



## cluster analysis

- maximize inter-cluster distance and minimize intra-cluster distance  
 ↳ between clusters ↳ inside the cluster elements

## types of clustering

- **partitional** = non-overlapping clusters
- **hierarchical** = tree like: traditional / non-traditional (size first  $\rightarrow$  merge 1's first then 2's)
- **non-exclusive** = a point may belong more than one cluster
  - $\hookrightarrow$  fuzzy clustering  $\Rightarrow x \rightarrow 0.3$  cluster A,  $0.4$  cluster B,  $0.3$  cluster C

## types of clusters

- **well-separated** = any point in a cluster is closer to every point in the cluster than to any point not in the cluster   $a > b \rightarrow$  not-well separated
- **prototype-based** = "closer to the prototype / center of its cluster than others centers  
 ↳ centroid, medoid: most representative
- **contiguous (nearest neighbor / transitive)** = a point is closer to one or more other points in the cluster 

→ cluster A
→ cluster B
- **density-based** = dense regions are clusters, separated by low-density regions  
 used when clusters are irregular, where there are outliers
- **described by an objective function** = find clusters that min/maximize an objective function

## clustering algorithms

- **k-means** = complexity  $O(\text{\#points} \times \text{\#clusters} \times \text{\#iterations} \times \text{\#attributes})$   $O(N)$ 
  - ↳ repeat until sum of squared error - euclidean func reaches minimum
  - ↳ **k-means++** = select initial centroids = next centroid is the point having max distance from the nearest centroid (most probably - logn converge guarantee)
  - ↳ **bisecting kmeans** = split some of points into two, choose one → like hierarchical clustering
- **hierarchical clustering** = complexity  $O(N^2)$  space → proximity / similarity matrix,  $O(N^3)$  time → worst than k-means
  - ↳ dendrogram, agglomerative (merge,  $n \rightarrow 1$  cluster) / divisive (split at each step,  $1 \rightarrow n$  clusters)
  - ↳ **single linkage** = min dist → sensitive to noise
  - ↳ **complete linkage** = among max distances choose min one
  - ↳ **average linkage** = avg of distances between all datapoints among two clusters
  - ↳ **ward's method** = sum of squared error, variance → like avg dist, but do not divide and squared
- **DBSCAN** = density based, merge core points until no left, assign borders to clusters at the end
  - ↳ **core point** = if the point has at least  $N(\text{minN})$  points within circle
  - ↳ **border point** = neighbor of a core point
  - ↳ **noise point** = neither core nor border point

measurement of the goodness of the clusters

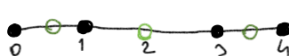
- **cluster cohesion** = how closely related are objects in a cluster  $\rightarrow$  SSE
- **cluster separation** = how distinct/well separated clusters are  $\rightarrow$  SS<sub>B</sub> between clusters
- **silhouette coefficient** = considers both cohesion and separation  
 $\hookrightarrow s = \frac{b-a}{\max(a,b)}$  ( $-1 \leq s \leq 1$ ),  $a$  = avg dist of a point to others in its cluster,  $b$  = avg dist in another cluster  
 $\rightarrow$  better (pointing to 1)
- **correlation** = measurement of the corr between proximity matrix and ideal similarity matrix  
 $\hookrightarrow$  must be high in magnitude (can be negative too)  $\rightarrow 1$ : if same clusters



$$k = 1$$

$$SSE = (0-2)^2 + (1-2)^2 + (3-2)^2 + (4-2)^2 = 10 \quad \} + 10$$

$$SS_B = 4(3-3)^2 = 0$$



$k = 2$

$$SSB = 2 \times (0 - 0.5)^2 + 2 \times (3 - 0.5)^2 = 9$$

$$SS_B = 2 \times (2 - 0.5)^2 + 2 \times (3.5 - 2)^2 = 9$$

$\rightarrow x(2-0.5) + \dots = x(0.5)$   
 $\downarrow$   $\downarrow$   $\downarrow$   
 in the centroid cluster center  
 cluster

$$SS_E + SS_B = \text{constant}$$