# Clustering

## Group Members

Rabadiya Utsav - 202201081

Ridham Patel – 202201430

## Problem Area

In today's highly competitive business landscape, **customer segmentation** plays a crucial role in understanding diverse customer groups and their behaviors. By segmenting customers into distinct groups, businesses can tailor their marketing strategies, personalize customer experiences, and improve product offerings. Segmentation can lead to increased sales, customer satisfaction, and loyalty by identifying key patterns in customer behavior.

**Clustering algorithms** are one of the primary methods used for customer segmentation. Among the many clustering algorithms, **k-means** is widely used due to its simplicity and efficiency. However, k-means is sensitive to outliers and noise, which can distort the final clusters. For example, a few extreme data points (outliers) can shift the cluster centroids significantly, resulting in inaccurate clustering outcomes.

To overcome these limitations, the **k-medoid algorithm**, a more robust variation of k-means, is used. Unlike k-means, which computes the centroid (average) of all points in a cluster, k-medoid selects a representative **medoid**, which is an actual data point that minimizes the sum of dissimilarities between itself and all other points in the cluster. This characteristic makes **k-medoid more resistant to outliers** and provides more reliable clustering in noisy datasets.

For this project, we aim to:

1. **Implement the k-medoid algorithm** for clustering on **Apache Spark**, a distributed computing framework.
2. Apply the k-medoid algorithm to a suitable dataset for **customer segmentation**.
3. **Analyze and compare** the results with the k-means algorithm to demonstrate the robustness of k-medoid in real-world data scenarios.

Apache Spark is chosen as the platform for the algorithm's implementation because of its ability to handle large-scale data efficiently through distributed computing. Spark's **MLlib** library provides robust support for various machine learning and clustering algorithms, making it a suitable platform for processing large datasets. The implementation will draw upon the strategies outlined in the research paper **"Parallelizing k-means-based clustering on Spark"** by Wang et al. (2016). While the paper focuses on k-means, the distributed processing strategies will be adapted for k-medoid, ensuring efficient performance on large datasets.

## Expected Outcome

The project aims to deliver the following outcomes:

1. **Implementation of k-medoid on Spark**: The primary goal is to successfully implement the k-medoid algorithm on Apache Spark. Spark's **in-memory** processing capabilities will ensure efficient execution even on large customer datasets. The implementation will follow the principles outlined by Wang et al. for parallel processing, ensuring the k-medoid algorithm can scale efficiently on Spark's distributed architecture.
2. **Customer segmentation on a real-world dataset**: The k-medoid algorithm will be applied to a relevant customer dataset (e.g., retail transaction data, e-commerce behavior, or customer demographics). The dataset will be chosen based on its suitability for segmentation, focusing on attributes like purchasing patterns, geographic location, customer preferences, and more. The output will be clusters representing distinct customer segments.
3. **Comparison with k-means**: To highlight the strengths of k-medoid over k-means, we will compare both algorithms' performance in terms of:
   o **Clustering quality**: Evaluating how well each algorithm segments the data.
   o **Robustness to noise**: Testing the algorithms on datasets with and without outliers.
   o **Execution performance**: Analyzing how efficiently the algorithms run on Spark, especially in distributed environments.
4. **Business insights from customer segmentation**: The ultimate goal is to provide actionable business insights based on the customer segments identified by the k-medoid algorithm. By understanding distinct customer groups, businesses can:
   o Develop personalized marketing campaigns for each segment.
   o Optimize product recommendations and improve customer satisfaction.
   o Identify high-value customers for loyalty programs or targeted advertising.
5. **Documentation and demonstration**: The final outcome will include a detailed report documenting the implementation process, the challenges faced, and the results obtained. A demonstration of the algorithm running on Spark with the chosen dataset will be included, showcasing its ability to handle large datasets and its effectiveness in customer segmentation.

## Selected Readings

The project will draw on the following key readings to guide the implementation and analysis:

1. **Wang, Bowen, et al. "Parallelizing k-means-based clustering on Spark."** 2016 International Conference on Advanced Cloud and Big Data (CBD). IEEE, 2016.
   o This paper presents a detailed strategy for parallelizing the k-means clustering algorithm on Apache Spark. While the paper specifically deals with k-means, its distributed processing techniques will be adapted for k-medoid clustering. This paper will guide the implementation, focusing on how Spark can be leveraged for efficient clustering across large datasets.

2. **Kaufman, Leonard, and Peter J. Rousseeuw. "Finding Groups in Data: An Introduction to Cluster Analysis."** Wiley, 2005.
   - o This foundational text provides an in-depth exploration of clustering methods, including k-medoid. The book offers theoretical insights into how k-medoid operates and its advantages over other clustering methods. This resource will be essential in understanding the nuances of the k-medoid algorithm and applying it effectively to customer segmentation.

3. **Apache Spark MLlib Documentation**: Spark's official documentation will be a critical resource for understanding how to implement clustering algorithms, particularly using the MLlib library. Spark's distributed computing framework and its optimization techniques will be integral to the project's success.

4. **Xu, Rui, and Donald Wunsch. "Clustering." IEEE Transactions on Neural Networks 16.3 (2005): 645-678.**
   - o This paper offers a comprehensive overview of different clustering methods, including k-means and k-medoid. It provides useful insights into how these algorithms can be applied to various real-world datasets and compares their effectiveness in different scenarios. This comparison will help us understand the strengths and limitations of k-medoid for customer segmentation.

5. **Further articles and tutorials on Apache Spark and distributed machine learning**:
   - o Additional tutorials and articles from reputable sources will be referenced as needed to ensure efficient and effective implementation of the k-medoid algorithm on Spark.