# Unit-5

# Clustering

# Introduction to clustering

### **Q** Introduction to Clustering

Clustering is a fundamental technique in **unsupervised machine learning** that involves grouping a set of data points into **clusters**  $\square$ , such that data points in the same group are more similar to each other than to those in other groups. It is widely used in **data analysis**, **pattern recognition**, **image processing**, **bioinformatics**  $\square$ , and **market segmentation**, among many other fields.

Unlike supervised learning  $\mathfrak{G}$ , where the data is labeled, clustering works without labels and instead finds hidden patterns or structures in the data  $\square$ .

# **Q** Key Concepts in Clustering

- Clusters **6**: Groups of similar data points.
- Centroids **②**: The center or average of a cluster (especially in K-Means).
- **Distance Metrics** \( \Sigma : Methods like Euclidean distance or cosine similarity used to measure similarity between points.

# **Common Clustering Algorithms**

# 1. K-Means Clustering **(**

- Partitions data into *K* clusters.
- Simple and fast but requires you to choose *K* beforehand.

# 2. Hierarchical Clustering 🗟

- o Builds a tree-like structure (dendrogram).
- o No need to specify the number of clusters upfront.

#### 3. **DBSCAN** □

o Density-based method that finds clusters of arbitrary shapes.

o Great for noise and outlier detection **A**.

# 4. Gaussian Mixture Models (GMM)

- o Uses probability distributions to assign points to clusters.
- o Can model more complex data distributions.

## **Applications of Clustering**

- **Customer segmentation** in marketing.
- **Ext or document clustering** for organizing data.
- **Image segmentation** in computer vision.
- **N** Anomaly detection in fraud and cybersecurity.
- **Gene/protein grouping** in bioinformatics.

# Partitioning of data

## **©** What is Partitioning in Clustering?

Partitioning is a **clustering strategy** where a dataset is **divided into non-overlapping subsets** (clusters), with the goal of:

- Grouping similar data points together in the **same cluster** □.
- Ensuring that different clusters are **distinct** from each other □.
- Each data point belongs to exactly one cluster there is no overlap or ambiguity.

This is in contrast to other clustering methods like hierarchical clustering or density-based clustering, which can allow nested or arbitrarily shaped clusters.

### ☐ How Does Partitioning Work? (Detailed Steps)

Here's how a typical partitioning algorithm like **K-Means** works:

#### **Step 1: Choose the Number of Clusters (K)**

You start by deciding how many clusters (**K**) you want the data to be divided into. This is often based on prior knowledge, heuristics (like the **elbow method**), or trial-and-error.

#### **Step 2: Initialize Cluster Centers**

You randomly select **K** data points to act as the initial cluster centers (also called centroids in K-Means).

#### **Step 3: Assign Points to the Nearest Cluster**

Each data point is assigned to the cluster whose center is **closest** to it (based on a **distance metric** like Euclidean distance  $\diamondsuit$ ).

#### **Step 4: Update Cluster Centers**

Once all points are assigned, the algorithm calculates the **new centroid** of each cluster — usually the average of all data points in that cluster.

### **Step 5: Repeat Until Convergence**

Steps 3 and 4 are repeated until either:

- The cluster assignments no longer change, or
- A maximum number of iterations is reached.

# 

Let's say we have a dataset of customers, and we want to segment them based on their annual income and spending score.

We choose **K** = **3**. K-Means will:

- 1. Randomly pick 3 customers as the starting cluster centers.
- 2. Assign every customer to one of the 3 clusters based on how close they are to each center.
- 3. Recalculate each center as the average location of all customers in that cluster.
- 4. Repeat until the groupings stabilize.

Eventually, you get 3 distinct customer segments like:

- High income, high spenders.
- High income, low spenders.

# ≅ Partitioning vs Other Clustering Methods

Feature	Partitioning Clustering	Hierarchical Clustering	Density-Based Clustering (e.g., DBSCAN)
Number of clusters	Must be specified (K) □	Not needed	Not needed
Shape of clusters	Spherical/round	Arbitrary	Arbitrary
Handles noise/outliers	Poorly 🕱	Poorly	Very well ✔
Example Algorithms	K-Means, K- Medoids	Agglomerative, Divisive	DBSCAN, OPTICS

## Advantages of Partitioning

- Simple and computationally efficient.
- Works well with large datasets.
- Z Easy to interpret and visualize.

# 

- X You must know or estimate K ahead of time.
- X Sensitive to **initial positions** of centroids.
- X Assumes clusters are **spherical and equal-sized** (in K-Means).
- X Doesn't handle non-globular shapes or noise well.

## **X** Use Cases

- Marketing: Customer segmentation by behavior.
- **Healthcare**: Grouping patients with similar symptoms.
- Finance: Categorizing transaction patterns.
- Retail: Recommending products based on shopper clusters.

### Matrix factorization

### □ Matrix Factorization — A Complete Explanation

**Matrix Factorization** is a technique in linear algebra and machine learning used to break down a **large matrix into smaller components** that are easier to analyze and interpret. It's a key method for **data compression**, **pattern discovery**, **dimensionality reduction**, and especially **recommendation systems**.

