

Clustering

Can You Group These Gems?

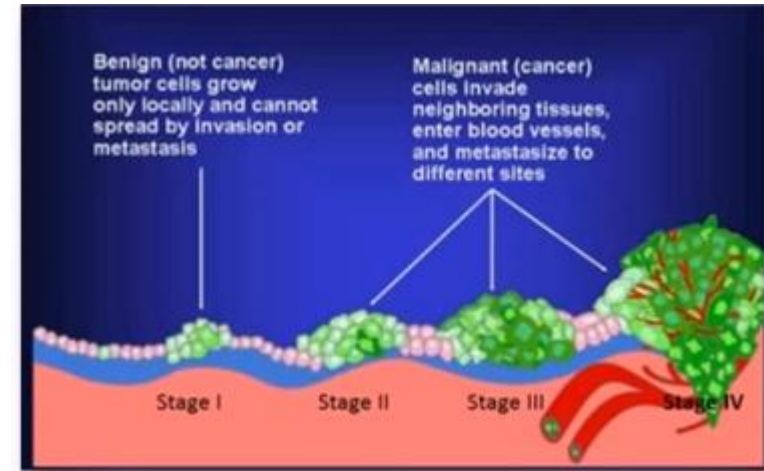


What is Clustering?

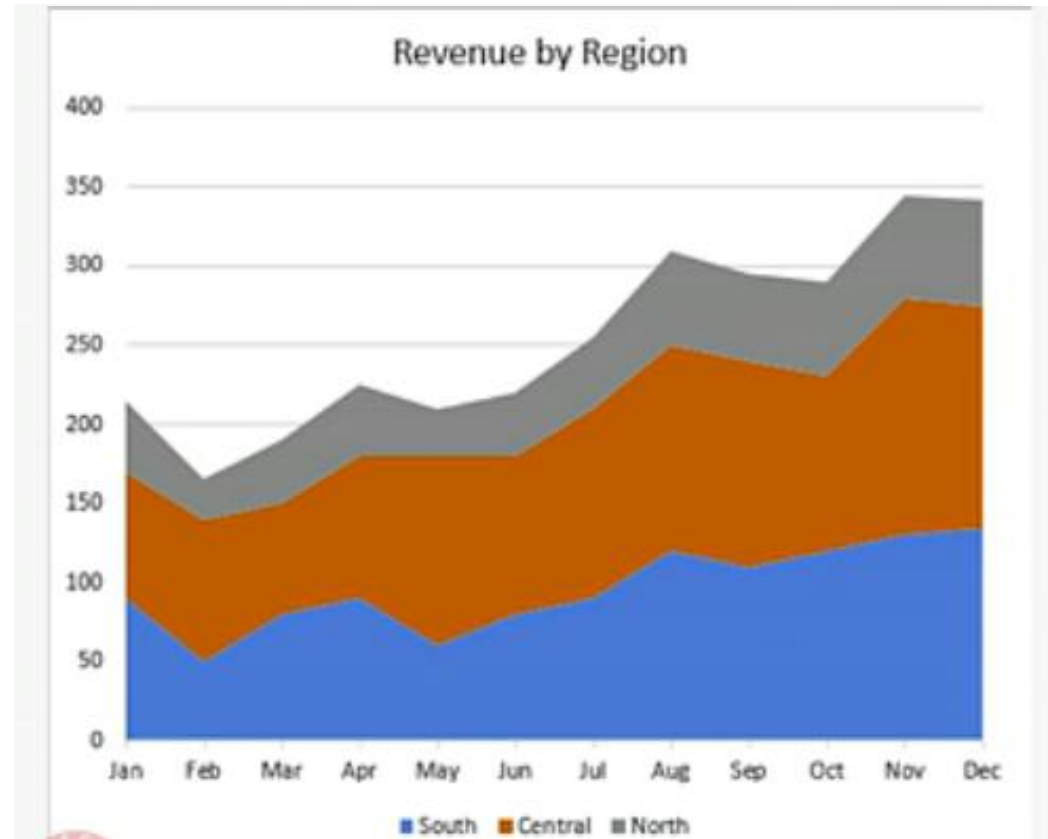
- Clustering is an algorithm which finds natural groupings inside your data when these groupings are not obvious.
- It finds the hidden variable that accurately segment your data.
- A cluster is a collection of data instances which are similar to each other and are dissimilar to data instances in other instances.

