



**SECP2753 – DATA MINING**

**SEMESTER 2 2021/2022**

---

**GROUP 2**

**TOPIC 4 - CLUSTERING**

**Members:**

1. **GROUP LEADER:** NAYLI NABIHAH JASNI (A20EC0105)
2. MIKHEL ADAM BIN MUHAMMAD EZRIN (A20EC0237)
3. MUHAMMAD DINIE HAZIM BIN AZALI (A20EC0084)
4. MADINA SURAYA BINTI ZHARIN (A20EC0203)

# TABLE OF CONTENT

<b>1.0 Definition</b>	<b>3</b>
1.1 Cluster	3
1.2 Cluster Analysis	3
1.3 Clustering	3
<b>2.0 What is Similarity</b>	<b>4</b>
2.1 Definition	4
2.2 Euclidean	4
2.3 Non-Euclidean	5
2.3.1 Jaccard Distance (Jaccard Index)	5
2.3.2 Hamming Distance	6
2.3.3 Manhattan Distance	6
<b>3.0 Applications of Clustering</b>	<b>8</b>
<b>4.0 Types of Clustering</b>	<b>9</b>
4.1 Hierarchical Algorithm	9
4.2 Partitional Clustering	10
<b>5.0 Clustering Algorithms</b>	<b>11</b>
<b>6.0 How does the K-Means clustering algorithm work?</b>	<b>13</b>
<b>7.0 Classification vs. Clustering</b>	<b>15</b>
<b>8.0 Reference</b>	<b>16</b>

## **1.0 Definition**

### **1.1 Cluster**

A data cluster is a subset of a larger dataset in which each data point is closer to the cluster center than to other cluster centers in the dataset, as determined via cluster analysis, which involves iteratively decreasing squared distances (Glenn, 2022).

### **1.2 Cluster Analysis**

Cluster analysis is often one of the first steps in the analysis of data, as such, it is an effort at unsupervised learning usually in the context of very little a priori knowledge (Davies et al., 1979).

Cluster analysis is the process of creating data clusters by minimizing the distance between data points and a reference (Glenn, 2022).

### **1.3 Clustering**

Clustering is one of the common methods of unsupervised learning in which data is segmented based on the similarity. Clustering is assigning a set of instances to subgroups class clusters so that each cluster is very similar with other clusters (Teng et al., 2022).

Clustering identifies groups of related records that can be used as a starting point for exploring further relationships. It is the first step in data mining analysis (Bijuraj, 2013).

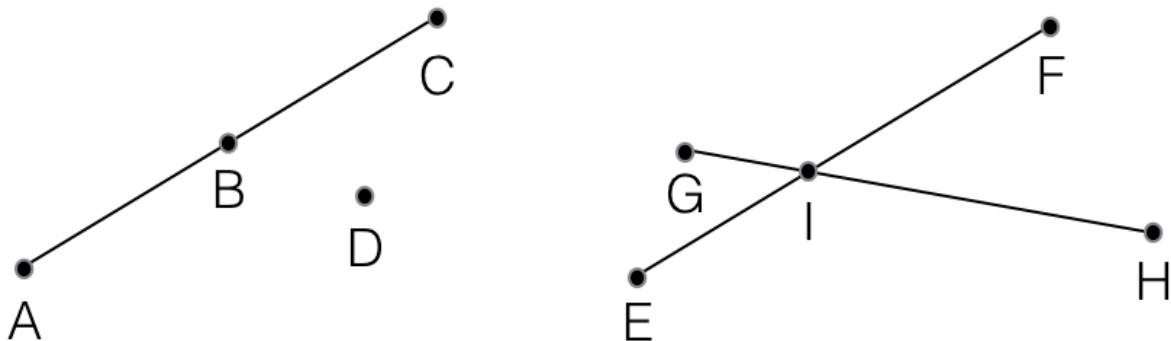
## 2.0 What is Similarity

### 2.1 Definition

Similarity, also known as **distance measures**, are core components of distance-based clustering algorithms used to group similar data points into the same clusters, while dissimilar or distant data points are placed into different clusters. The distance between various data points can be defined as a similarity measure. Similarity is an amount that reflects the strength of relationship between two data items. Meanwhile, dissimilarity is the measurement of divergence between two data items.

