# Unsupervised Learning: Clustering Algorithms

Clustering is a foundational task in unsupervised machine learning that focuses on discovering natural groupings within data. Unlike supervised learning, it operates without predefined labels. The primary goal is to partition data points into distinct groups, or clusters, where points within the same cluster are more similar to each other than to those in other clusters.

## 1. K-Means Clustering

K-Means is one of the most popular and straightforward clustering algorithms. It is a centroid-based method, which means it aims to find cluster centers (centroids) that represent the central point of different data regions.

#### **How it Works**

- 1. **Initialization**: First, you choose the number of clusters, \$K\$, and randomly initialize \$K\$ centroids. These points act as the initial centers for the clusters.
- 2. **Assignment Step**: Each data point is assigned to its nearest centroid, typically based on Euclidean distance. This process forms \$K\$ distinct clusters.
- 3. **Update Step**: The centroid of each cluster is recalculated by computing the mean of all data points assigned to it.
- 4. **Iteration**: The assignment and update steps are repeated until the centroids no longer shift significantly, indicating that the clusters have stabilized.

**Example**: Imagine plotting customer data by age and spending score. K-Means would identify \$K\$ central points (like "young high-spenders" or "old low-spenders") and group the surrounding customers around them.

Pros	Cons
Fast and scalable for large datasets.	You must specify the number of clusters (\$K\$) beforehand.
Simple to understand and implement.	It is sensitive to the initial placement of centroids.
	Struggles with non-spherical clusters or clusters of varying sizes and densities.

## 2. Hierarchical Clustering

This method constructs a tree-like hierarchy of clusters known as a **dendrogram**. A key advantage is that it does not require you to pre-specify the number of clusters.

There are two primary approaches:

- Agglomerative (Bottom-up): This is the more common method. It starts by treating each data point as its own cluster. In each step, it merges the two closest clusters until only a single cluster remains.
- **Divisive (Top-down)**: This method starts with all data points in one large cluster. It then recursively splits a cluster into two at each step until every data point is its own cluster.

#### **Linkage Criteria**

This determines how the distance between clusters is measured:

- **Single Linkage**: The distance between the closest points in the two clusters.
- **Complete Linkage**: The distance between the farthest points in the two clusters.
- Average Linkage: The average distance between all pairs of points across the two clusters
- Ward's Linkage: Merges clusters in a way that minimizes the increase in variance within the new merged cluster.

Pros	Cons
No need to specify the number of clusters upfront.	Computationally expensive, with a complexity of at least \$O(n^2)\$, making it difficult for large datasets.
The resulting dendrogram is highly informative and helps in understanding the data's structure.	Can be sensitive to noise and outliers.

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

DBSCAN is a density-based algorithm that groups closely packed points together. It is particularly effective at identifying outliers, which it marks as noise points lying in low-density regions.

#### **Key Concepts**

- **Epsilon (\$\epsilon\$)**: A specified distance radius.
- **Minimum Points (minPts)**: The minimum number of points required within the \$\epsilon\$ radius for a point to be considered a core point.
- Core Point: A point that has at least minPts within its \$\epsilon\$ radius (including itself).
- **Border Point**: A point within the \$\epsilon\$ radius of a core point that doesn't meet the minPts requirement itself.
- Noise Point (Outlier): Any point that is neither a core nor a border point.

#### **How it Works**

The algorithm starts with an arbitrary point and, if it's a core point, creates a new cluster. It

then finds all density-connected points and adds them to the cluster, repeating the process until every point has been visited.

Pros	Cons
Can find arbitrarily shaped clusters.	Struggles with clusters of varying densities.
Robust to outliers and explicitly identifies them.	Performance is dependent on the choice of \$\epsilon\$ and minPts.
Does not require specifying the number of clusters.	

## 4. OPTICS (Ordering Points To Identify Clustering Structure)

OPTICS is an extension of DBSCAN designed to overcome one of its main weaknesses: identifying clusters in data with varying densities.

#### **How it Works**

Rather than producing a simple cluster assignment, OPTICS generates an augmented ordering of the data points that contains information about the density structure. It calculates two key values for each point:

- Core Distance: The smallest \$\epsilon\$ that would make a point a core point.
- **Reachability Distance**: The distance from a point to its nearest core point.