### Core Idea

Given a large matrix **R**, matrix factorization aims to approximate it as the product of two or more smaller matrices:

$$R pprox U imes V^T$$

#### Where:

- **R** is the original matrix (e.g., user-item ratings).
- **U** is the **user-feature matrix** (latent user preferences).
- **V** is the **item-feature matrix** (latent item attributes).
- The product U×VTU \times V^TU×VT reconstructs the original matrix as closely as possible.

# Real-Life Example: Movie Recommendation System

#### Scenario:

You have a matrix **R** with:

- Rows = Users
- Columns = Movies
- Entries = Ratings (1–5), with many missing values

	Movie-1	Movie-2	Movie-3
User-1	5	?	3
User-2	4	2	?
User-3	?	5	4

Matrix Factorization fills in the blanks by learning:

- What kind of movies each user likes (e.g., comedy, action)
- · What kind of genre each movie represents

These are **latent features** not explicitly given in the data.

## Common Types of Matrix Factorization

Method	Description	
SVD (Singular Value Decomposition)	Decomposes into orthogonal components (U, $\Sigma$ , V <sup>t</sup> )	
NMF (Non-negative Matrix Factorization)	Ensures all matrix values remain positive (good for interpretability)	
PCA (Principal Component Analysis)	Finds directions (components) that capture maximum variance	
ALS (Alternating Least Squares)	Used in large-scale recommender systems	

# Mathematically (SVD):

Let's say R is an m× n matrix.

Domain

Then:

$$R = U \Sigma V^T$$

- **U**: m×k orthogonal matrix (left singular vectors)
- Σ: k×k diagonal matrix (singular values)
- **V**: n×k orthogonal matrix (right singular vectors)

In practice, we **truncate** the matrices to retain only the top k components (low-rank approximation).

**Use Case** 

 000 00.00
Recommender systems (Amazon, Netflix)

Domain	Use Case
<b></b> NLP	Topic modeling (Latent Semantic Analysis)
Computer vision	Image compression, face recognition
☐ Bioinformatics	Gene expression pattern analysis
■ Data Science	Dimensionality reduction for visualization

## ∧ Advantages & Disadvantages

## Advantages:

- Reveals hidden (latent) structures
- Helps handle sparse data (e.g., missing values)
- · Improves performance of ML algorithms
- Reduces dimensionality

### **X** Disadvantages:

- May require tuning parameters (like rank k)
- Not interpretable if you use SVD (better with NMF)
- Sensitive to noise or outliers

# Clustering of patterns

### ☐ What is Clustering of Patterns?

Clustering is the **unsupervised learning** process where we automatically discover natural groupings in data **without using labeled training data**. It's called "unsupervised" because we don't tell the algorithm what the clusters should be — it finds them on its own based on pattern similarities.

# Real-life Analogy:

Imagine you're sorting a box of mixed buttons  $\square$  by color and size, without anyone telling you what categories to use. You group them based on what *looks* similar — that's clustering!

# **௸ Why Do We Cluster?**

#### Clustering helps to:

- Discover hidden structures
- Summarize data
- Reduce dimensionality □
- Preprocess for other machine learning tasks (like classification)

### ☐ Deep Dive into Common Algorithms:

## 

- Objective: Divide data into **K groups** so each point belongs to the group with the nearest mean (centroid).
- Process:
  - 1. Randomly pick K centers.
  - 2. Assign points to the nearest center.
  - 3. Recalculate the centers based on the new assignments.
  - 4. Repeat until stable.
- Best For: Well-separated, spherical clusters.
- \( \frac{\lambda}{\text{Limitation: You need to decide **K** beforehand.
- **Example**: Grouping customers by purchasing behavior.

# 2. Hierarchical Clustering 🗟

- Builds a hierarchy of clusters.
- Two types:
  - Agglomerative (bottom-up): Start with each point as a cluster and merge them step-by-step.
  - Divisive (top-down): Start with one cluster and divide it.
- No need to specify the number of clusters!
- **Visual Output**: Dendrogram a tree structure showing merges/splits.
- **Example**: Grouping books by themes and sub-themes.

### 3. DBSCAN (Density-Based Spatial Clustering of Applications with Noise)



- Groups points that are **densely packed** together.
- Good for arbitrarily shaped clusters.
- Can identify noise/outliers.
- Parameters:
  - ε (epsilon): radius of neighborhood
  - MinPts: minimum points in a neighborhood to form a cluster
- **Example**: Detecting traffic congestion zones or fraud transactions.

## 4. Gaussian Mixture Models (GMM) 🚱

- Assumes data points are generated from multiple Gaussian distributions (bell-shaped curves).
- Assigns probabilities to data points for each cluster soft clustering.
- Uses Expectation-Maximization (EM) algorithm.
- **A Example**: Voice pattern analysis or speaker identification.

## P Distance & Similarity Measures:

- Euclidean Distance \( \script{\circ} : Straight-line distance in space. \)
- Manhattan Distance : Distance along axes (like city blocks).
- Cosine Similarity : Angle between vectors, great for text data.
- Jaccard Index ♥: Overlap between sets, useful for binary data.

# Applications Across Fields:

- **Marketing** : Segmenting customers for targeted ads.
- **Healthcare** Grouping patients based on symptoms or genetics.
- **Astronomy (**: Classifying galaxies or star clusters.
- **Education** : Clustering students based on performance and learning style.