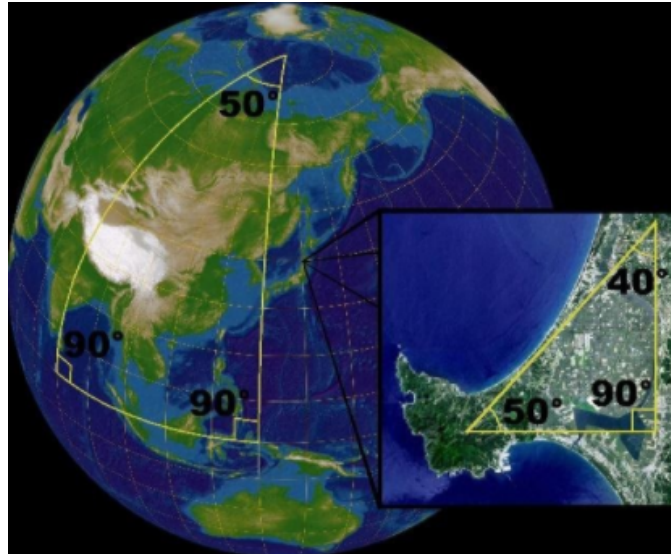
### 2.2 Euclidean

Euclidean geometry is the study of geometrical shapes (plane and solid) and figures based on different axioms and theorems. It is basically introduced for flat surfaces or plane surfaces. For geometrical problems, Euclidean distance is considered the standard metric. It's simply the ordinary distance between two points. Euclidean distance is also extensively used in clustering problems. It is also the default distance measure used by the K-means algorithm. The Euclidean distance determines the root of square differences between the coordinates of a pair of objects. The most basic terms of geometry are a point, a line, and a plane. A point has no dimension (length or width), but it does have a location. A line is straight and extends infinitely in the opposite directions. A plane is a flat surface that extends indefinitely.



## 2.3 Non-Euclidean

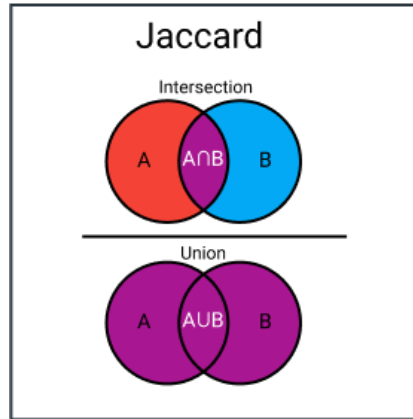
Non-Euclidean spaces are basically any geometry that is not the same as Euclidean geometry. An example of Non-Euclidean geometry can be seen by drawing lines on a sphere, that is straight lines that are parallel at the equator can meet at the poles. This “triangle” has an angle sum of  $90 + 90 + 50 = 230$  degrees.



For non-Euclidean spaces, there are several distance measures used. Some of them include Jaccard distance, Hamming distance, and Manhattan distance.

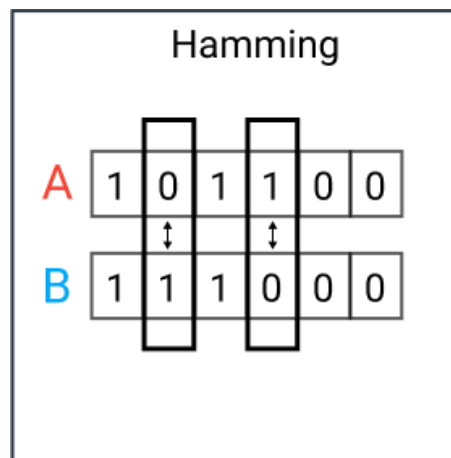
### 2.3.1 Jaccard Distance (Jaccard Index)

The Jaccard Index (also known as Intersection over Union) is a metric for calculating sample set similarity and diversity. It is equal to the intersection size divided by the sample set union size. In practice, it is the total number of similar entities between sets divided by the total number of entities. For example, if two sets have 1 entity in common and there are 5 different entities in total, then the Jaccard index would be  $1/5 = 0.2$ .



