



# **Text Clustering**

## **- Similarity based Methods**

(CPEG 657: Search and Data Mining)

Hui Fang

Department of Electrical and Computer Engineering  
University of Delaware

## Similarity-based Clustering

- Given a similarity function that can be used to measure similarity between two documents,
- Find a partition to
  - Maximize intra-cluster similarity
  - Minimize inter-cluster similarity
- How to find the partition?
  - Hierarchically group similar documents
  - Search by starting at a random partition

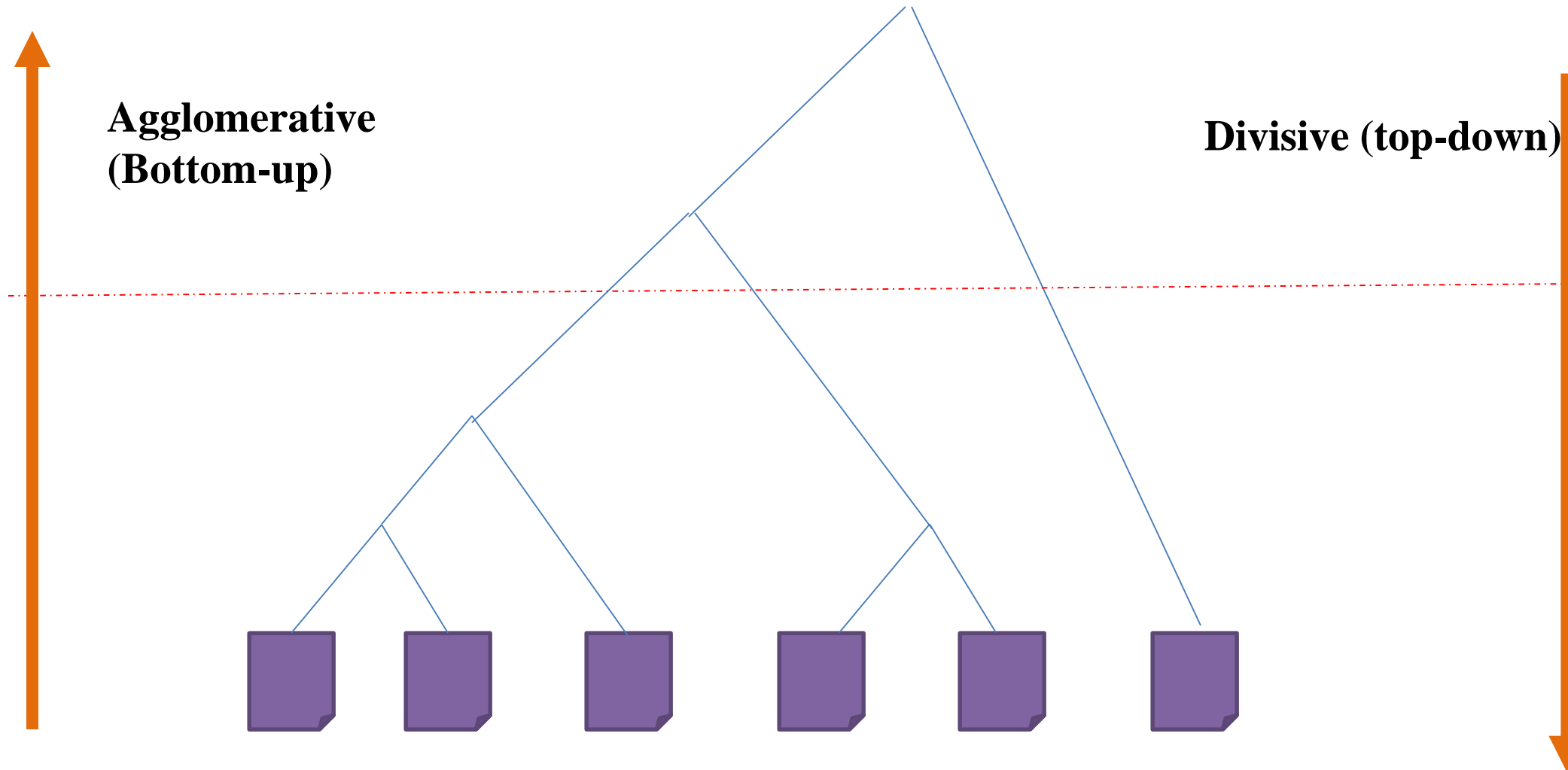
