# Introduction to K-Means Clustering

## 1    What is K-Means Clustering?

K-means clustering is a method used to group similar data points together. Imagine you have a bunch of points on a piece of paper, and you want to organize them into groups (called **clusters**) so that points in the same group are close to each other. K-means helps you do this automatically.

## 2    Why is it Useful?

K-means clustering is useful because:

- It helps you find patterns in data. For example:

    - Grouping customers based on their shopping habits.
    - Organizing pictures of animals into cats, dogs, and birds.
    - Finding groups of similar genes in biology.

- It's an **unsupervised learning** technique, which means you don't need to know the groups beforehand. The algorithm figures it out for you.

## 3    How Does K-Means Work?

Let's break it down step by step:

### 3.1    1. Choose the Number of Clusters ($k$)

- You decide how many groups ($k$) you want. For example, if you're organizing animals, you might choose $k = 3$ for cats, dogs, and birds.

- This is the only input you need to give the algorithm.

### 3.2    2. Place $k$ Centroids Randomly

- A **centroid** is like the "center point" of a cluster.

- At the start, the algorithm randomly places $k$ centroids on your data (like dropping $k$ pins on a map).

### 3.3    3. Assign Points to the Nearest Centroid

- The algorithm looks at each point and assigns it to the nearest centroid. For example:

    - If a point is closer to Centroid A than Centroid B, it joins Cluster A.

- This creates $k$ groups of points.

### 3.4    4. Move the Centroids to the Center of Their Clusters

- After assigning points, the algorithm recalculates the centroid for each cluster by finding the average position of all the points in that cluster.

- This moves the centroid to the "center" of the cluster.

### 3.5  5. Repeat Until Convergence

- The algorithm repeats the assignment and centroid update steps until the centroids stop moving (or until the changes are very small).

- At this point, the clusters are finalized, and the algorithm stops.

## 4  Example

Imagine you have the following 2D points on a plane:

$$(1,1), (1,2), (2,1), (5,4), (6,5), (6,4)$$

If you choose $k = 2$, the algorithm might:

1. Randomly place two centroids, say at $(1,1)$ and $(6,5)$.

2. Assign points closer to $(1,1)$ to Cluster 1 and points closer to $(6,5)$ to Cluster 2.

3. Recalculate the centroids as the average of the points in each cluster.

4. Repeat until the centroids stabilize.

## 5  Key Points to Remember

- K-means is simple and fast but requires you to choose $k$ in advance.

- It works best when the data is naturally grouped into spherical clusters.

- The results can vary depending on where the centroids are initially placed.

## 6  Conclusion

K-means clustering is a powerful and easy-to-understand algorithm for grouping data into clusters. By following the steps of initialization, assignment, and updating, it organizes data points into meaningful groups without needing prior knowledge of the groups.