## **M** Key Challenges:

- **② Choosing the right number of clusters** (especially in K-Means).
- | Handling high-dimensional data (curse of dimensionality).
- Noisy data & outliers can mislead clusters.
- M Interpretability: Clusters might not have meaningful labels.
- **Scalability**: Performance issues with big data.

## **m** Bonus: Evaluating Clusters

After clustering, you need to assess how "good" your clusters are. Common metrics:

- **Silhouette Score** (ranges from -1 to 1): Measures how similar an object is to its own cluster compared to others.
- Davies-Bouldin Index: Lower is better measures intra- vs inter-cluster distance.
- **Inertia** (K-Means): Sum of squared distances from points to their cluster centers.

# Divisive clustering

Divisive Clustering (Top-Down Hierarchical Clustering)

**Divisive clustering** is a **top-down** approach to hierarchical clustering. Unlike agglomerative clustering, which starts with individual data points and merges them, divisive clustering begins with **all data points in one big cluster** and **recursively splits** them into smaller sub-clusters until certain stopping conditions are met.

### ☐ Step-by-Step Workflow:

- 1. Start with all data in one cluster.
- 2. Select a way to split the cluster into two:
  - Use distance-based techniques (e.g., K-means, spectral clustering).
- 3. **Measure intra-cluster dissimilarity** (how different the members are).
- 4. **Split the cluster** to maximize the **between-cluster distance** and minimize within-cluster distance.

- 5. **Pick the next cluster** to divide (often the one with highest internal dissimilarity).
- 6. **Repeat** steps 2–5 until:
  - o All points are in their own clusters, OR
  - o A certain number of clusters is reached, OR
  - o A threshold for cluster quality is met.

### **Example:**

Let's say you have customer data with the following features:

- Total purchase amount (§)
- Frequency of visits
- Product categories bought

#### **Divisive Process:**

- 1. Start with all customers in one group.
- 2. Use K-means (K=2) to split into:
  - Group A: High purchase, frequent visits.
  - o Group B: Low purchase, infrequent visits.
- 3. Now take Group A and split again into:
  - Group A1: Electronics buyers
  - Group A2: Fashion buyers
- 4. Continue until groups are small and meaningful.

# **Q** Common Splitting Methods:

- 1. K-Means:
  - Often used at each step to split clusters into 2.
  - o Fast and effective but assumes spherical clusters.
- 2. PDDP (Principal Direction Divisive Partitioning):
  - Uses PCA to find the main axis of data spread.
  - Splits data along this axis.

### 3. Spectral Clustering:

- Constructs a similarity graph and splits based on graph partitioning techniques.
- Effective for non-convex shapes and complex structures.

## **A** Characteristics of Divisive Clustering:

Feature	Description
Approach	Top-down (starts with one cluster)
Result	A dendrogram (tree) showing how data is split
Output	Set of nested clusters
Stopping Criteria	Predefined number of clusters, distance threshold, or statistical cutoff
Cluster Shape	Depends on splitting technique (K-means = spherical, spectral = flexible)
Performance	Often <b>slower</b> than agglomerative due to repeated global operations

### ☐ When to Use Divisive Clustering:

### ✓ Use it when:

- You want a **global-to-local** understanding of data.
- You prefer to break a dataset into meaningful top-level categories first.
- You have enough computational resources (it can be resource-intensive).

#### 

- You need a fast, scalable solution for huge datasets.
- Your data is very noisy or contains many outliers divisive methods may be sensitive to these.

# **Ⅲ** Real-World Applications:

Domain Example

**Healthcare** Clustering patient records to identify disease subtypes

E-commerce Segmenting customers for personalized marketing

Cybersecurity 

Breaking down network traffic into normal vs. suspicious clusters

**Text Analysis** ≅ Separating large document sets into topics, then subtopics

```
[All Data]
/ \
[Cluster1] [Cluster2]
/ \
```

[Subcluster1] [Subcluster2]

### X Techniques Used:

- **Spectral Clustering**  $\square$ : Often used in divisive methods; uses eigenvectors of similarity matrices to perform splitting.
- **K-Means-based Split** : Sometimes K-Means is applied to split the data into two clusters at each step.
- **Principal Direction Divisive Partitioning (PDDP)**: Uses principal component analysis (PCA) to project data and find optimal split.

# Advantages:

- Captures a global view first, then drills down.
- Can be more accurate in some cases, as it considers the entire dataset's structure before splitting.
- Flexible stopping criteria.

# 

- Computationally expensive 😓 it must evaluate all possible splits.
- Not as popular or well-optimized as agglomerative clustering.

Sensitive to the method used for splitting.

### Real-World Example:

Let's say you're analyzing customer behavior:

- Start with all customers in one group.
- Split them into two major types (e.g., frequent buyers vs. infrequent buyers).
- Further split frequent buyers into loyal customers and deal hunters, and so on.

#### □ Applications:

- 🗁 Document classification
- ☐ Gene expression analysis
- ¶ Intrusion detection systems
- Market segmentation

# Agglomerative clustering

# **Q** What is Agglomerative Clustering?

Agglomerative clustering is the **most common form of hierarchical clustering**, where the goal is to build a **hierarchy of clusters**.