Can You Identify These Groups?





Cluster for Analysis



Cluster Techniques

Partition Methods

K-Means

K-Medoids

Hierarchical Methods

Bottom-Up (Agglomerative)

Top-Down (Divisive)

Density-Based Methods

DBSCAN

OPTICS

DENCLUE

Grid-Based Methods

STING

CLIQUE

Prototype-Based Methods

Fuzzy Clustering

SOM

Graph-Based Methods

MST

OPOSSUM

Jarvis-Patrick

SNN

Scalable

BIRCH

CURE

KM

K-Means: Within and Between Cluster

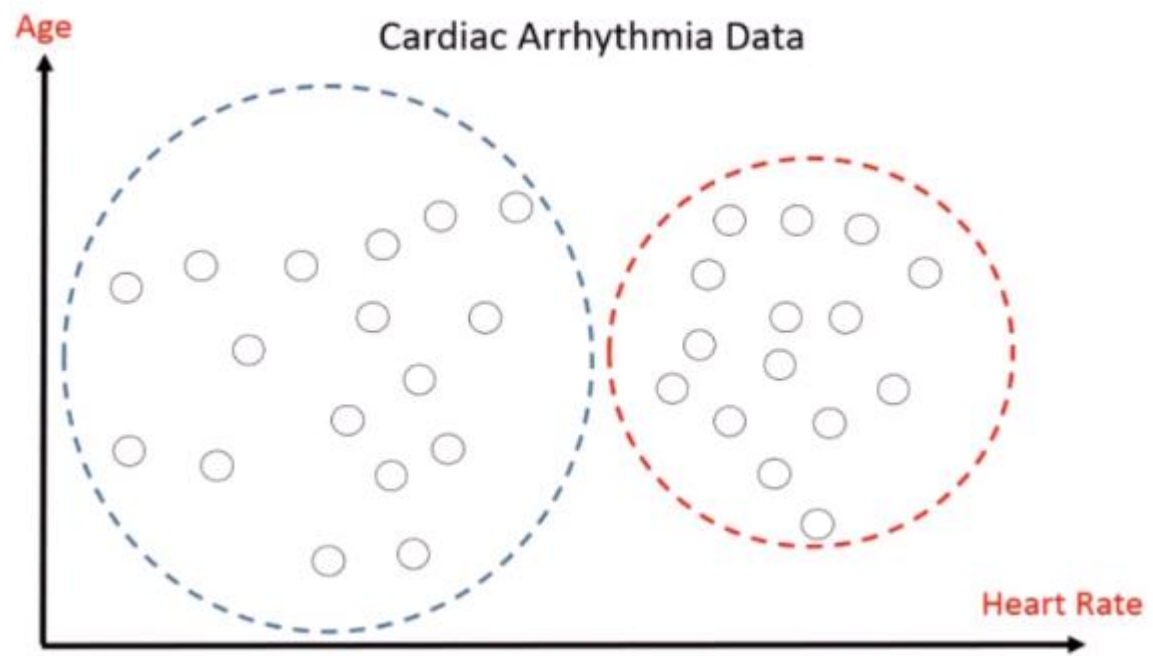
$$T = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d(x_i, x_j)$$

$$T = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \left(\sum_{C(j)=k} d(x_i, x_j) + \sum_{C(j) \neq k} d(x_i, x_j) \right)$$

