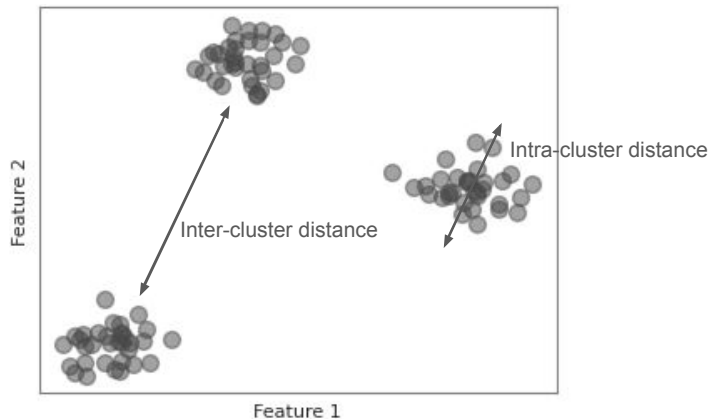# What is Clustering

**Definition:** Clustering is a type of unsupervised machine learning technique where the goal is to group similar data points into clusters or subgroups based on the inherent patterns or similarities within the data.

**Objective:** The primary objective of clustering is to group data in such a way that there is minimum distance among data points within a cluster and a maximum distance among the clusters themselves.



**Intercluster distance:** Intercluster distance refers to the distance between different clusters in a clustering algorithm. It measures the dissimilarity or separation between clusters. In other words, it quantifies how distinct or separate different clusters are from each other.

**Intracluster distance:** Intracluster distance, also known as intragroup distance, refers to the average distance between all pairs of points within the same cluster. It measures the compactness or cohesion of a cluster, indicating how tightly the data points within a cluster are clustered together.

# Types of Clustering

Clustering algorithms can be broadly categorized into several types based on their clustering approach. Here are some common types of clustering:

1. **Partitioning Methods**
   Partitioning methods are algorithms that divide data into separate clusters with no overlap. They begin with an initial grouping and then iteratively adjust data points to enhance a chosen criterion. One of the prominent algorithms in this category is K-Means.

2. **Hierarchical Clustering**
   Hierarchical clustering forms cluster hierarchies via bottom-up (agglomerative) or top-down (divisive) methods. Agglomerative starts with individual points, merging them into larger clusters, while divisive begins with a single cluster, splitting it into smaller ones.

3. **Density-Based Methods**
   Density-based methods identify clusters as areas of higher density than their surroundings, treating sparse regions as noise or border points. A key example is DBSCAN, which excels in finding clusters of varied shapes and distinguishing dense clusters from sparse noise.

# Applications of Clustering

**Customer Segmentation:** This strategy groups customers by shared traits like buying habits, spending, and preferences, helping businesses customize marketing and products for each segment.

**Document Clustering:** This technique organizes documents into clusters based on content similarity, aiding in information discovery and organization in natural language processing and information retrieval.

**Anomaly Detection:** This entails spotting data anomalies that stray from normal patterns, crucial for identifying issues in areas like fraud, network security, and fault detection.

**Genomic Clustering:** This involves categorizing genes by similar expression or functions. Through analyzing gene expression under various conditions, scientists can find gene clusters with shared regulation or roles, aiding in understanding gene functions and disease processes.
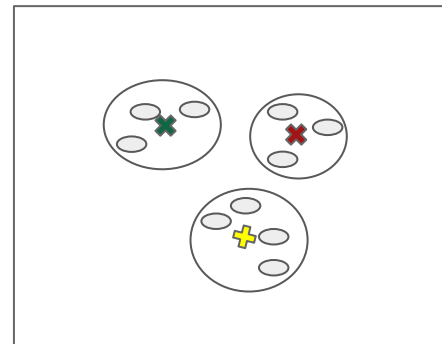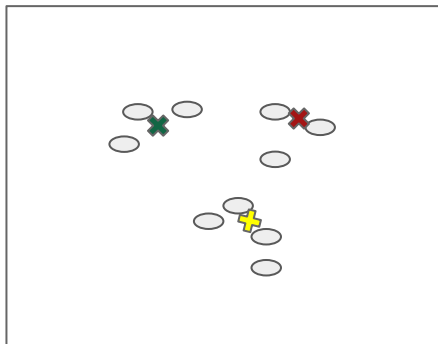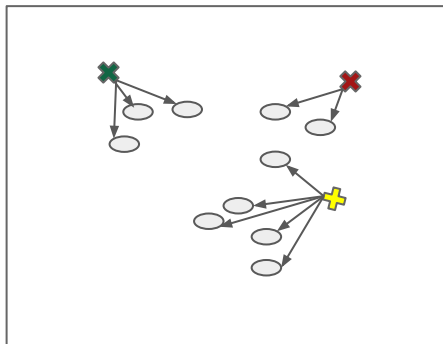
# K-Means Clustering

K-Means clustering is a widely used unsupervised learning algorithm in data mining and machine learning. Its goal is to partition a set of observations into a predefined number k of clusters, where each observation belongs to the cluster with the nearest mean (centroid). The algorithm operates through the following steps:

**Initialization:** Randomly select k initial cluster centroids.

**Assignment:** Assign each data point to the cluster whose centroid is closest (typically based on Euclidean distance).

**Update Centroids:** Calculate the new centroids of the clusters by taking the mean of all data points assigned to each cluster.

