

Introduction

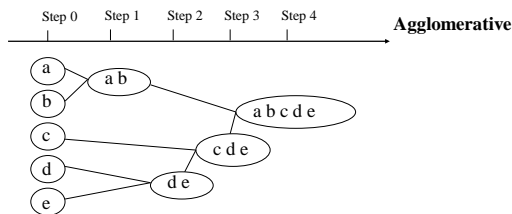
Agglomerative Clustering

HAC

- Hierarchical Clustering Approach
 - A typical clustering analysis approach by partitioning data set **sequentially**
 - Construct nested partitions layer by layer by grouping objects into a tree of clusters (**without the need to know the number of clusters in advance**)
 - Use (generalised) distance matrix as clustering criteria

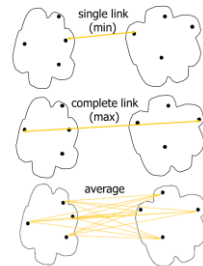
Introduction

• Example



Cluster Distance Measures

- **Single link:** smallest distance between an element in one cluster and an element in the other, i.e., $d(C_i, C_j) = \min\{d(x_{ip}, x_{jq})\}$
- **Complete link:** largest distance between an element in one cluster and an element in the other, i.e., $d(C_i, C_j) = \max\{d(x_{ip}, x_{jq})\}$
- **Average:** avg distance between elements in one cluster and elements in the other, i.e., $d(C_i, C_j) = \text{avg}\{d(x_{ip}, x_{jq})\}$



Cluster Distance Measures

Example: Given a data set of five objects characterised by a single continuous feature, assume that there are two clusters: $C_1: \{a, b\}$ and $C_2: \{c, d, e\}$.

	a	b	c	d	e
Feature	1	2	4	5	6

1. Calculate the distance matrix

	a	b	c	d	e
a	0				
b	1	0			
c	3	2	0		
d	4	3	1	0	
e	5	4	2	1	0

2. Calculate three cluster distances between C_1 and C_2 .

$$\text{dist}(C_1, C_2) = \min\{d(a, c), d(a, d), d(a, e), d(b, c), d(b, d), d(b, e)\} = \min\{3, 4, 5, 2, 3, 4\} = 2$$

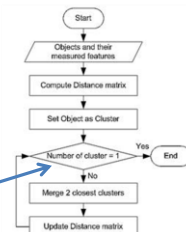
$$\text{dist}(C_1, C_2) = \max\{d(a, c), d(a, d), d(a, e), d(b, c), d(b, d), d(b, e)\} = \max\{3, 4, 5, 2, 3, 4\} = 5$$

$$\text{dist}(C_1, C_2) = \frac{d(a, c) + d(a, d) + d(a, e) + d(b, c) + d(b, d) + d(b, e)}{6} = \frac{3 + 4 + 5 + 2 + 3 + 4}{6} = \frac{21}{6} = 3.5$$

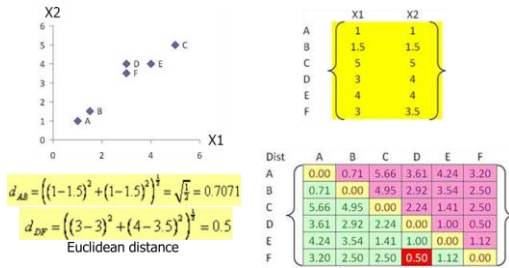
HAC-Algorithm

- The *Agglomerative* algorithm is carried out in three steps:

- 1) Convert all object features into a distance matrix
- 2) Set each object as a cluster (thus if we have N objects, we will have N clusters at the beginning)
- 3) Repeat until number of cluster is one (or known # of clusters)
 - Merge two closest clusters
 - Update "distance matrix"

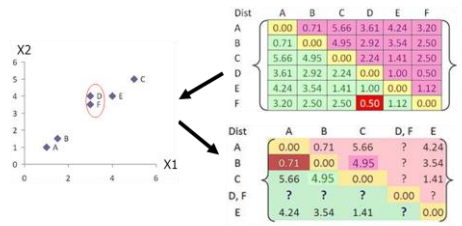


Example

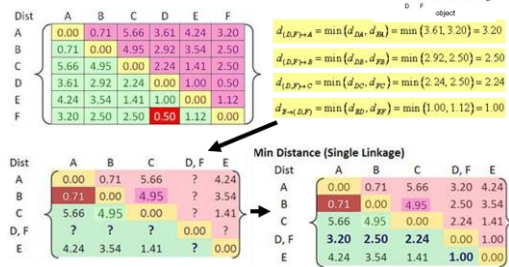


Example

- Merge two closest clusters (iteration 1)