$$T = W(C) + B(C)$$

Within
Cluster

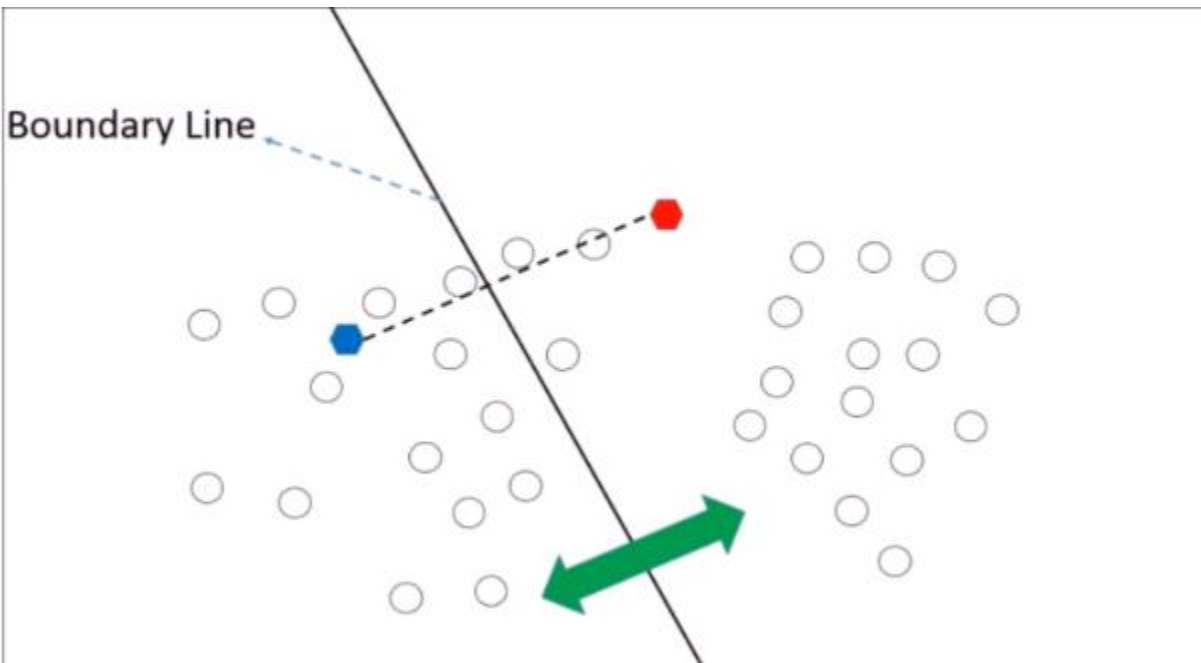
Between
Clusters



Initialization

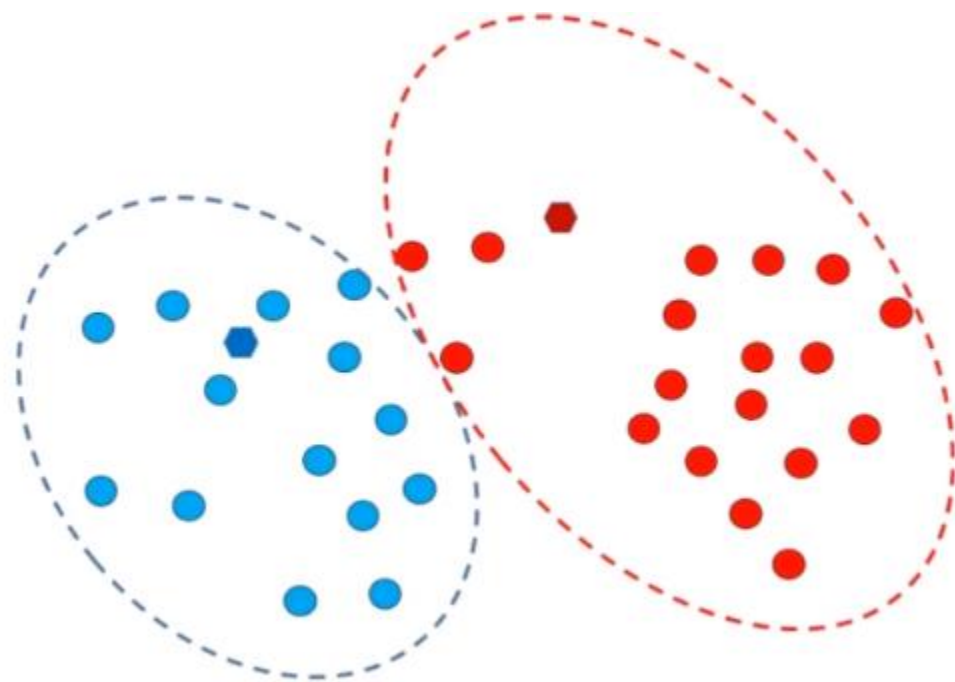
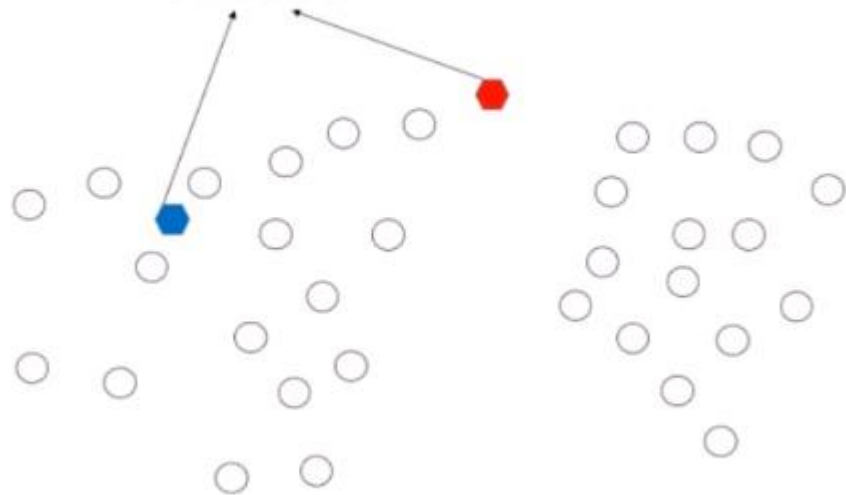