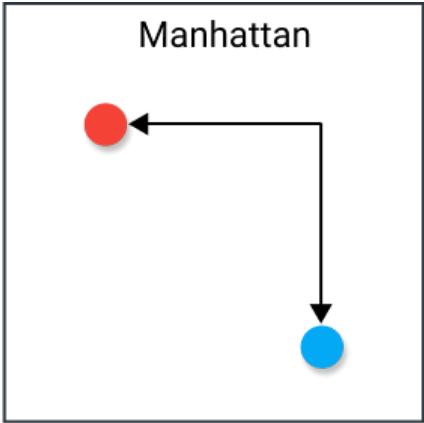
### 2.3.2 Hamming Distance

The Hamming distance is the difference in values between two vectors. It's usually used to compare two binary strings of the same length. It can also be used to compare the similarity of strings by calculating the number of characters that differ.



### 2.3.3 Manhattan Distance

The Manhattan distance, also known as the Taxicab distance or City Block distance, is a calculation that determines the distance between two real-valued vectors. Consider vectors that describe objects on a chessboard-like uniform grid. The Manhattan distance is the distance between two vectors if they could only move in the same direction. The distance is calculated without any diagonal movement.



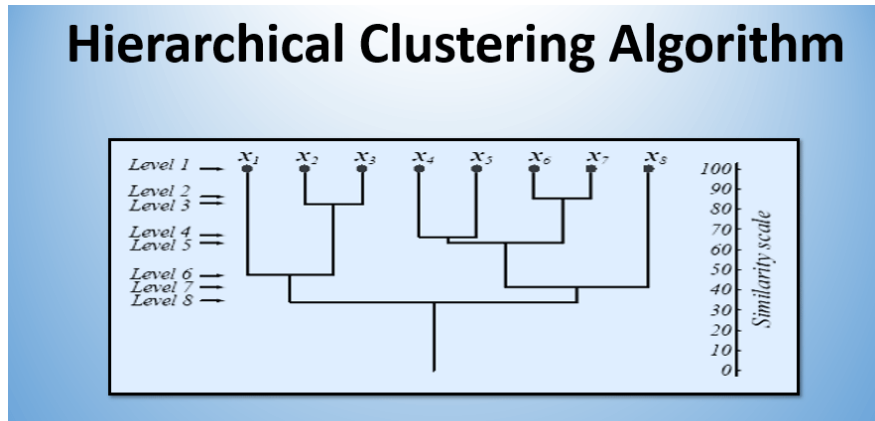
### 3.0 Applications of Clustering

Applications	Algorithms	Area
KNIME	<ul style="list-style-type: none"><li>- K-Means</li><li>- K-Medoids</li><li>- Hierarchical Clustering</li><li>- Fuzzy c-Means</li><li>- SOTA (self organizing tree algorithm)</li></ul>	<ul style="list-style-type: none"><li>- Voting patterns</li><li>- Market segmentation</li></ul>
Orange	<ul style="list-style-type: none"><li>- K-Means</li><li>- SOM (self organizing maps)</li><li>- Hierarchical Clustering</li><li>- MDS (multidimensional scaling)</li></ul>	<ul style="list-style-type: none"><li>- Anomaly detection</li><li>- Medical imaging</li></ul>
RapidMiner Community Edition	<ul style="list-style-type: none"><li>- Hierarchical Clustering</li><li>- Support Vector Clustering</li><li>- Top-Down Clustering</li><li>- K-Means</li><li>- K-Medoids</li></ul>	<ul style="list-style-type: none"><li>- Social network analysis</li><li>- Market segmentation</li></ul>
Tanagra	<ul style="list-style-type: none"><li>- K-Means</li><li>- SOM</li><li>- LVQ (Learning Vector Quantizers)</li><li>- Hierarchical Clustering</li></ul>	<ul style="list-style-type: none"><li>- Biological models of neural systems</li></ul>
Weka	<ul style="list-style-type: none"><li>- DBSCAN</li><li>- COBWEB</li><li>- K-Means</li><li>- EM (Expectation maximization)</li></ul>	<ul style="list-style-type: none"><li>- Economic</li><li>- Recommendation engines</li></ul>



