

clustering: similarity is hard to judge in high-dimension space. { hierarchical
point assignment

1. Hierarchical Clustering. { represent cluster = euclidean case / non-euclidean case
determine "nearness":
when to stop combining

1. euclidean case: represent cluster by centroid measure by distance of centroid.
→ smallest max-dist / avg-dist / sum-square-dist

2. non-euclidean case: cluster / data point closest to other points, measure by distance of cluster.

cohesion: merge clusters whose union is most cohesive { diameter of merged cluster → max-dist in cluster
avg-dist between points in cluster
density-based approach: diameter / avg, then / # of point

3. complexity: n data point, at most $(n-1)$ step of merge, naive: $O(n^2)$ to form dist-matrix $\Rightarrow O(n^3)$

→ using priority queue: $O(n^2 \log n)$. → still too expensive.

4. when to stop?
b. measure: (1) sum square error (SSE): $\sum (\text{dist}^2(\text{point}, \text{it's centroid}))$
as k clusters

2. k-means Alg. 1. Step: (1) pick k points randomly, then place each point to cluster whose current centroid is nearest → update location of centroid after all points assigned → reassign → repeat
2. how to select k : avg-dist to centroid $\uparrow \rightarrow k$

3. BFR Alg. 1. assumption: (1) clusters are normally distributed around a centroid in Euclidean. → memory required: $O(k)$

2. step: read disk → clustering to k cluster → points in cluster summarized as μ and σ .
(DS: discard set). → second read → update μ and σ and dump points.

3. types of points: (1) DS (discard set): close enough to a centroid, to be summarized and chopped
(2) CS (compression set): points close together but not close to any existing centroid.
to be summarized but not assigned to cluster

(3) RS (retained set): isolated points waiting to be assigned to CS → outliers.

4. For each cluster, DS is summarized by (1) N : # of point, (2) vector SUM: $\text{sum}[i]$ is sum of coordinates of points in i th dimension, (3) vector SUMSQ: $\text{sumsq}[i]$ is sum of square of coordinate in i th dimension. $\text{Var}(X) = E(X^2) - E(X)^2$

5. (2d+1) values represent any size cluster in d dimension, centroid in i dimension = $\frac{\text{sum}_i}{N}$
var of a cluster's DS in i th dimension = $\frac{\text{sumsq}_i}{N} - \left(\frac{\text{sum}_i}{N}\right)^2$

b. Mahalanobis distance: normalized Euclidean dist from centroid

$u_i = \frac{x_i - c_i}{\sigma_i}$ $d(x, c) = \sqrt{\sum_{i=1}^d u_i^2} = \sqrt{\sum_{i=1}^d \left(\frac{x_i - c_i}{\sigma_i}\right)^2} \Rightarrow$ if clusters are normally distributed in d

dimension, then after transformation, one standard deviation = \sqrt{nd} if less than 500

7. When do two compressed set combined: calculate var of combined clusters ^{threshold}

10. CURIE Alg. - for arbitrary shape - use collection of representative points to represent cluster

1. problem of BFR / k-means: assumption & axes are fixed.

2. step. (1) Pass 1: pick random points \rightarrow cluster these points hierarchically \rightarrow for each cluster, pick

a sample of points as representative points (as dispersed as possible) \rightarrow then remove xx%.

toward the centroid. (2) pass 2: rescan dataset and place data point to "closest cluster"

3. BFR vs CURIE.

closest to one of representative points