**Method 1 (Similarity-based):**

**Agglomerative Hierarchical Clustering**

# Agglomerative Hierarchical Clustering

- Given a similarity function to measure similarity between two documents
- Gradually group similar documents together in a bottom-up fashion
- Stop when some stopping criterion is met

# Agglomerative Hierarchical Clustering



## How to Compute Group Similarity?



Given two groups  $g1$  and  $g2$ ,

**Single-link algorithm:**  $s(g1, g2)$  = similarity of the **closest** pair

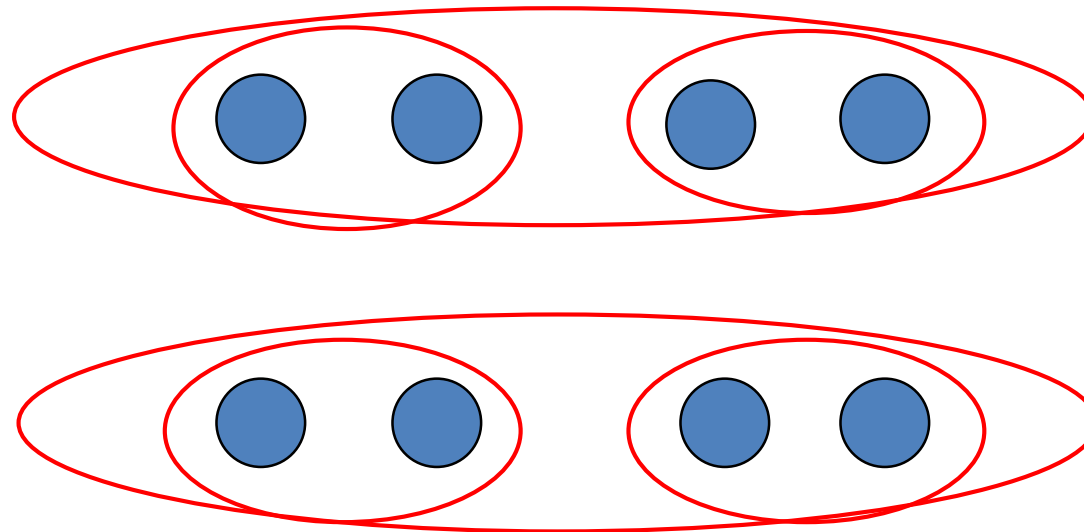
**complete-link algorithm:**  $s(g1, g2)$  = similarity of the **farthest** pair