- It's a bottom-up approach.
- Starts with each data point in its own cluster.
- Repeatedly merges the most similar pair of clusters.
- The process continues until **only one cluster** remains (or a defined number is reached).

#### ■ Why Use It?

Agglomerative clustering is useful when:

- You want a multi-level view of how data naturally groups.
- You're unsure of the number of clusters and want flexibility.

 You care about interpretability — dendrograms help visualize merging decisions.

# How Does It Determine Which Clusters to Merge?

At each step, the algorithm merges the **two closest clusters**. But how is "closest" defined?

## **►** Distance (Proximity) Measures:

- Euclidean Distance: Straight-line distance (default in most tools).
- Manhattan Distance: Sum of absolute differences (useful for grid-like data).
- Cosine Distance: Measures the angle between vectors (great for text or highdimensional data).
- **Correlation-based**: How related the patterns are (used in bioinformatics or time-series).

# Linkage Criteria (how to measure distance between clusters):

Linkage Type	Description	Behavior
Single Linkage	Distance between <b>closest pair</b> of points in two clusters	Forms long, chain-like clusters
Complete Linkage	Distance between <b>farthest pair</b> of points	Forms compact, tight clusters
Average Linkage	Average distance between all pairs across clusters	Balanced
Ward's Method	Increases in variance when clusters are merged	Optimizes internal compactness

# ⚠ Mathematics Behind It (Ward's Method Example)

Ward's linkage minimizes the total within-cluster variance:

• At each merge, it looks for the pair of clusters that **causes the smallest** increase in total variance.

$$\sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

Objective: Minimize

### ☐ Real-Life Applications

Field	Use Case
Marketing	Grouping customers based on purchasing behavior
Genomics □	Clustering gene expression data
Psychology 🏻	Grouping survey responses to identify personality types
Text Analysis 僵	Clustering articles or social media posts by theme
Finance 💸	Grouping stocks by behavior or risk

# Pros and Cons

## ✓ Pros:

- No need to specify K (can extract clusters by cutting the dendrogram).
- **Deterministic** (same input → same output).
- Provides rich insights via hierarchy.
- Works well with **non-spherical clusters** if proper linkage is chosen.

#### X Cons:

- Scales poorly with large datasets.
- Greedy merging: can't undo bad merges.
- Sensitive to:
  - Distance measure.
  - Noise/outliers.
  - o Choice of linkage method.

# **VS** Agglomerative vs Divisive Clustering

Feature	Agglomerative	Divisive
Approach	Bottom-up	Top-down
Start With	Each point as its own cluster	All points in one cluster
Merge/Split Process	Merge closest clusters	Split most dissimilar cluster
Common Usage	More commonly used and supported	Less common, more computationally expensive
Flexibility	Hard to revise once merged	Can avoid bad merges with better global view

# Partitional clustering

# **6** What is Partitional Clustering?

Partitional clustering is a non-hierarchical method that:

- Divides the data into non-overlapping subsets (clusters),
- Each data point belongs to exactly one cluster,
- And the clusters are created simultaneously not incrementally (like hierarchical methods).

Unlike hierarchical clustering (which builds a tree), partitional clustering just splits the dataset into a flat set of clusters.

# How Partitional Clustering Works

- 1. Select the number of clusters (K) you want to find.
- 2. **Initialize cluster centroids** randomly or using some method (like K-means++ for better starting points).
- 3. Assign each data point to the nearest cluster center.
- 4. **Update the centroids** based on the current membership.
- 5. **Repeat** steps 3–4 until convergence (i.e., no more changes in assignments or centroids).

# Most Common Partitional Algorithm: K-Means Clustering

### ☐ K-Means Steps:

- 1. Choose K, the number of clusters.
- 2. Initialize K centroids randomly.
- 3. For each point:
  - Assign it to the nearest centroid (Euclidean distance).
- 4. For each cluster:
  - Recalculate the centroid as the mean of all points assigned to it.
- 5. Repeat until:
  - o Centroids stop changing, or
  - Max iterations reached.

### □ Objective:

Minimize the Within-Cluster Sum of Squares (WCSS):

$$\sum_{k=1}^{K} \sum_{x_i \in C_k} \|x_i - \mu_k\|^2$$

# **?** Other Partitional Algorithms

Algorithm Key Feature

**K-Means** Fast, uses means to represent clusters

K-Medoids(PAM)Uses actual data points as centers (medoids)

**CLARANS** Sampling-based improvement over K-Medoids

Fuzzy C-Means

Soft clustering: points can belong to multiple clusters with a degree of membership

#### □ Characteristics of Partitional Clustering

Feature Description

**Number of clusters** Must be specified in advance

Feature Description

Shape of clusters Works best with spherical, equally sized clusters

**Scalability** Very efficient for large datasets

**Interpretability** Simple, easy to understand and visualize

Cluster Assignment Each point is assigned to one cluster (hard clustering)

### **居 Example: Customer Segmentation**

Imagine you have a retail dataset:

• Features: Age, Income, Spending score

You apply K-means with K = 3:

- Cluster 1: Young, high spenders
- Cluster 2: Middle-aged, average spenders
- Cluster 3: Older, low spenders 

   □

Each customer belongs to one of these groups — that's partitional clustering!

