# **Clustering & Association Rule Mining**

## **Clustering & Association Rule Mining**

These methods help in discovering patterns in large datasets, supporting decision-making in domains like marketing, healthcare, and finance.

# **Clustering Algorithms**

Clustering is an **unsupervised learning** technique that groups similar data points together. It is used in customer segmentation, anomaly detection, and pattern recognition.

### Why is Clustering Important?

- Helps identify hidden patterns in massive datasets.
- Facilitates efficient data organization, reducing search complexity.
- Supports market segmentation by grouping customers with similar purchasing behavior.

## **Types of Clustering Algorithms**

- 1. Hierarchical Algorithms
  - These algorithms create a tree-like structure of nested clusters.
  - They allow a gradual breakdown of large groups into smaller subgroups.

#### Examples:

- Agglomerative Clustering: Starts with individual data points and merges them step-by-step.
- **Divisive Clustering**: Starts with a single large group and divides it iteratively.

#### 2. Partitional Algorithms

- These divide the dataset into non-overlapping clusters directly.
- They require pre-specifying the number of clusters.

#### **Examples:**

- K-Means: Assigns data points to k predefined clusters based on similarity.
- K-Medoids: Works like K-Means but minimizes the influence of outliers.

### **Clustering Large Databases**

For massive datasets, specialized algorithms handle scalability and efficiency.

#### **Key Algorithms:**

BIRCH Algorithm in Data Mining

**BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies)** is an unsupervised data mining algorithm designed to perform hierarchical clustering on large datasets. It is particularly effective for clustering large datasets by creating a compact summary of the data, which can then be clustered using other algorithms<sup>1</sup>.

### Key Concepts

Clustering Feature (CF)

BIRCH summarizes large datasets into smaller, dense regions called Clustering Feature (CF) entries. A CF entry is defined as an ordered triple (N, LS, SS):

- N: Number of data points in the cluster.
- **LS**: Linear sum of the data points.
- **SS**: Squared sum of the data points<sup>1</sup>.

### CF Tree

The CF tree is a height-balanced tree that gathers and manages clustering features. Each leaf node contains a sub-cluster, and every entry in a CF tree contains a pointer to a child node and a CF entry made up of the sum of CF entries in the child nodes<sup>1</sup>.

Algorithm Stages

- Building the CF Tree: BIRCH summarizes large datasets into smaller, dense regions called CF entries. The CF tree is built by comparing the location of each record with the location of each CF in the root node and passing the record to the closest CF<sup>1</sup>.
- Global Clustering: Applies an existing clustering algorithm on the leaves of the CF tree. This step can be refined to improve clustering quality<sup>1</sup>.
- Advantages and Limitations

## Advantages:

- Efficiency: BIRCH is efficient as it does not require scanning all data points and existing clusters for each clustering decision<sup>1</sup>.
- Memory Usage: It uses available memory to derive the finest possible sub-clusters while minimizing I/O costs<sup>1</sup>.
- **Incremental**: BIRCH can incrementally and dynamically cluster incoming data points<sup>1</sup>.

## Limitations:

 Metric Attributes: BIRCH can only process metric attributes, meaning it cannot handle categorical data<sup>2</sup>.

### Parameters

- Threshold: Maximum number of data points a sub-cluster in the leaf node of the CF tree can hold<sup>2</sup>.
- **Branching Factor**: Maximum number of CF sub-clusters in each node (internal node)<sup>2</sup>.

• **n\_clusters**: Number of clusters to be returned after the entire BIRCH algorithm is complete<sup>2</sup>.

BIRCH is a powerful algorithm for clustering large datasets efficiently and can be used in conjunction with other clustering algorithms to improve clustering quality  $\frac{12}{2}$ .

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

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is an unsupervised clustering algorithm that groups data points based on their density. Unlike other clustering algorithms like K-Means, DBSCAN does not require the number of clusters to be specified beforehand and can identify clusters of arbitrary shapes, making it highly effective for datasets with noise and outliers<sup>1</sup>.

