

**Scalable data**  $\Rightarrow$  Data can grow in size, systems or algo can still handle it efficiently without crashing.

## **K-Means** (Flat Clustering)

- Clustering  $\rightarrow$  organizing objects into groups whose members turn out to be similar.

Also called

data segmentation (form of learning by ~~observing~~ observations rather than by example).

- Approach  $\Rightarrow$  collect data / find clusters / Measure quality of clusters.

### • Cluster Analysis $\Rightarrow$ Requirements

- Scalability  $\rightarrow$  ~~We~~ need high scalable ~~data~~ algorithms which can process high data efficiently.

- Ability to deal with  $\rightarrow$  • Algo should perform for binary, categorical or ordinal data well.  
different type of attributes.  
+ noisy data

- Identify every shape of cluster  $\rightarrow$  • Most Algo tend to find only spherical clusters.

- High dimensionality  $\rightarrow$  • Algo should understand all dimension and then make clusters

- Clustering results are sensitive to <sup>input</sup> parameters. Selecting parameters is difficult.

### • Issues in clusters $\rightarrow$

- How many clusters  $\rightarrow$  • You can fix K before clustering
  - let number depend on some quality measure
  - Right choice depends on problem you want to solve.

Flat (Find all clusters at once)

- Reallocate objects to improve clustering.

Hierarchical (found new cluster by joining previous)

Agglomerative (cluster 1 and merge them)  
Divisive (all objects in 1 cluster and then split them)



- Hard clustering (Each object belongs to exactly one cluster) | Simple + Common.
- Soft clustering (Objects have "fuzzy" membership in multiple clusters)

- Overall quality of clustering is measured by  $f$   
 $f$  is closely related to measure of distance.

Primary goals →

- Low inter-cluster similarity
- High intra-cluster similarity
- Avoid very small or large clusters
- Above all focuses on internal criteria. Compare it with external criteria (like hand crafted reference clustering)
- Naive approach & Heuristics → • Trying all possible clusterings is not practical as number of possible clustering grows exponentially with data size.
- Use heuristic methods that provide a good solution

- Flat clustering (K-means) → • Every cluster is measured from the center.
- K-medoids or PAM (Partition Around Medoids)  
 (cluster representation is from actual data point)

- Number of clusters ( $K$ ) are defined in advance.
- Data points represented as unit vectors.



- Centroid of a cluster is defined as

$$\mu(A) = \frac{1}{m} \sum_{i=1}^m d_i$$

- RSS (Residual Sum of Squares) of cluster.

↓

- measures how far points in cluster are from centroid.

- goal in k-mean is to minimize RSS.

$$RSS(A) = \sum ||d_i - \mu(A)||^2$$

- Quality of clustering & minimizing RSS

How well data points are clustered?

■ Measured by summing the RSS of each cluster.

■ k-means minimize total RSS

- k-Mean Algo (Lloyd's Algo)

■ Randomly select k data points as initial clusters centers (seeds or centroids) + create empty clusters.

■ Assign 1 centroid to each cluster.

■ Iterate over each data points + assign each data point to cluster

■ Check clustering good if not then from create empty clusters.

- When to decide k-means cluster is good?

■ Small change since previous iteration

■ Max no of iterations reached.

■ Set threshold for RSS.

- k-means is fairly efficient (near fast) + often terminates at local optimum.

$O(nkt)$  where  $n$  = objects,  $k$  = clusters,  $t$  = iterations.

$$K, t \ll n.$$



Similar Approaches → ① K-medoids (use median)  
② Fuzzy c-mean (soft clustering)  
③ Model Based Clustering

- Disadvantages → • cluster number are decided before.  
• only work on those whose mean is

calculated.

• not suitable for categorical data.

• do not handle noisy data or outliers efficiently.

• do not identify non-convex clusters (for searching by separate nhi kr skaty).