From this, a "reachability plot" can be created, where valleys indicate dense clusters, allowing for cluster extraction with different density parameters.

Pros	Cons

Handles clusters of varying densities better than DBSCAN.	More complex to understand and implement than DBSCAN.
Provides a more informative view of the data's structure.	Does not produce a simple clustering partition directly; it requires an extra step to extract clusters from the reachability plot.

## 5. Mean-Shift Clustering

Mean-Shift is a mode-finding algorithm that doesn't require you to specify the number of clusters. Its goal is to locate the peaks (modes) in the data's density function.

#### **How it Works**

- 1. A window (kernel) is placed around each data point.
- 2. The mean of the points within that window (the "mean-shift") is calculated.
- 3. The window's center is shifted to this new mean.
- 4. Steps 2 and 3 are repeated until the window's position converges at a peak of data density.
- 5. All points that converge to the same peak are grouped into the same cluster.

Pros	Cons
Automatically determines the number of clusters.	The choice of bandwidth is critical and can significantly affect results.
Can find arbitrarily shaped clusters.	Can be computationally expensive for large datasets.
Its only parameter, bandwidth (window size), can often be estimated automatically.	

### 6. Gaussian Mixture Models (GMM)

GMM is a probabilistic model that assumes data points are generated from a mixture of several Gaussian distributions (bell curves) with unknown parameters. It is considered a "soft clustering" method.

#### **How it Works**

GMM provides a probability for each point belonging to each cluster, rather than a hard assignment. It uses the **Expectation-Maximization (EM)** algorithm to determine the parameters of the Gaussian distributions.

- Expectation (E-step): Estimates the probability of each point belonging to each cluster.
- Maximization (M-step): Updates the Gaussian parameters based on those probabilities. These steps are repeated until the model converges.

Pros	Cons
Soft clustering provides more information about the uncertainty of assignments.	You must specify the number of components (clusters).
Can model elliptically shaped clusters, unlike K-Means which assumes spherical clusters.	Can be computationally slow and may converge to a local minimum.

## 7. BIRCH Clustering

BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies) is specifically designed for very large datasets where memory is a constraint.

#### **How it Works**

- 1. **Build a CF Tree**: The algorithm performs a single pass over the data to build a compact in-memory summary called a Clustering Feature (CF) Tree. Each leaf node in this tree represents a small, dense sub-cluster.
- 2. **Global Clustering**: Another clustering algorithm (like K-Means) is then applied to the leaf nodes of the CF Tree, not the entire original dataset.

Pros	Cons
Very fast and memory-efficient for large datasets.	Only works with numerical data.
Processes the dataset in a single pass.	The resulting clusters may not be as accurate as those from slower algorithms.

## 8. Affinity Propagation

This algorithm uses a unique approach where data points "vote" for other points to become their representatives, known as "exemplars".

#### **How it Works**

All data points are simultaneously considered potential exemplars. The algorithm works by exchanging messages between points until a high-quality set of exemplars and clusters emerges.

- **Responsibility**: A message from point \$i\$ to candidate exemplar \$k\$ reflects how well-suited \$k\$ is to be the exemplar for \$i\$.
- **Availability**: A message from candidate exemplar \$k\$ to point \$i\$ reflects how appropriate it would be for \$i\$ to choose \$k\$ as its exemplar.

Pros	Cons
Does not require specifying the number of clusters.	High computational complexity (\$O(n^2)\$), making it unsuitable for large datasets.
The chosen exemplars are actual data points, which makes the clusters easy to interpret.	Can be sensitive to the "preference" parameter.

## 9. Spectral Clustering

Spectral Clustering uses the connectivity and graph structure of the data to perform clustering. It is exceptionally powerful for finding clusters with complex, non-convex shapes.

#### **How it Works**

- 1. **Create an Affinity Matrix**: An affinity matrix is constructed where each entry \$(i, j)\$ represents the similarity between points \$i\$ and \$j\$.
- 2. Compute the Laplacian Matrix: A graph Laplacian is derived from the affinity matrix.
- 3. **Eigen-decomposition**: The eigenvectors of the Laplacian matrix are calculated.
- 4. **Cluster in Lower Dimension**: A simpler algorithm like K-Means is used on these eigenvectors to partition the data.