Key Concepts and Parameters

DBSCAN relies on two main parameters:

- **Epsilon (eps)**: Defines the radius of the neighborhood around a data point. Points within this distance are considered neighbors.
- MinPts: The minimum number of points required within the eps radius for a point to be considered a core point<sup>1</sup>.

Types of Points in DBSCAN

- **Core Point**: A point with at least MinPts neighbors within the eps radius.
- Border Point: A point that has fewer than MinPts neighbors but is within the neighborhood of a core point.
- **Noise Point**: A point that is neither a core point nor a border point<sup>1</sup>.

Steps in DBSCAN Algorithm

- 1. **Identify Core Points**: Find all points within eps radius and identify core points with at least MinPts neighbors.
- 2. Form Clusters: For each core point, if it is not already assigned to a cluster, create a new cluster.
- 3. **Expand Clusters**: Recursively find all density-connected points and assign them to the same cluster as the core point.
- Mark Noise: Iterate through the remaining unvisited points and mark those that do not belong to any cluster as noise<sup>1</sup>.

• Evaluation Metrics

To evaluate the performance of DBSCAN, we can use metrics like the **Silhouette Score** and the **Adjusted Rand Index (ARI)**:

- Silhouette Score: Measures how similar a point is to its own cluster compared to other clusters. It ranges from -1 to 1, with higher values indicating better clustering.
- Adjusted Rand Index (ARI): Measures the similarity between the true labels and the clustering labels. It ranges from 0 to 1, with higher values indicating better clustering<sup>1</sup>.

#### • When to Use DBSCAN

DBSCAN is particularly useful when:

- The data contains clusters of arbitrary shapes.
- The number of clusters is not known beforehand.
- The dataset contains noise and outliers<sup>1</sup>.

### • CURE (Clustering Using Representatives)

 The CURE (Clustering Using Representatives) algorithm is a hierarchical clustering technique that addresses the limitations of traditional clustering algorithms like K-means and BIRCH. It is particularly effective for identifying clusters with non-spherical shapes and varying sizes, and it is more robust to outliers<sup>12</sup>.

**Key Principles** 

Random Sampling

The algorithm begins by randomly selecting a subset of data points from the dataset. This ensures that the sample is representative of different regions of the data space<sup>2</sup>.

Hierarchical Clustering

The selected points are then clustered hierarchically using either agglomerative or divisive techniques. Agglomerative clustering merges similar representatives until one central representative per cluster is reached, while divisive clustering splits them based on dissimilarities<sup>2</sup>.

Cluster Shrinkage

Once clusters are formed, each cluster's size is reduced by shrinking the outlier points towards the centroid. This helps eliminate noise and focuses on relevant patterns within each cluster<sup>2</sup>.

Final Data Point Assignment

After shrinking the clusters, all remaining non-representative points are assigned to their nearest existing representative based on distance measures like Euclidean distance<sup>2</sup>.

Steps in CURE Algorithm

- 1. **Random Sampling**: Select a subset of data points from the dataset<sup>2</sup>.
- 2. **Hierarchical Clustering**: Cluster the sampled points using hierarchical methods<sup>2</sup>.
- 3. **Distance Measures**: Compute distances between clusters during merging operations<sup>2</sup>.
- Merging Representative Points: Merge representative points from various sub-clusters into a unified set<sup>2</sup>.
- Cluster Refinement and Splitting: Refine clusters by exchanging outliers and splitting when necessary<sup>2</sup>.
- Final Membership Assignment: Assign remaining objects to their nearest representative points<sup>2</sup>.

Conclusion

The **CURE algorithm** is a powerful clustering technique that overcomes the limitations of traditional algorithms by using representative points and hierarchical clustering. Its ability to handle non-spherical clusters and outliers makes it suitable for various applications in data analysis and machine learning<sup>12</sup>.