# Advantages

- Fast and scalable for large datasets.
- Works well when clusters are well-separated and convex.
- | Simple to implement and understand.

### **X** Disadvantages

- ↓ Requires you to know K in advance.
- Doesn't work well with non-convex or varying-size clusters.
- Sensitive to:
  - Outliers (can skew centroids),
  - Initialization (random starting points can affect results).

### □ Applications of Partitional Clustering

Field	Use Case
E-commerce 🖺	Customer segmentation
Healthcare 🔡 🤇	Grouping patients by symptoms or diagnosis
Image Processing	Color quantization (reducing colors)
Biology □	Grouping organisms by genetic traits
Social Media 🖺	Community detection in user behavior

### □ Summary Table

Feature	Agglomerative	Divisive	Partitional (e.g., K-means)
Approach	Bottom-up	Top-down	Simultaneous partitioning
Structure	Hierarchical (tree)	Hierarchical (tree)	Flat (non-hierarchical)
Output	Dendrogram	Dendrogram	Cluster assignments only
Flexibility Ca	an explore many levels	Same	Must define K in advance
Efficiency	Slower	Slowest	Fastest

# k-means clustering

# **6** What is K-Means Clustering?

**K-Means** is an **unsupervised machine learning algorithm** that divides a dataset into **K distinct**, **non-overlapping clusters**. Each cluster is defined by its **centroid** (the arithmetic mean of the points in that cluster).

☑ It works by **iteratively refining** the position of the centroids and the assignment of points to those centroids.

# Core Concepts

• **K**: The number of clusters you want to divide the data into.

- **Centroid**: The center point of a cluster (mean of all points in the cluster).
- Distance metric: Usually Euclidean distance is used to assign points to centroids.

### Step-by-Step Working of K-Means

Let's break it down step-by-step:

- 1. Initialize K centroids randomly from the data points.
- 2. Assign each point to the nearest centroid.
  - This forms K clusters.
- 3. Recalculate the centroids.
  - Take the mean of all data points in each cluster.
- 4. Repeat steps 2 and 3:
  - Until centroids don't change significantly (i.e., convergence),
  - Or a maximum number of iterations is reached.

# □ Mathematical Objective

 K-Means aims to minimize the intra-cluster variance (also called Within-Cluster Sum of Squares, WCSS):

$$\text{WCSS} = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|^2$$

Where:

- Ck = Cluster k
- µk= Centroid of cluster k
- x<sub>i</sub> = Data point in cluster C<sub>k</sub>

# ☑ Visualization Example

Imagine plotting 2D data points on a graph:

K-means starts with random cluster centers,

- Points are **grouped** around the nearest center,
- Centers are **updated** to the middle of the group,
- This repeats until clusters **stabilize**.

# ✓ Advantages of K-Means

Benefit	Description
← Fast & Scalable	Works efficiently on large datasets
□ Easy to Understand	Simple mathematical model
Minimizes Variance	Good at reducing intra-cluster spread
	in most ML libraries (scikit-learn, TensorFlow, etc.)

# **X** Disadvantages

Limitation	Description		
	You have to choose the number of clusters		
Assumes Spherical Clusters	Doesn't handle irregular shapes well		
⊠ Sensitive to Outliers	A few extreme values can distort centroids		
☐ Random Initialization	May lead to <b>local optima</b> (can use K-means++ to fix)		

## 

Domain Use Case

Retail Customer segmentation

☐ Bioinformatics Grouping gene expressions

Computer Vision Image compression, color quantization

Travel Grouping destinations by visitor interest

□ NLP Document clustering (if preprocessed into vectors)

## **⋈** Advanced Variations

Variant Description

K-Means++ Smarter initialization

MiniBatch K-

Means

Faster for very large datasets (uses random mini-batches)

K-Medoids Uses medoids (actual data points) instead of mean — more

robust to outliers

Fuzzy C-Means Each point belongs to multiple clusters with probabilities

# Summary: Strengths & Weaknesses

Aspect	Pros 🗸	Cons 🗙
Speed	Fast, especially for large data	Can still be slow with very large K
Accuracy	Often good when assumptions hold	Poor if data is noisy or non- spherical
Simplicity	Easy to implement and interpret	Requires specifying K
Cluster Shapes	Works well with spherical clusters	Not ideal for complex shapes

Aspect	Pros 🗸	Cons X
Outlier	Sensitive to extreme values	Can be improved with other
Handling	Sensitive to extreme values	methods

# Soft partitioning

### **Soft Partitioning (Fuzzy Clustering)**

Soft partitioning is an approach to clustering where each data point does **not** belong to just **one cluster**. Instead, it can belong to multiple clusters with a degree of membership or **fuzziness**. This is in contrast to **hard partitioning** (like K-Means), where each data point is strictly assigned to a single cluster.

The most common algorithm for soft partitioning is Fuzzy C-Means (FCM).