Boundary Line

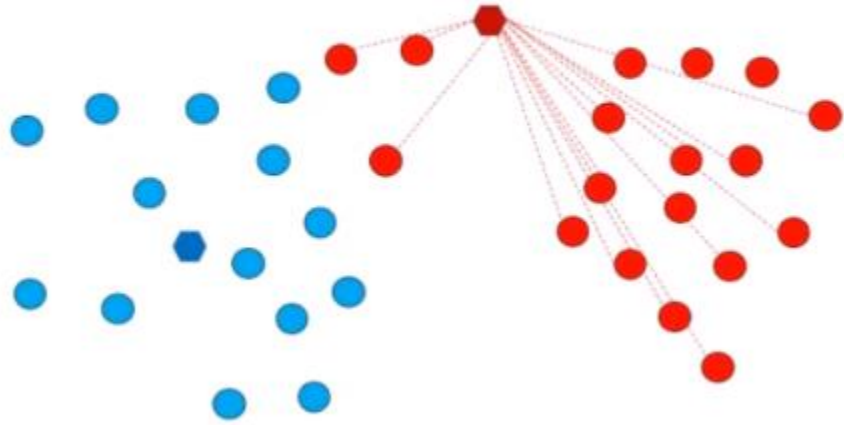


Centroids

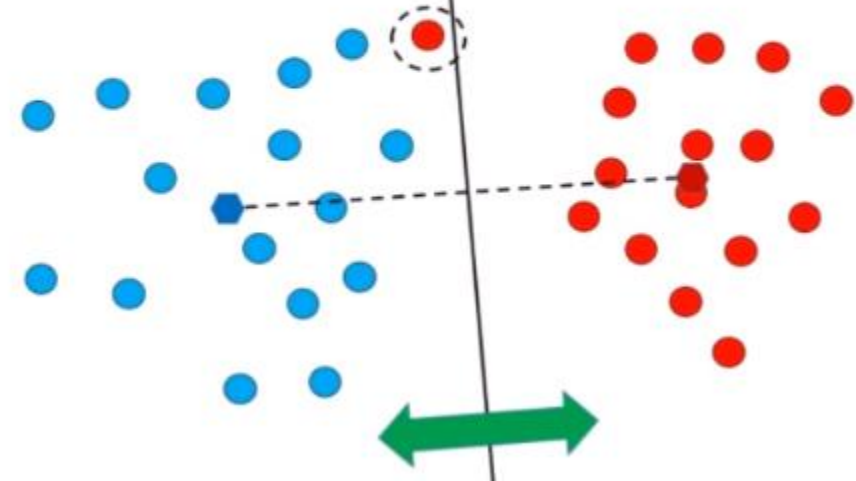
Initialization



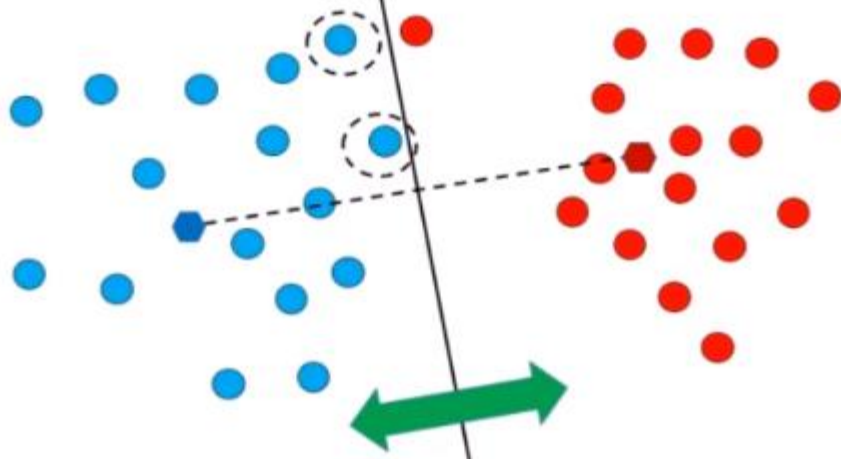