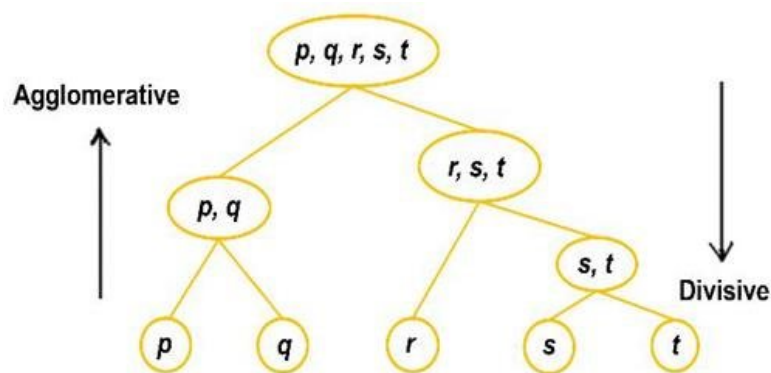
## 4.0 Types of Clustering

### 4.1 Hierarchical Algorithm



Hierarchical algorithm is an unsupervised Machine Learning technique. Its aim is to determine natural groupings based on the data's features. Hierarchical algorithm aims to find nested groups of data by constructing a hierarchy. It's similar to the plant or animal kingdom's biological taxonomy. The hierarchical tree known as a dendrogram is commonly used to describe hierarchical groups.

There are 2 types of Hierarchical algorithm:



#### a. Divisive

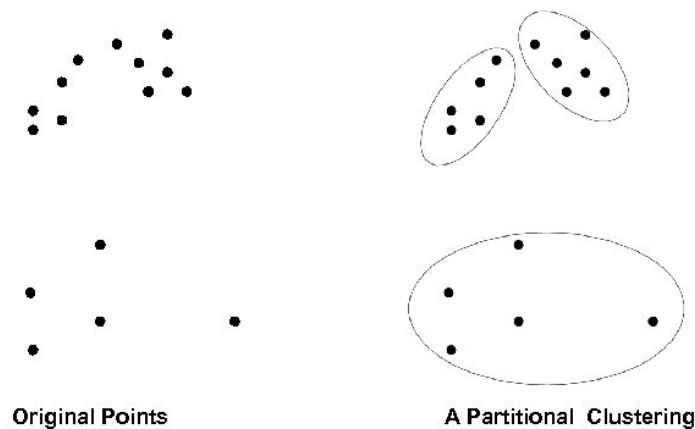
This is a top-down technique, in which the full data set is first considered as a single group, and then subgroups are created repeatedly. When the number of clusters in a hierarchical clustering algorithm is known, the division process ends when that number is

reached. Otherwise, the procedure ends when the data can no longer be split, implying that the current iteration's subgroup is identical to the prior iteration's (one can also consider that the division stops when each data point is a cluster).

b. Agglomerative

It's a bottom-up technique that relies on cluster merging. The data is first divided into  $m$  singleton clusters (where  $m$  is the number of samples/data points). Iteratively, two clusters are merged into one, lowering the number of clusters in each iteration. When all clusters have been merged into one or the desired number of clusters has been reached, the process of merging clusters comes to an end.

## 4.2 Partitional Clustering



This clustering method classifies the information into multiple groups based on the characteristics and similarity of the data. The number of clusters that must be formed for the clustering procedures is specified by the data analysts.

When a database( $D$ ) contains multiple( $N$ ) objects, the partitioning method creates user-specified( $K$ ) data partitions, each of which represents a cluster and a specific region. Many algorithms fall within the partitioning method category, including K-Mean, PAM(K-Medoids), CLARA algorithm (Clustering Large Applications), and others.

## 5.0 Clustering Algorithms

Algorithm	Function	Advantages	Disadvantages	Field Area	Application
K-Means	Used to cluster numerical data attributes.	<ul style="list-style-type: none"> <li>- Widely applied by variety of packages</li> <li>- Fast convergence when apply clustering to small datasets</li> <li>- Easy to implement</li> </ul>	<ul style="list-style-type: none"> <li>- Computationally expensive for large datasets (k will be large)</li> <li>- Occasionally, it will be hard to choose the initial value for the number of clusters (k)</li> <li>- Strong sensitivity to outliers</li> </ul>	Marketing	Tableau R
Agglomerative	A type of hierarchical clustering used to group objects based on their similarity	<ul style="list-style-type: none"> <li>- Easy to implement and understand</li> <li>- Able to produce an ordering of objects</li> <li>- No need to pre-specify the number of clusters</li> </ul>	<ul style="list-style-type: none"> <li>- It gives the best result only in some cases</li> <li>- Cannot undo what was done previously</li> <li>- The various distance metrics for measuring distances between clusters may produce different results.</li> </ul>	Bioinformatics	XLSTAT Excel
DBSCAN	Used in density-based clustering	<ul style="list-style-type: none"> <li>- Does not need to specify the number of clusters beforehand</li> </ul>	<ul style="list-style-type: none"> <li>- Sometimes, determining an appropriate distance of neighborhood (eps) is hard and requires domain</li> </ul>	Bio-Medical	XLSTAT