### Mathematical Foundation of Fuzzy C-Means (FCM)

### **Objective Function (Revisited)**

Fuzzy C-Means clustering aims to minimize the **objective function** J(U,V) which is a weighted sum of squared distances between the data points and the centroids. The goal is to minimize this function to find the **optimal centroids** and **membership** values.

$$J(U, V) = \sum_{i=1}^{n} \sum_{k=1}^{K} u_{ik}^{m} \|x_{i} - v_{k}\|^{2}$$

#### Where:

- U<sub>ik</sub>: Membership degree of point x i for cluster k.
- m: Fuzziness exponent (typically m=2).
- V<sub>k</sub>: Centroid of cluster k.
- X<sub>i</sub>: Data point i.
- ||x<sub>i</sub>-v<sub>k</sub>||<sup>2</sup>: Euclidean distance between point x<sub>i</sub> and centroid v<sub>k</sub>.

## Membership Calculation (Fuzzy)

To calculate the membership values, FCM uses the following formula, which is based on the relative distances between data points and centroids:

$$u_{ik} = \frac{1}{\sum_{j=1}^{K} \left(\frac{\left\|x_i - v_k\right\|}{\left\|x_i - v_j\right\|}\right)^{\frac{2}{m-1}}}$$

This equation is designed to make the membership values depend on the **distance** between data points and centroids:

- If point  $x_i$  is closer to centroid  $v_k$ , the membership  $u_{ik}$  will be closer to 1.
- If point x<sub>i</sub> is far from v<sub>k</sub>, the membership u<sub>ik</sub> will be closer to 0.

# ✓ Advantages of Soft Partitioning (Fuzzy Clustering)

**♀** Advantage

**Explanation** 

1. Flexible Membership

Each data point can belong to multiple clusters with varying degrees, which reflects real-world ambiguity better than hard clustering.

2. Good for Overlapping Data Ideal when clusters are not clearly separable and overlap significantly (e.g., human behavior, medical symptoms).

3. Robust to Boundary Uncertainty

Handles borderline or ambiguous data points more naturally by assigning fractional membership.

4. Realistic in Complex Systems

More representative in systems where objects naturally belong to more than one group (e.g., a hybrid car in vehicle clustering).

5. Smooth Transitions
Between Clusters

Unlike hard assignments, the transitions between clusters are gradual, making the results more nuanced.

6. Better Interpretability in Some Contexts

Membership values can provide additional insights into the level of belonging or likelihood of association.

7. Useful in Decision Support Systems

Helpful in systems requiring probabilistic or fuzzy decision making (e.g., expert systems, medical diagnosis).

# ★ Disadvantages of Soft Partitioning (Fuzzy Clustering)

# 

# **Explanation**

1. More Computationally Intensive

Requires more memory and time compared to hard clustering due to extra membership calculations.

2. Sensitive to Initialization

Like K-Means, results can vary depending on the initial membership values or centroids.

3. Parameter Sensitivity (Fuzziness 'm')

Performance depends heavily on the fuzziness parameter mmm; choosing the wrong value can lead to poor results.

4. Less Intuitive Interpretation

While flexible, interpreting degrees of membership can be harder than discrete labels.

5. Struggles with Outliers

FCM is sensitive to outliers and noise unless modified (like with possibilistic clustering).

6. Assumes Spherical Clusters

Similar to K-Means, FCM works best when clusters are roughly spherical and similar in size.

Minima

7. Convergence to Local May not find the global optimum and can converge to a poor local minimum if not carefully initialized.

# Fuzzy c-means clustering

# **Q** What is Fuzzy C-Means (FCM)?

Fuzzy C-Means is a **clustering algorithm** that allows each data point to **partially** belong to multiple clusters. This contrasts with traditional clustering methods (like K-Means) where each point can belong to **only one cluster**.

#### Think of it like this:

If you're classifying colors, the color "turquoise" might be:

- 60% blue
- 40% green

Traditional clustering would force "turquoise" into either blue or green — but Fuzzy C-Means allows it to be both, with degrees.

Let's break it into steps with more intuitive explanation:

#### 1. Initialization

- Choose:
  - o c: number of clusters.
  - o m: fuzziness parameter (commonly m=2, must be > 1).
- Randomly assign a membership matrix U:
  - o Each element u ik defines how much data point xi belongs to cluster k.
  - o For each point x<sub>i</sub>, the total membership across all clusters is 1:

$$\sum_{k=1}^{c} u_{ik} = 1$$

### 2. Compute Centroids

 For each cluster k<sub>i</sub>, calculate its center (centroid) using all data points, weighted by their membership:

$$v_k = \frac{\sum_{i=1}^{n} u_{ik}^m \cdot x_i}{\sum_{i=1}^{n} u_{ik}^m}$$

•

So, if a point strongly belongs to a cluster (i.e., high u<sub>ik</sub>, it has a
greater influence on the centroid.

### 3. Update Membership Matrix

• Update uik based on how close xi is to each cluster center:

$$u_{ik} = \frac{1}{\sum_{j=1}^{c} \left(\frac{\left\|x_{i} - v_{k}\right\|}{\left\|x_{i} - v_{j}\right\|}\right)^{\frac{2}{m-1}}}$$

- The closer a point is to a centroid, the **higher** its membership to that cluster.
- The farther it is, the lower the membership.

### 4. Repeat Steps 2 and 3 Until Convergence

• Stop when the **change in membership values** is less than a small threshold (e.g.,  $\epsilon$ =0.001\epsilon = 0.001 $\epsilon$ =0.001).

• This means the algorithm has found a **stable configuration** of clusters.