**average-link algorithm:**  $s(g1, g2)$  = **average** of similarity of all pairs



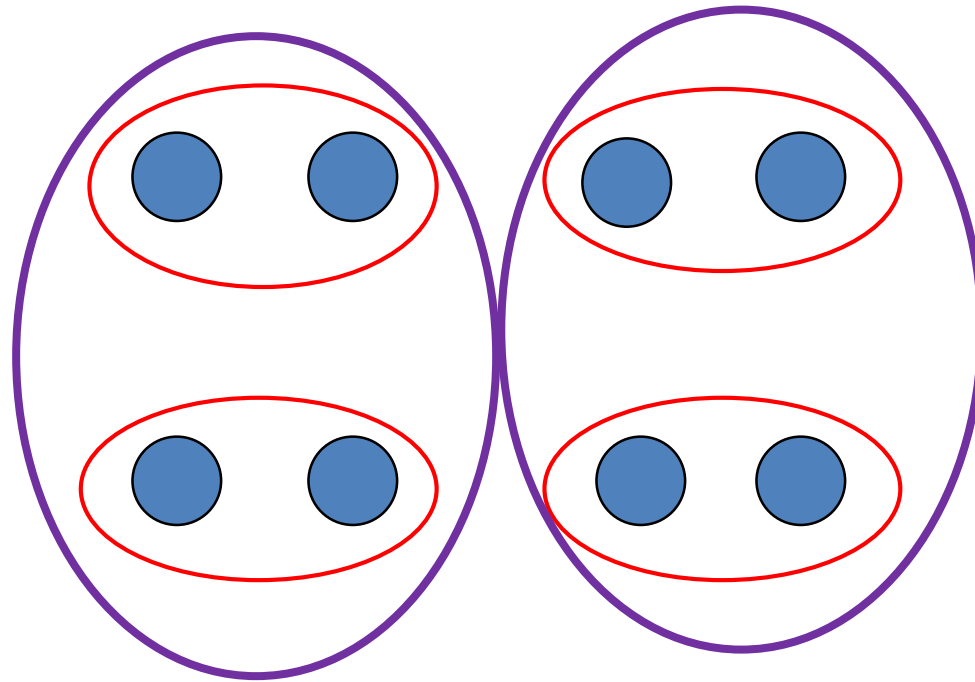
# Three Methods Illustrated

## Single Link Example





## Complete Link Example



## Comparison of the Three Methods

- Single-link
  - Loose clusters
  - Individual decision, sensitive to outliers
- Complete-link
  - Tight clusters
  - Individual decision, sensitive to outliers
- Average-link
  - Group decision, insensitive to outliers



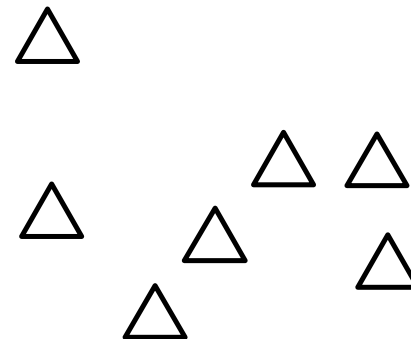
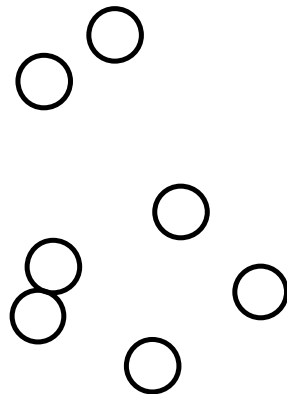
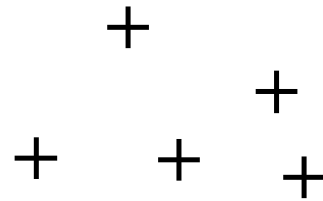
**Method 2 (similarity-based):**