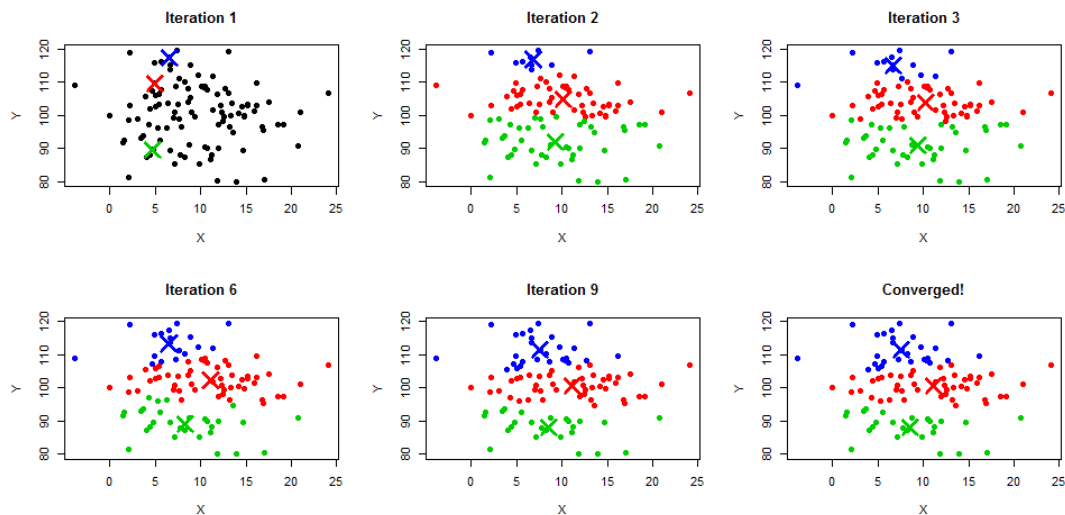
		<ul style="list-style-type: none"> <li>- Performs well with arbitrary shapes clusters</li> <li>- Robust to outliers and can detect the outliers</li> </ul>	<p>knowledge</p> <ul style="list-style-type: none"> <li>- If clusters are very differ in terms of in-cluster densities, DBSCAN is not well suited to define the clusters</li> </ul>		
Self-Organizing-Map (SOM)	Used to help to understand high dimensional data by reducing the dimensions of data to a map	<ul style="list-style-type: none"> <li>- Data mapping is easily interpreted</li> <li>- Capable of organizing large and complex data sets</li> </ul>	<ul style="list-style-type: none"> <li>- Hard to determine what input weights to use</li> <li>- Mapping can resulting to divided clusters</li> <li>- Nearby points need to behave similarly</li> </ul>	Economy	XLSTAT
Fuzzy C-Means (FCM)	Used to cluster multidimensional data by assigning each point a membership in each cluster center	<ul style="list-style-type: none"> <li>- Gives out the best result for overlapped data compared to k-means</li> <li>- Low time complexity</li> </ul>	<ul style="list-style-type: none"> <li>- Sensitive to the initial values of k and p</li> <li>- Sensitive to outliers</li> </ul>	Marketing	MATLAB

## 6.0 How does the K-Means clustering algorithm work?

K-means clustering uses “centroids”, K different randomly-initiated points in the data, and assigns every data point to the nearest centroid. After every point has been assigned, the centroid is moved to the average of all of the points assigned to it. Then the process repeats: every point is assigned to its nearest centroid, centroids are moved to the average of points assigned to it. The algorithm is done when no point changes the assigned centroid. K-means clustering distinguishes itself from Hierarchical since it creates K random centroids scattered throughout the data (Robinson, 2022).