The core idea is that the eigenvectors project the data into a lower-dimensional space where clusters become much more separable.

Pros	Cons
Excellent for finding non-convex clusters (e.g., concentric circles).	Requires you to specify the number of clusters.
Has a strong theoretical foundation in	Can be computationally expensive for large datasets due to the eigen-decomposition

graph theory.	step.

# Project: Segmenting Customers Based on Purchase Behavior

This project demonstrates using clustering to identify meaningful customer groups in a retail environment, enabling data-driven business decisions.

#### Goal

The main objective is to segment supermarket customers into distinct groups based on their purchasing habits. This allows the business to implement:

- Targeted Marketing: Create campaigns tailored to specific groups.
- Personalized Offers: Send relevant promotions to increase sales and loyalty.
- Better Inventory Planning: Stock products based on demand from key customer segments.

#### **Techniques to Use**

- **Clustering Algorithm**: K-Means is a great starting point for its simplicity and efficiency. DBSCAN is a good alternative if outliers or non-spherical clusters are suspected.
- Data Analysis: Exploratory Data Analysis (EDA) and visualization.

### Methodology & Steps

#### 1. Data Collection & Understanding

We'll define a hypothetical dataset at the customer level. Key features would include:

- CustomerID: Unique identifier for each customer.
- Recency: Days since the customer's last purchase.
- Frequency: Total number of transactions.
- MonetaryValue: Total amount spent.
- AvgBasketSize: Average number of items per transaction.
- ProductCategoryDiversity: Number of unique product categories purchased.

#### 2. Data Preprocessing

- **Handling Missing Values**: Decide on a strategy like imputation or removing the customer record.
- **Feature Scaling**: This is critical for distance-based algorithms like K-Means. Use StandardScaler or MinMaxScaler to ensure features with larger scales don't dominate the process.

#### 3. Exploratory Data Analysis (EDA)

- Visualize feature distributions with histograms and box plots to find skewness and outliers.
- Use scatter plots to observe initial relationships between feature pairs (e.g., Frequency vs. MonetaryValue).

#### 4. Determining the Optimal Number of Clusters (K)

For K-Means, \$K\$ must be chosen. Two common methods are:

- **The Elbow Method**: Plot the within-cluster sum of squares (WCSS) for different values of \$K\$. The "elbow" point, where the rate of decrease slows, is a good estimate for \$K\$.
- **Silhouette Score**: Measures how similar a point is to its own cluster versus others. A score closer to 1 indicates better-defined clusters.

#### 5. Applying the K-Means Algorithm

- Initialize and train the K-Means model with the chosen \$K\$ on the preprocessed data.
- Assign each customer to one of the resulting clusters.

#### 6. Analyzing and Profiling the Clusters

Analyze the average feature values for each cluster to understand its characteristics. Potential segments could include:

#### • Cluster 0: Loyal Champions

- **Profile**: High Frequency, high MonetaryValue, low Recency. They are your best customers who shop often, spend a lot, and visited recently.
- Action: Reward them with a loyalty program, exclusive offers, and early access to new products.

#### • Cluster 1: Potential Loyalists

o **Profile**: Moderate Frequency, moderate Monetary Value, low Recency. They are

recent shoppers with the potential to become champions.

 Action: Engage them with personalized recommendations and rewards to increase their visit frequency.

#### • Cluster 2: Budget-Conscious Shoppers

- Profile: High Frequency but low MonetaryValue. They shop often but spend little, likely on discounted essential items.
- o Action: Target them with "buy one, get one free" offers and weekly discount flyers.

#### • Cluster 3: At-Risk Customers

- **Profile**: Low Frequency, high Recency (haven't visited in a while), and moderate Monetary Value. They used to spend well but have stopped visiting.
- Action: Launch a "We miss you!" re-engagement campaign with a special discount to bring them back.

#### 7. Implementation & Conclusion

Once the segments are understood, marketing and inventory teams can implement these targeted strategies. The model should be periodically retrained with new data to ensure the segments remain relevant as customer behavior changes.