	Why	Use	FCM	Instead	of	K-Means?
--	-----	-----	-----	---------	----	----------

Feature	K-Means	Fuzzy C-Means
Point belongs to	One cluster only	Multiple clusters
Result type	Hard labels	Membership scores
Good for	Well-separated clusters	Overlapping clusters
Handles uncertainty	Poorly	Well
Example use case	Clustering clear categories	Fuzzy diagnoses, image segmentation

# Real-World Analogy

Imagine a student taking multiple courses:

- In **hard clustering**, you'd assign the student to **one major** e.g., Computer Science.
- In **fuzzy clustering**, you acknowledge the student might be:
  - o 70% Computer Science
  - o 30% Math

This reflects **real human behavior** more accurately — we often don't fit neatly into one category.

# **%** Advantages of FCM (Explained)

## 1. Captures Real-World Ambiguity

Many datasets have overlapping patterns — fuzzy logic handles that elegantly.

### 2. More Informative

Instead of just "Cluster A," you get **how strongly** a point belongs to each cluster.

## 3. **Better for Noisy Data**

Outliers don't dominate a single cluster — they get **low memberships** across all clusters.

## 4. Versatile Applications

Used in medical diagnosis, image processing, text mining, etc.

## 

#### 1. X Slower than K-Means

More complex math means more computational time.

### 2. X Sensitive to Outliers

Very distant points still affect centroids unless additional logic is added.

#### 3. X Parameter Sensitivity

Choosing the right number of clusters (c) and fuzziness (m) is critical.

## Rough clustering

### ■ What Is Rough Clustering (In Depth)?

Rough Clustering is a soft-computing technique rooted in Rough Set Theory (RST) — developed by Zdzisław Pawlak. It's a way to deal with vagueness, uncertainty, and imprecise data by forming clusters with boundaries.

Instead of assigning a data point **definitely** or **probabilistically** to a single cluster (like K-Means or Fuzzy C-Means), **Rough Clustering uses two key approximations**:

## **Q** Key Concepts

Term Meaning

Lower
Approximation
The set of objects that definitely belong to the cluster

Upper
Approximation
The set of objects that possibly belong to the cluster

**Boundary Region** The difference between upper and lower approximations (i.e., uncertain region)

So each cluster is defined as:

Cluster = (lower approximation, upper approximation)

# Core Definitions

Term	Symbol	Meaning	
Lower Approximation	<u>C</u>	Data points that definitely belong to cluster C. No ambiguity.	
Upper Approximation	C_	Data points that possibly belong to cluster C.	
Boundary Region	C- <u>C</u>	Area of <b>uncertainty</b> , where it's unclear if the data point belongs to C.	
		$C = (\underline{C}, \overline{C})$	

So instead of just "belongs" or "doesn't belong," a data point can be:

- Clearly inside
- Clearly outside
- Possibly inside (boundary region)

## How Does Rough Clustering Work? (Step-by-Step)

Here's a more detailed algorithm outline:

#### 1. Data Representation

You start with a dataset that may contain **incomplete or uncertain values** (e.g., missing entries, noisy measurements).

#### 2. Define Similarity or Distance Function

You measure how **similar** or **different** each data point is to others — using a function like Euclidean distance, cosine similarity, or even decision table-based methods.

#### 3. Set a Threshold

Choose a threshold value that determines:

- If a point is **certainly in** a cluster (low distance → lower approximation)
- If it is **possibly in** a cluster (moderate distance → upper approximation)

This threshold is often called the **indiscernibility threshold**.

#### 4. Assign Approximations

- Assign a point to the **lower approximation** of a cluster if it meets the strong similarity criterion.
- Assign it to the **upper approximation** if it meets a relaxed criterion.
- If it doesn't meet either it's **excluded**.

#### 5. Update Clusters

Iteratively update the approximation sets and recompute similarity measures until convergence.

## Comparison with Other Clustering Techniques

Feature	K-Means	Fuzzy C-Means	Rough Clustering
Data Belonging Type	Hard	Soft (probability)	Rough (set-based)
Belonging Output	One cluster	Degrees of membership	Two approximations
Uncertainty Representation	<b>X</b> No	Yes (probabilistic)	Yes (explicit boundary)
Handles Missing/Noisy Data	× Poorly		✓ Well (core strength)
Basis of Logic	Distance	Fuzzy logic	Rough set theory

# Real-World Analogy

Imagine classifying animals:

- Lion: Clearly in the category "Big Cats" → lower approximation
- Jaguar: Maybe a Big Cat, but also overlaps with "Tree Dwellers" → boundary region
- **Dog**: Not a Big Cat → excluded

Rough clustering allows you to capture that ambiguity.

# **✓** When To Use Rough Clustering

#### Situation

#### Why Rough Clustering Works

Uncertain or noisy datasets

It doesn't force hard labels — gives space for

ambiguity

**Medical diagnostics** 

Patients might show symptoms overlapping multiple

conditions

Social science or behavior data

Human behavior is often fuzzy and not clearly separable

Text/document classification

Some documents may belong to more than one topic

# Rough k-means clustering algorithm

The **Rough K-Means Clustering Algorithm** is a hybrid clustering technique that combines concepts from **Rough Set Theory** and **K-Means Clustering**. It is particularly useful when dealing with **uncertain**, **imprecise**, **or vague data**. Here's a breakdown of what it is and how it works:

### ☐ Key Concepts

- Rough Set Theory (introduced by Zdzisław Pawlak): Handles vagueness by using lower and upper approximations of a set.
  - Lower Approximation: Objects that definitely belong to a cluster.
  - Upper Approximation: Objects that possibly belong to a cluster.
- K-Means Algorithm: A classical clustering technique that partitions data into k clusters by minimizing the distance between points and their cluster centroids.

