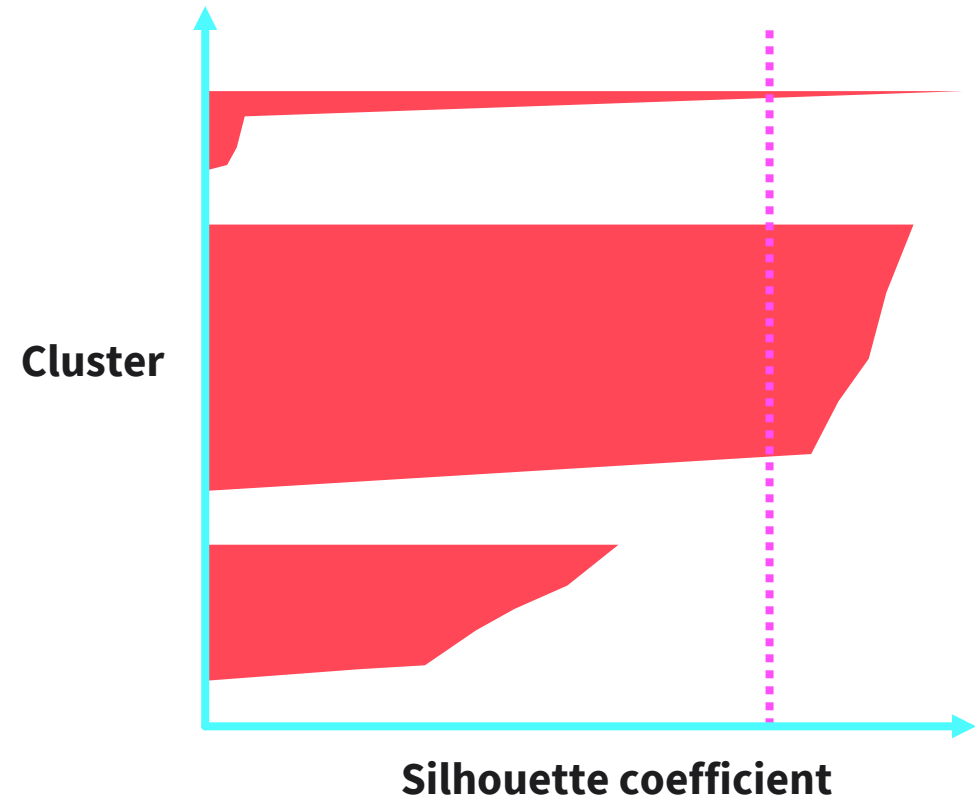


EVALUATION METRICS

Silhouette Score

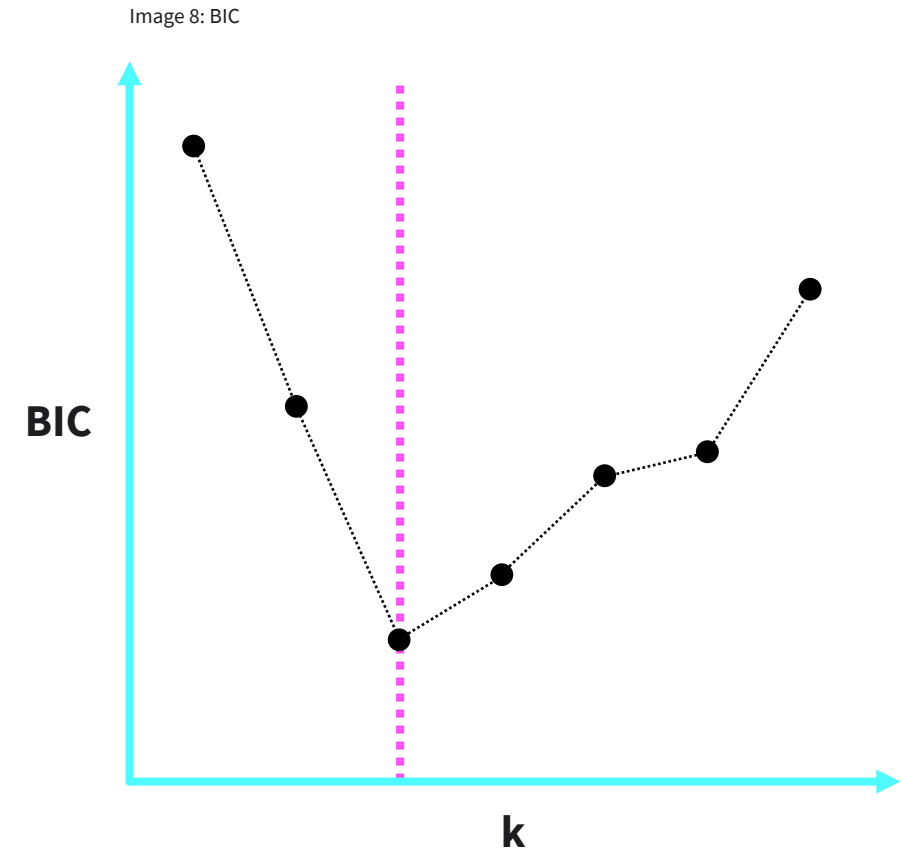
- Cohesion
- Separation
- Range $[-1, 1]$
- For each data point
- As overall metric

