

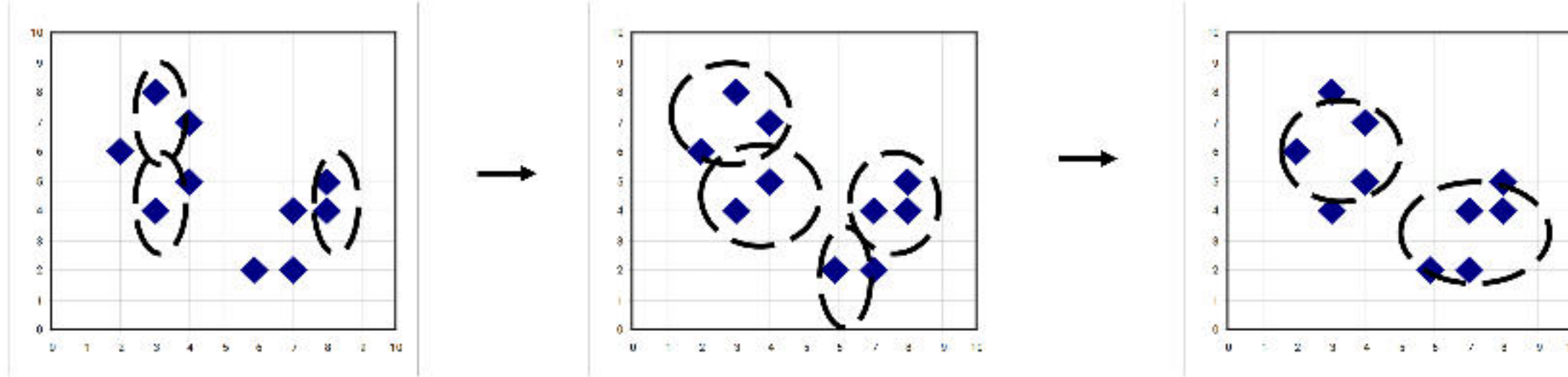


Agglomerative Clustering Algorithms

Agglomerative Clustering Algorithm

□ AGNES (AGglomerative NESTing) (Kaufmann and Rousseeuw, 1990)

- Use the single-link method and the dissimilarity matrix
- Continuously merge nodes that have the least dissimilarity
- Eventually all nodes belong to the same cluster



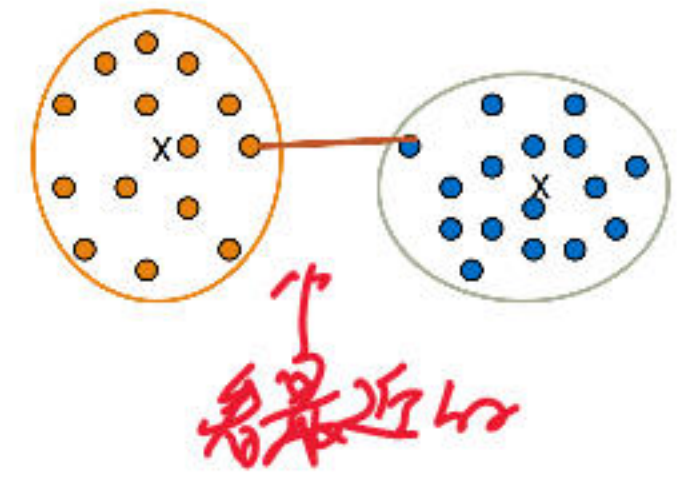
□ Agglomerative clustering varies on different similarity measures among clusters

- Single link (nearest neighbor)
- Average link (group average)
- Complete link (diameter)
- Centroid link (centroid similarity)