**K-Means**

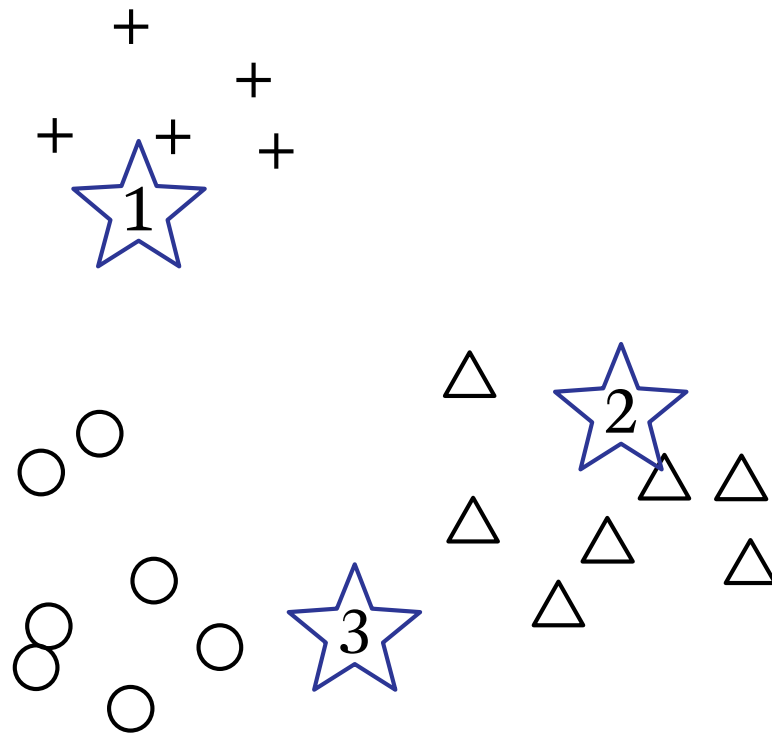
## K-Means Clustering

- Start with  $k$  randomly selected documents, and assume they are the centroids of  $k$  clusters
- Assign every document to a cluster whose centroid is the closest to the document
- Re-compute the centroid for each cluster
- Repeat this process until the centroids converge