Image 7: Silhouette score



Bayesian Information Criterion (BIC)

- $BIC = \ln(n) \cdot p - 2\ln(L)$
- n = number of samples
- p = number of parameters
- L = Maximum Likelihood





- Explain the **functioning principal of clustering** approaches and how they work.
- **Implement** a clustering approach.
- **Test and evaluate the quality** of the obtained clusters.
- **Choose the clustering approach** with respect to the challenges and constraints of the dataset.

SESSION 2

TRANSFER TASK

TRANSFER TASKS

A start-up that sells **sustainable products in smaller stores** has been very successful in recent years. As a result, more stores are to be opened worldwide.

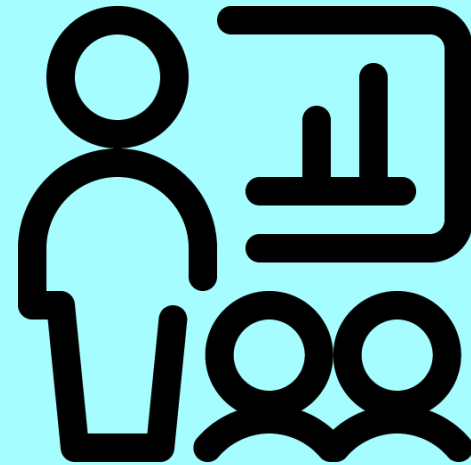
To keep an **overview of the offered products**, you and your team of Data Scientists are tasked to **define homogeneous groups of products** to facilitate ordering, marketing, and distribution. There are **different use cases** for your results, and you should use different methods appropriate to each use case:

1. The customer base **constantly changes**, and the clustering must be conducted **as quickly as possible**.
2. Once a month, a **more thorough analysis** should be conducted. **Not all features** seem to be **equally informative** to differentiate the customers into groups.
3. The **number of clusters** has to be **adapted on-the-fly** for the ordering process to quickly assess how many different products should be ordered in bulk.

TRANSFER TASK
PRESENTATION OF THE RESULTS

Please present your
results.

The results will be
discussed in plenary.





1. What does the elbow criterion consider when assessing the quality of clusters?
 - a) the cohesion of the clusters
 - b) the separability of the clusters
 - c) the cohesion and separability of the clusters
 - d) the non-convex shape of the clusters



2. A silhouette score indicates a high quality of clusters when the value is...
- a) ... close to 0.
 - b) ... close to -1 .
 - c) ... larger than 1.
 - d) ... close to 1.



3. Which of the following propositions is correct when a Gaussian mixture model is used?
- a) A data point has 1 as a membership value to one cluster and 0 for the other clusters.
 - b) A data point has probability membership values to the different clusters.
 - c) The provided clusters do not depend on the initialization.
 - d) There is no need to define the number of clusters in advance.

LIST OF SOURCES

Images

- Müller-Kett, 2020.
- Müller-Kett, 2021.
- Müller-Kett, 2023.
- Microsoft Archive.

© 2022 IU Internationale Hochschule GmbH

This content is protected by copyright. All rights reserved.

This content may not be reproduced and/or electronically edited, duplicated, or distributed in any kind of form without written permission by the IU Internationale Hochschule GmbH.