Iteration 1



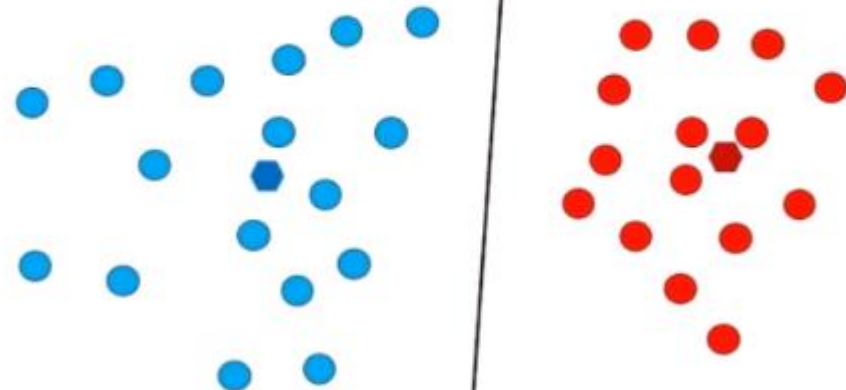
Iteration 2



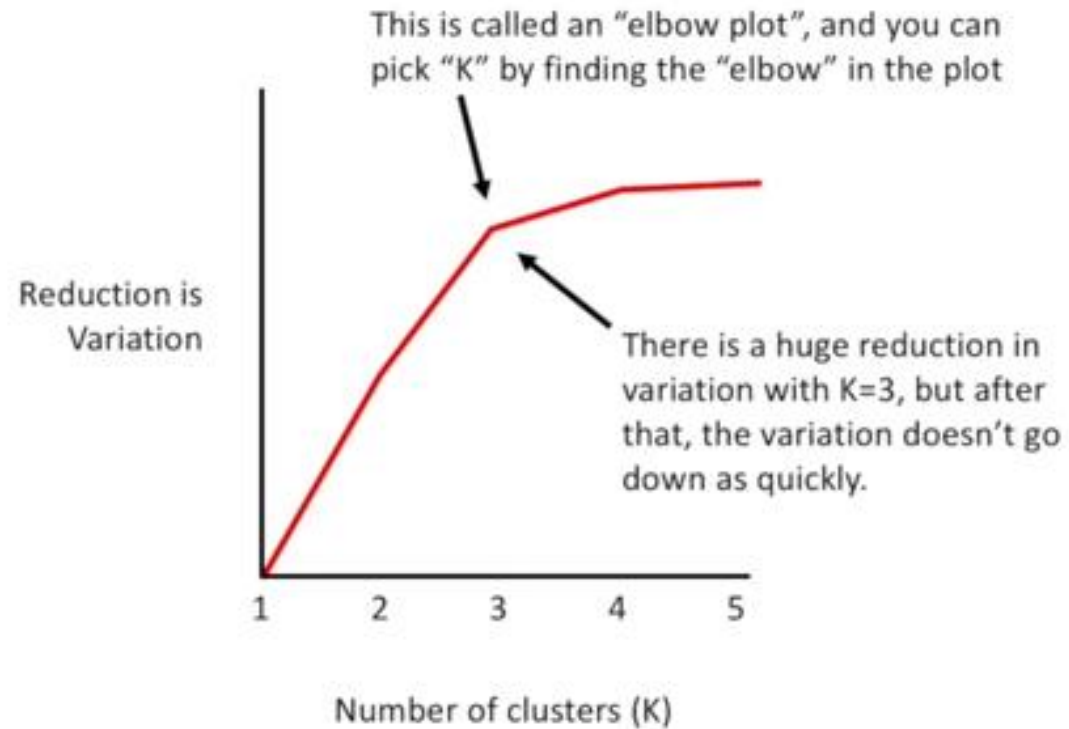
Iteration 1



Iterations
5,6,7...