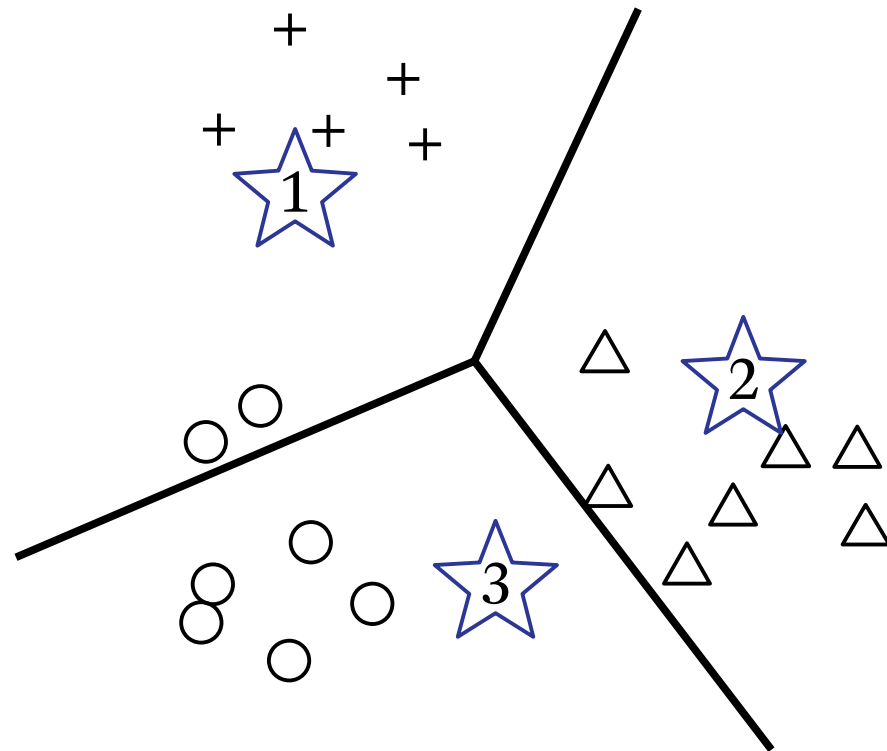
# An Example of K-Means



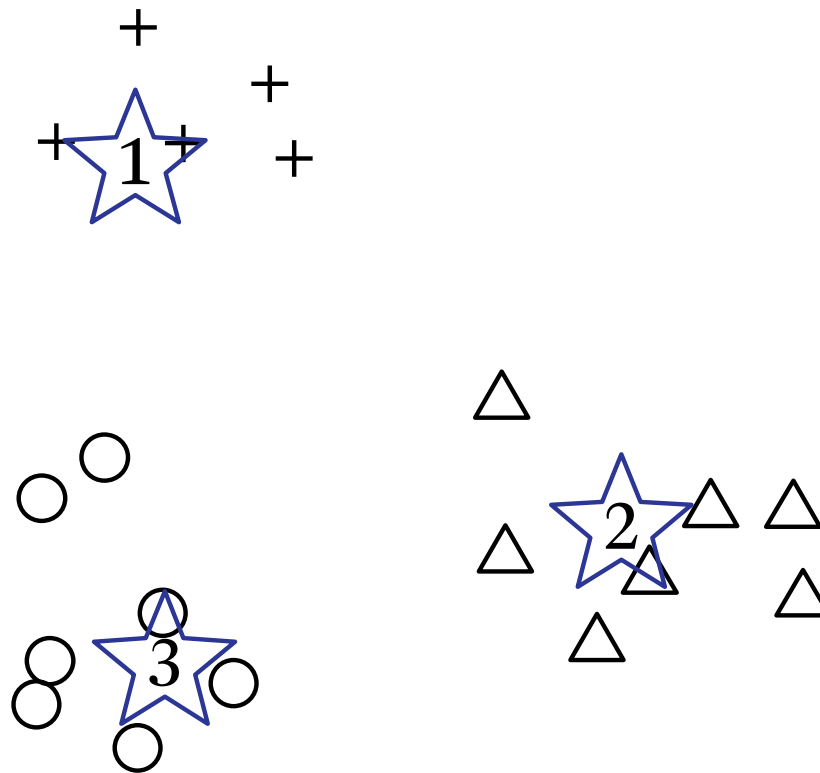
# An Example of K-Means



# An Example of K-Means

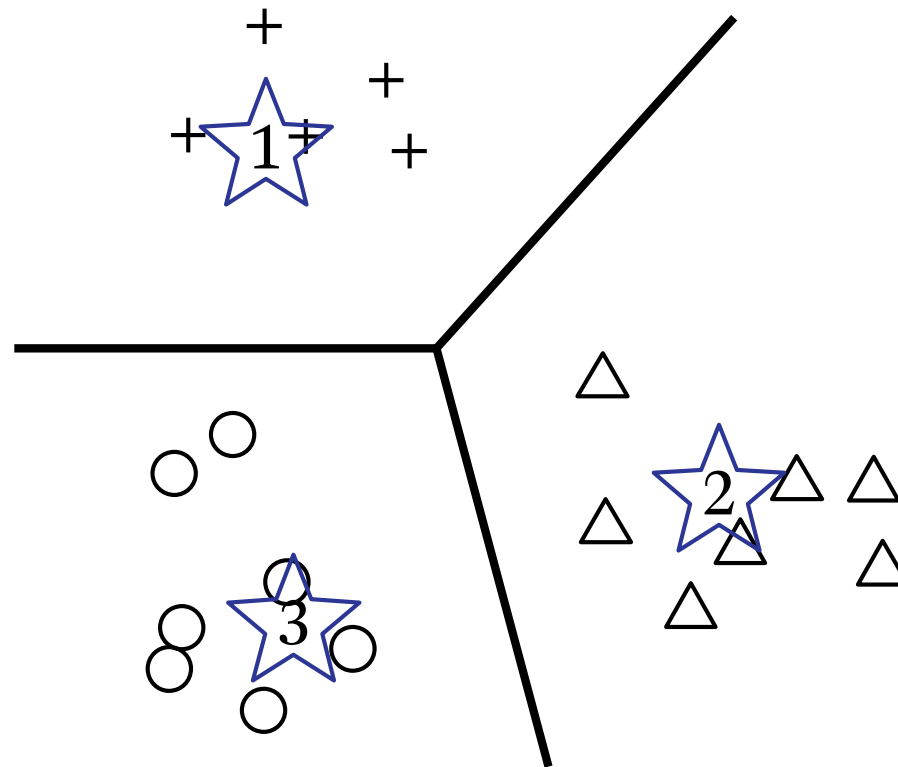


# An Example of K-Means





# An Example of K-Means



**Results can vary based on random seed selection.**