Single Link vs. Complete Link in Hierarchical Clustering

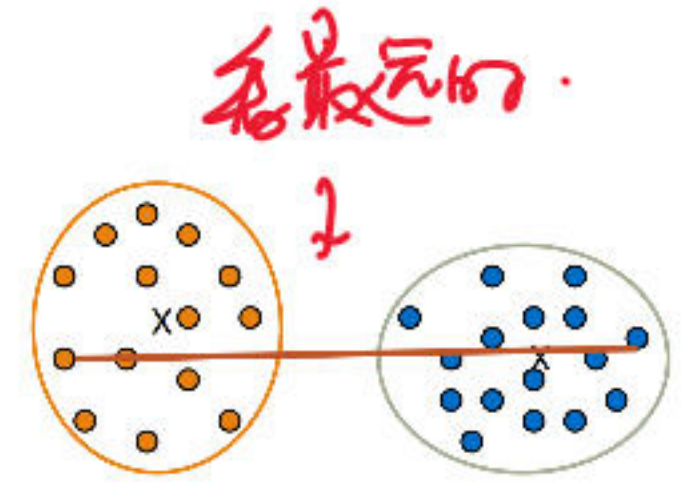
Single link (nearest neighbor)

- The similarity between two clusters is the similarity between their most similar (nearest neighbor) members
- Local similarity-based: Emphasizing more on close regions, ignoring the overall structure of the cluster
- Capable of clustering non-elliptical shaped group of objects
- Sensitive to noise and outliers



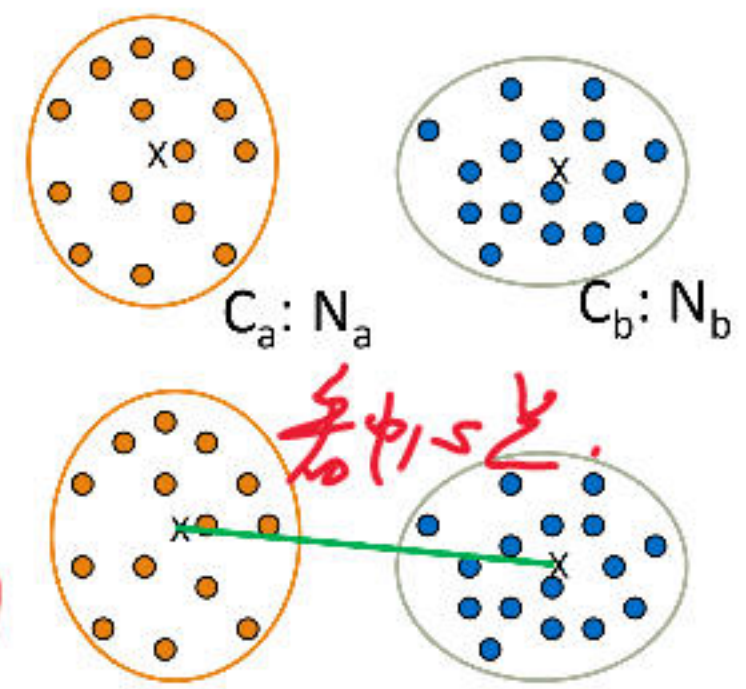
Complete link (diameter)

- The similarity between two clusters is the similarity between their most dissimilar members
- Merge two clusters to form one with the smallest diameter
- Nonlocal in behavior, obtaining compact shaped clusters
- Sensitive to outliers



Agglomerative Clustering: Average vs. Centroid Links

- Agglomerative clustering with average link 开销太大.
 - Average link:** The average distance between an element in one cluster and an element in the other (i.e., all pairs in two clusters)
 - Expensive to compute



- Agglomerative clustering with centroid link
 - Centroid link:** The distance between the centroids of two clusters

- Group Averaged Agglomerative Clustering (GAAC)** 考虑权重
 - Let two clusters C_a and C_b be merged into $C_{a \cup b}$. The new centroid is:

N_a is the cardinality of cluster C_a , and c_a is the centroid of C_a

- The similarity measure for GAAC is the average of their distances

$$c_{a \cup b} = \frac{N_a c_a + N_b c_b}{N_a + N_b}$$

- Agglomerative clustering with **Ward's criterion**

- Ward's criterion:** The increase in the value of the SSE criterion for the clustering obtained by merging them into $C_a \cup C_b$:

$$W(C_{a \cup b}, c_{a \cup b}) - W(C, c) = \frac{N_a N_b}{N_a + N_b} d(c_a, c_b)$$