**Repeat:** Iterate between steps 2 and 3 until convergence. Convergence is reached when the centroids no longer change significantly between iterations, or when a maximum number of iterations is reached.

# Centroid Initialization Problem

The random initial selection of centroids in K-means clustering has a significant impact on the algorithm's outcome and performance. There are several challenges associated with this initial centroid selection:

**Convergence to Local Minima:** Poor initial centroid placement can lead K-means to converge to suboptimal local minima, affecting the quality of clustering.

**Difficulty with Complex Cluster Shapes:** The algorithm struggles with non-spherical or unevenly sized clusters due to its simplistic initial centroid selection method.

**Influence of Outliers:** Outliers can distort initial centroid positioning, leading to clusters that don't accurately reflect the data distribution.

# Centroid Initialization Solutions

**K-means++:** This initialization technique selects initial centroids that are spread out from each other. The first centroid is chosen randomly from the data points, and subsequent centroids are chosen from the remaining points with probability proportional to the square of the distance to the nearest existing centroid. This method aims to place initial centroids in a way that encourages convergence to a better global solution.
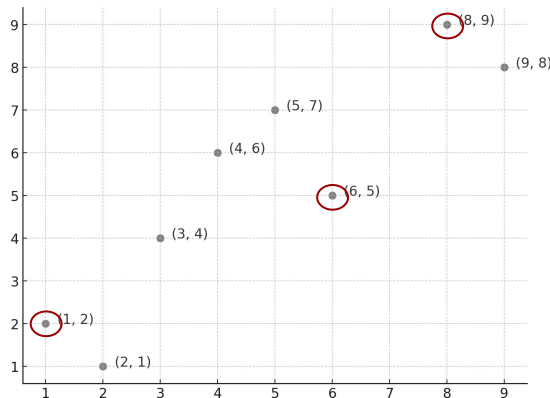
**Multiple Runs with Different Seeds:** Running K-means multiple times with different random initializations and selecting the best outcome based on a criterion can improve the chances of finding a better clustering solution.

**Using External Information:** If domain knowledge or external information about the data is available, it can guide the initial selection of centroids, potentially leading to improved clustering outcomes.

# K-means++

Data points

1. [1, 2]

2. [2, 1]

3. [3, 4]

4. [5, 7]

5. [8, 9]

6. [9, 8]

7. [6, 5]

8. [4, 6]



**Selected Centers:**

**First Center:** [6, 5] (chosen randomly as the initial step)

**Second Center:** [8, 9] (chosen based on the squared distance probability distribution from the first center)

**Third Center:** [1, 2] (chosen similarly, considering the distances from all previously selected centers)
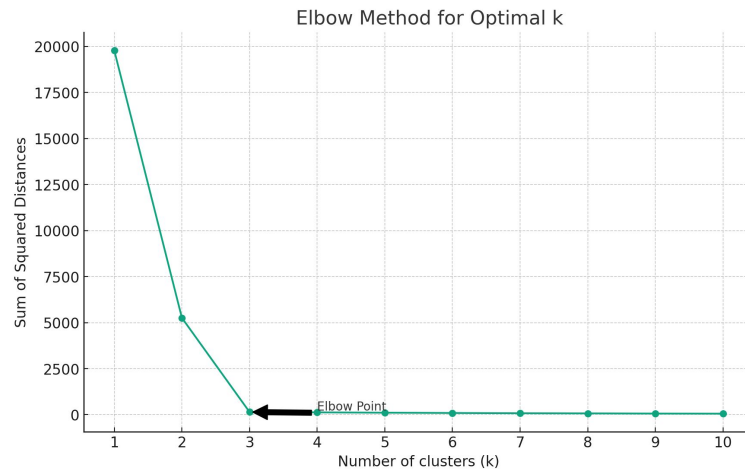
These selected centers using the KMeans++ method ensure a good spread across the data space, aiming for a better starting point for the KMeans clustering process.

There is a setup cost for k-means++, but it converges better and faster.

# Finding Optimal K

**Elbow Method:**

The Elbow Method involves running the K-means algorithm on the dataset for a range of k values (say, 1 to 10) and calculating the sum of squared distances (SSD) between data points and their respective cluster <span style="color:red">centroids</span> for each k. The goal is to find the k at which the SSD begins to decrease at a slower rate, creating an "elbow" in the plot of SSD versus k. This point is considered to be the optimal number of clusters because increasing k beyond this point does not improve clustering performance significantly.

### Elbow Method for Optimal k

Sum of Squared Distances

Elbow Point

Number of clusters (k)

The Elbow Method plot shows at k=3 the sum of squared distances (SSD) begins to decrease at a slower rate, indicating that k=3 is the optimal number of clusters for this dataset.