# **Association Rule Mining**

Association rule mining helps find **relationships** between variables in datasets. It is widely used in **market basket analysis**, where businesses determine products frequently bought together.

# Why is Association Rule Mining Important?

- Helps optimize retail strategies by analyzing purchase patterns.
- Supports fraud detection by uncovering suspicious activity.
- Assists in **recommender systems** by suggesting relevant items.

# Parallel & Distributed Association Rule Mining Algorithms

Parallel and distributed association rule mining algorithms are designed to efficiently discover relationships between items in large-scale datasets. These algorithms enhance traditional **association rule mining (ARM)** by leveraging **parallel computing** and **distributed systems** to handle massive data volumes.

# **1. Association Rule Mining (ARM) Overview**

Association rule mining identifies patterns in transactional databases, such as:

- Frequent itemsets (e.g., "Customers who buy bread often buy butter").
- Association rules (e.g., "If a customer buys milk, they are likely to buy cereal").

Traditional ARM algorithms, like **Apriori**, suffer from scalability issues when dealing with large datasets. **Parallel and distributed approaches** solve this problem by distributing computations across multiple processors or machines.

# **2. Parallel Association Rule Mining Algorithms**

Parallel ARM algorithms divide the dataset into smaller chunks and process them simultaneously across multiple processors.

# 2.1 Parallel Apriori Algorithm

- Apriori is a classic ARM algorithm that generates frequent itemsets iteratively.
- Parallel Apriori distributes computations across multiple processors to speed up execution.
- Uses data partitioning and task parallelism to reduce memory overhead.

# 2.2 FP-Growth in Parallel

- FP-Growth is an alternative to Apriori that builds a frequent pattern tree (FP-tree).
- Parallel FP-Growth splits the FP-tree across processors for faster mining.
- Works well for high-dimensional datasets.

### **2.3 Eclat Algorithm**

- Uses depth-first search instead of breadth-first search.
- Parallel Eclat distributes itemset computations across multiple cores.
- Efficient for dense datasets.

# **3. Distributed Association Rule Mining Algorithms**

Distributed ARM algorithms process data across multiple machines in a network.

## 3.1 Distributed Apriori Algorithm

- Divides the dataset across multiple nodes.
- Each node computes frequent itemsets independently.
- Results are merged to form global association rules.

### 3.2 MapReduce-Based ARM

- Uses Hadoop MapReduce for large-scale ARM.
- Mapper phase: Generates local frequent itemsets.
- Reducer phase: Aggregates results across nodes.
- Works well for **big data applications**.

#### **3.3 Spark-Based ARM**

- Uses Apache Spark for in-memory processing.
- Faster than MapReduce due to distributed caching.
- Ideal for real-time association rule mining.

# 4. Comparison of Parallel vs. Distributed ARM

Aspect	Parallel ARM	Distributed ARM
Execution	Runs on multiple processors within a single machine	Runs across multiple machines in a network
Data Handling	Works well for moderate-sized datasets	Handles massive datasets efficiently
Scalability	Limited by hardware constraints	Highly scalable
Examples	Parallel Apriori, FP-Growth, Eclat	Distributed Apriori, MapReduce, Spark

# **5. Applications of Parallel & Distributed ARM**

- Retail & E-commerce Market basket analysis for product recommendations.
- Healthcare Identifying disease correlations in medical records.
- Cybersecurity Detecting fraud patterns in financial transactions.
- Social Media Analyzing user behavior for targeted advertising.

### **Parallel Computing**

- All processors work on the same problem simultaneously, breaking it into smaller tasks.
- Shared memory is often used, meaning all processors access the same dataset.
- Example: Training a large neural network—each processor handles a fraction of the calculations.

# **Distributed Computing**

- Different nodes work on different parts of the same problem by communicating over a network.
- Each node has its own memory and processing power, meaning they work independently but contribute to a unified result.
- Example: A search engine like Google—different servers index webpages separately but contribute to a single search database.