# Boundary=Upper Approximation-Lower Approximation

# Rough K-Means Clustering: How It Works

The Rough K-Means algorithm modifies the traditional K-Means algorithm by incorporating lower and upper approximations for each cluster.

# Algorithm Steps

- 1. Initialization:
  - Choose the number of clusters k.
  - Initialize centroids randomly.
  - o Assign data points to **lower**, **upper**, or **boundary** regions of clusters.

### 2. Assignment Step:

- For each data point, calculate its distance to all cluster centroids.
- o Assign each point:
  - To lower approximation if it is clearly closest to one cluster.
  - To upper approximation (i.e., boundary region) if it is almost equally close to two or more clusters.

### 3. Centroid Update Step:

Update cluster centroids using a weighted combination:

$$C_i = \frac{w_l \sum_{x \in L_i} x + w_u \sum_{x \in B_i} x}{w_l |L_i| + w_u |B_i|}$$

- L<sub>i</sub>: Lower approximation of cluster i
- B<sub>i</sub>: Boundary region of cluster i
- $W_1, w_u$ : Weights for lower and upper approximations (typically  $w_1>w_u$ )

#### 4. Repeat:

 Reassign points and update centroids until convergence (e.g., centroids stop changing).

#### 5. Stopping Criterion:

 Algorithm stops when the assignment of points or centroid positions stabilizes.

# Advantages

- Handles **imprecise or overlapping clusters** better than standard K-means.
- Suitable for noisy or vague datasets.
- Incorporates uncertainty in cluster memberships.

### **X** Disadvantages

- Slightly more complex than K-Means.
- · Requires choosing weights wlw lwl and wuw uwu appropriately.
- Computationally heavier due to dual approximations.

### Applications

- Medical diagnosis
- Image segmentation
- Market segmentation
- · Text categorization with ambiguity

# Expectation maximization – based clustering

### **Expectation-Maximization (EM) Based Clustering – In Depth**

**Expectation-Maximization (EM)** is a **probabilistic clustering algorithm** that is widely used when data is assumed to come from a mixture of different distributions — typically **Gaussian Mixture Models (GMMs)**. Unlike K-means, which assigns points to the nearest cluster **deterministically**, EM assigns each point to **clusters probabilistically**.

### ☐ Key Concepts

- Latent variables: Hidden (unobserved) variables that represent the cluster each data point belongs to.
- **Expectation-Maximization**: An iterative algorithm to find **maximum likelihood estimates** when data has hidden variables.
- **Gaussian Mixture Model (GMM)**: A weighted sum of multiple Gaussian distributions, each representing a cluster.

#### 2. Initialization

- Choose the number of clusters k.
- Initialize:
  - Means μ<sub>k</sub> randomly
  - Covariance matrices Σ<sub>k</sub>
  - $\circ$  Mixing coefficients  $\pi_k = 1/k$

#### 3. EM Iteration

#### A. E-Step (Expectation)

• For each data point  $x_i$ , calculate the **responsibility** (posterior probability) that cluster k generated  $x_i$ :

$$\gamma_{ik} = \frac{\pi_k \cdot \mathcal{N}(x_i \mid \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \cdot \mathcal{N}(x_i \mid \mu_j, \Sigma_j)}$$

- $\square$  N(x<sub>i</sub> |  $\mu_k$ ,  $\Sigma_k$ ): Gaussian probability density function

### B. M-Step (Maximization)

Update parameters based on the responsibilities:

Update mean:

$$\mu_k = \frac{\sum_{i=1}^N \gamma_{ik} x_i}{\sum_{i=1}^N \gamma_{ik}}$$

Update covariance:

$$\Sigma_k = \frac{\sum_{i=1}^{N} \gamma_{ik} (x_i - \mu_k) (x_i - \mu_k)^T}{\sum_{i=1}^{N} \gamma_{ik}}$$

· Update mixing coefficient:

$$\pi_k = \frac{1}{N} \sum_{i=1}^{N} \gamma_{ik}$$

The log-likelihood is:

$$\log L = \sum_{i=1}^{N} \log \left( \sum_{k=1}^{K} \pi_k \cdot \mathcal{N}(x_i | \mu_k, \Sigma_k) \right)$$

# Advantages

- Handles overlapping clusters well.
- More flexible clusters can be elliptical, not just spherical.
- **Probabilistic** more informative than hard assignments.
- Can handle missing data better via latent variables.

## **X** Disadvantages

- Sensitive to initialization (like K-means).
- Computationally heavier (matrix inversion in covariance).
- Can converge to **local maxima**.
- Choosing the number of clusters kkk is still required (though model selection like AIC/BIC can help).

## ☐ Real-World Applications

- Image segmentation
- Speaker recognition
- Document clustering
- **Medical diagnostics** (e.g., modeling patient data distributions)