The algorithm will look like below:

- Initialize K random centroids. You could pick K random data points and make those your starting points. Otherwise, you pick K random values for each variable.
- For every data point, look at which centroid is nearest to it. Using some sort of measurement like Euclidean or Cosine distance.
- Assign the data point to the nearest centroid.
- For every centroid, move the centroid to the average of the points assigned to that centroid.
- Repeat the last three steps until the centroid assignment no longer changes. The algorithm is said to have “converged” once there are no more changes.



- Iteration 1 shows the random centroid centers.
- Iteration 2 shows the new location of the centroid centers.
- Iteration 3 has a handful more blue points as the centroids move.
- Jumping to iteration 6, we see the red centroid has moved further to the right.

5. Iteration 9 shows the green section is much smaller than in iteration 2, blue has taken over the top, and the red centroid is thinner than in iteration 6.
6. The 9th iteration's results were the same as the 8th iteration, so it has "converged".

## 7.0 Classification vs. Clustering

Aspects	Classification	Clustering
Type of Machine Learning	Supervised	Unsupervised
How does it work	It uses algorithms to categorize the new data as per the observations of the training set.	It uses statistical concepts in which the data set is divided into subsets with the same features.
Number of class	Known	Unknown
Training Data	Required	Not required
Complexity	More complex than clustering	Less complex
Example of Algorithm	<ul style="list-style-type: none"><li>- Naive Bayes Classifier</li><li>- Decision Tree</li><li>- Random Forests</li></ul>	<ul style="list-style-type: none"><li>- K-Means</li><li>- Mean-Shift Clustering</li><li>- Gaussian (EM) Clustering</li></ul>

## 8.0 Reference

- Chen, J., Zhang, H., Pi, D., Kantardzic, M., Yin, Q., & Liu, X. (2021). A Weight Possibilistic Fuzzy C-Means Clustering Algorithm. Scientific Programming, 2021.
- Dalby, A. R. (2015, December 11). *A Comparison Study on Similarity and Dissimilarity Measures in Clustering Continuous Data*. National Center for Biotechnology Information.  
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4686108/#:%7E:text=Similarity%20or%20distance%20measures%20are,are%20placed%20into%20different%20clusters>
- DBSCAN (density-based spatial clustering of applications with noise). (n.d.). XLSTAT, Your data analysis solution. Retrieved May 17, 2022, from  
<https://www.xlstat.com/en/solutions/features/dbscan-density-based-spatial-clustering-of-applications-with-noise>
- Edpresso Team. (2019, October 11). *Classification vs. clustering*. Educative: Interactive Courses for Software Developers. Retrieved May 17, 2022, from  
<https://www.educative.io/edpresso/classification-vs-clustering>
- Harmouch, M. (2021, May 2). *17 clustering algorithms used in data science & mining*. Retrieved May 17, 2022, from  
<https://towardsdatascience.com/17-clustering-algorithms-used-in-data-science-mining-49dbfa5bf69a#7e1d>
- ML | Classification vs clustering. (2019, October 3). GeeksforGeeks. Retrieved May 17, 2022, from  
<https://www.geeksforgeeks.org/ml-classification-vs-clustering/>
- Partitioning Method (K-Mean) in Data Mining. (2020, February 5). GeeksforGeeks. Retrieved May 17, 2022, from  
<https://www.geeksforgeeks.org/partitioning-method-k-mean-in-data-mining/>
- Pedamkar, P. (n.d.). *Hierarchical Clustering Algorithm | Types & Steps of Hierarchical Clustering*. eduCBA. Retrieved May 17, 2022, from  
<https://www.educba.com/hierarchical-clustering-algorithm/>



Prado, K. S. do. (2019, June 3). *How DBSCAN works and why should we use it?* Medium.  
Retrieved May 17, 2022, from  
<https://towardsdatascience.com/how-dbscan-works-and-why-should-i-use-it-443b4a191c80>

Robinson, D. (2022). *K-Means Clustering – What it is and How it Works – Learn by Marketing*.  
Learn by Marketing. Retrieved May 17, 2022, from  
<https://www.learnbymarketing.com/methods/k-means-clustering>