How to Calculate K

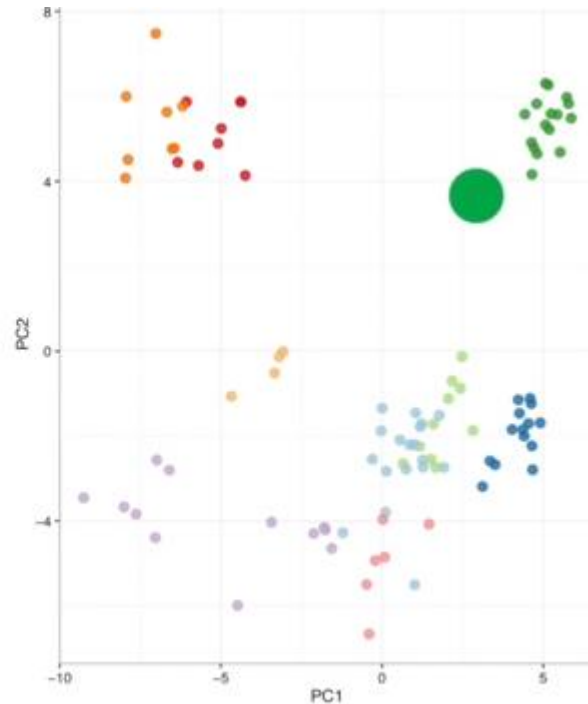


K-Nearest Neighbor (KNN)

In this case, the category is **GREEN**.

If $K=11$, we would use the 11 nearest neighbors.

In this case, the category is still **GREEN**.



If $K=11$ and the new cell is between two (or more) categories, we simply pick the category that "gets the most votes".

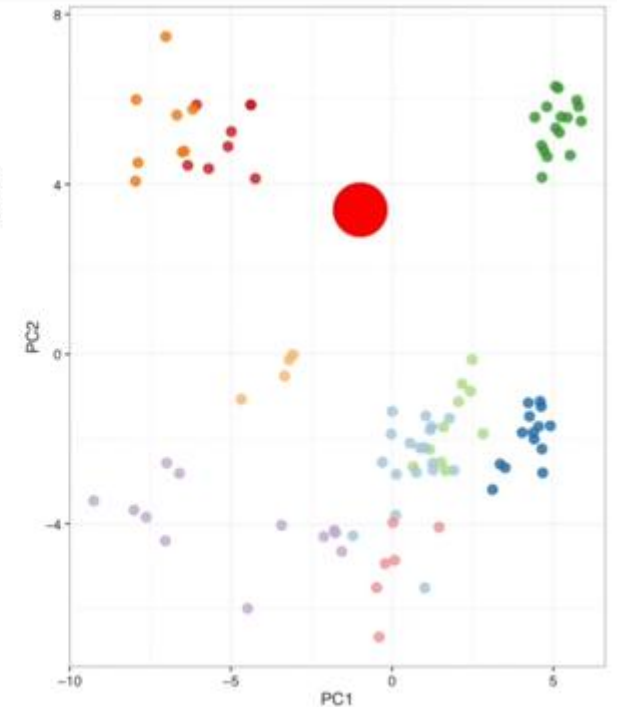
In this case....

7 nearest neighbors are **RED**.

3 nearest neighbors are **ORANGE**.

1 nearest neighbor is **GREEN**.

Since **RED** got the most votes, the final assignment is **RED**.