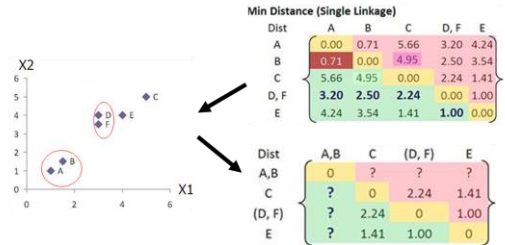


Example



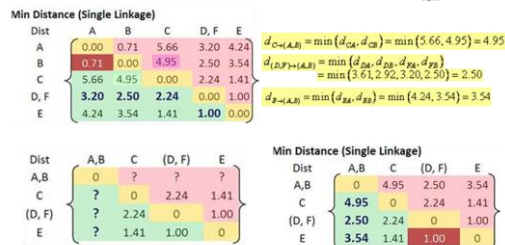
Example

- Merge two closest clusters (iteration 2)



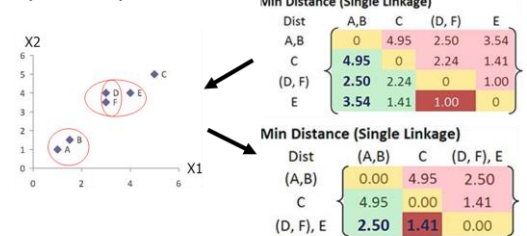
Example

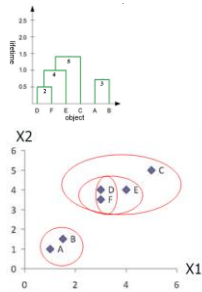
- Update distance matrix (iteration 2)



Example

- Merge two closest clusters/update distance (iteration 3)





Example

Min Distance (Single Linkage)

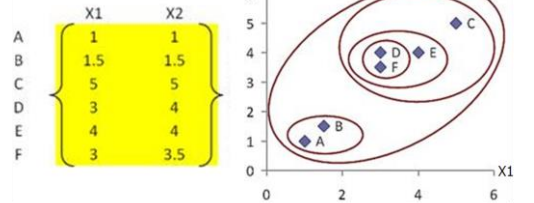
Dist	(A,B)	C	(D,F), E
(A,B)	0.00	4.95	2.50
C	4.95	0.00	1.41
(D,F), E	2.50	1.41	0.00

Min Distance (Single Linkage)

Dist	(A,B)	((D,F), E), C
(A,B)	0.00	2.50
((D,F), E), C	2.50	0.00

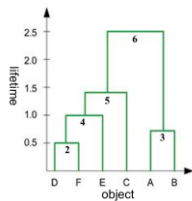
Example

- Final result (meeting termination condition)



Lifetime

Lifetime vs K-cluster Lifetime



Lifetime

The distance between that a cluster is created and that it disappears (merges with other clusters during clustering).
e.g. lifetime of A, B, C, D, E and F are 0.71, 0.71, 1.41, 0.50, 1.00 and 0.50, respectively, the life time of (A, B) is $2.50 - 0.71 = 1.79$,

K-cluster Lifetime

The distance from that K clusters emerge to that K clusters vanish (due to the reduction to K-1 clusters).
e.g.

5-cluster lifetime is $0.71 - 0.50 = 0.21$

4-cluster lifetime is $1.00 - 0.71 = 0.29$

3-cluster lifetime is $1.41 - 1.00 = 0.41$

2-cluster lifetime is $2.50 - 1.41 = 1.09$

Conclusion

- How to determine the number of clusters
 - If the number of clusters known, termination condition is given!
 - The *K-cluster lifetime* as the range of threshold value on the dendrogram tree that leads to the identification of K clusters
 - Heuristic rule: cut a dendrogram tree with maximum life time to find a "proper" K
- Major weakness of agglomerative clustering methods
 - Can never undo what was done previously
 - Sensitive to cluster distance measures and noise/outliers
 - Less efficient: $O(n^2 \log n)$, where n is the number of total objects