Hierarchical Clustering

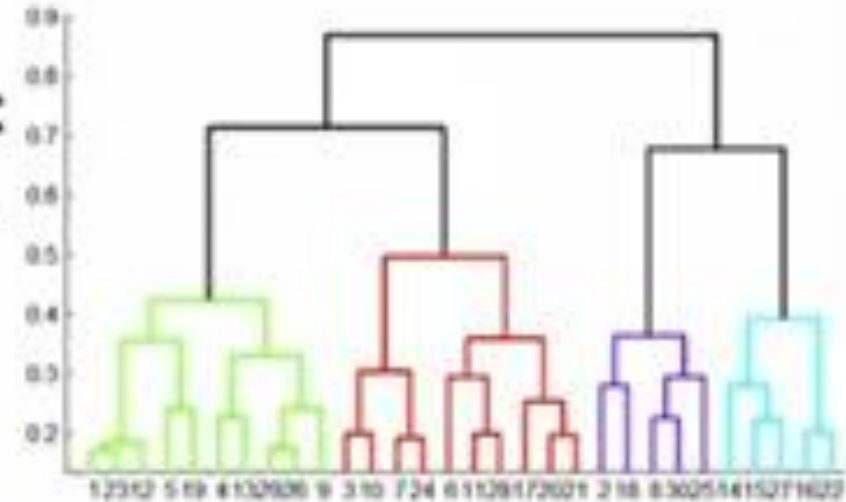
- **Hierarchical:**

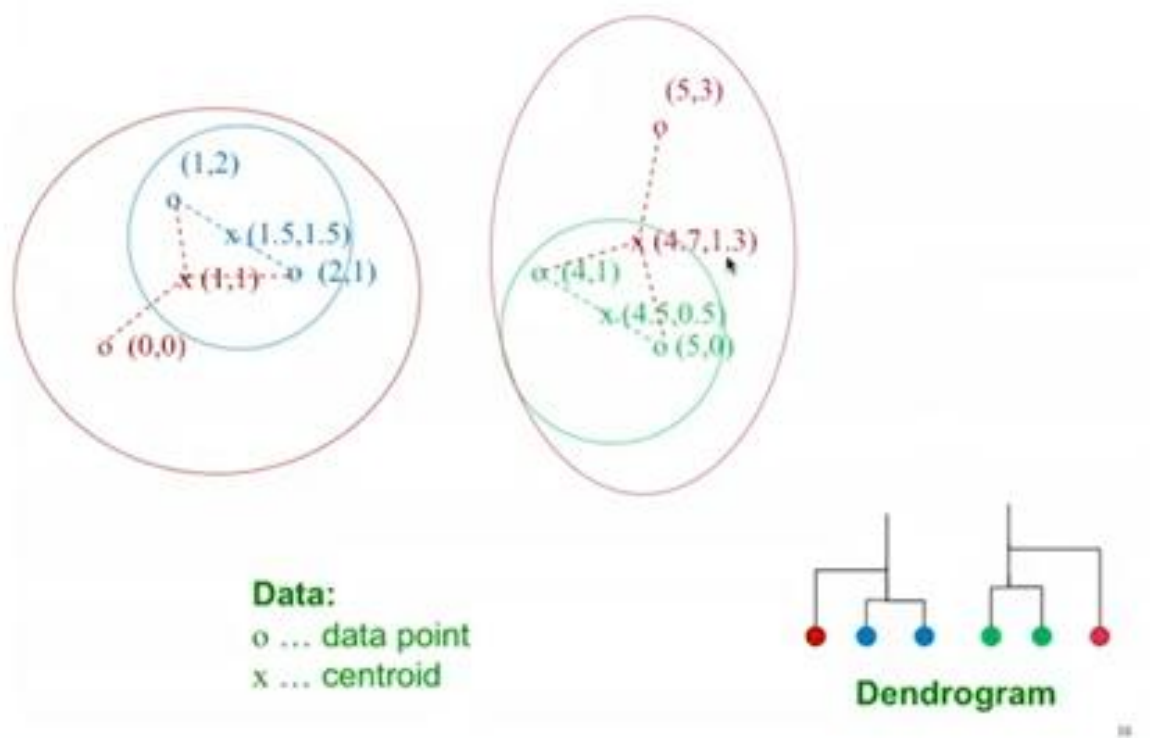
- **Agglomerative** (bottom up):

- Initially, each point is a cluster
 - Repeatedly combine the two “nearest” clusters into one

- **Divisive** (top down):

- Start with one cluster and recursively split it





Distance Functions

- Distance or similarity functions play a central role in all clustering algorithms.
- Numerous distance functions have been reported.
 - Euclidean Distance
 - Manhattan Distance
 - Minkowski Distance
 - Chebychev Distance
 - Cosine Distance