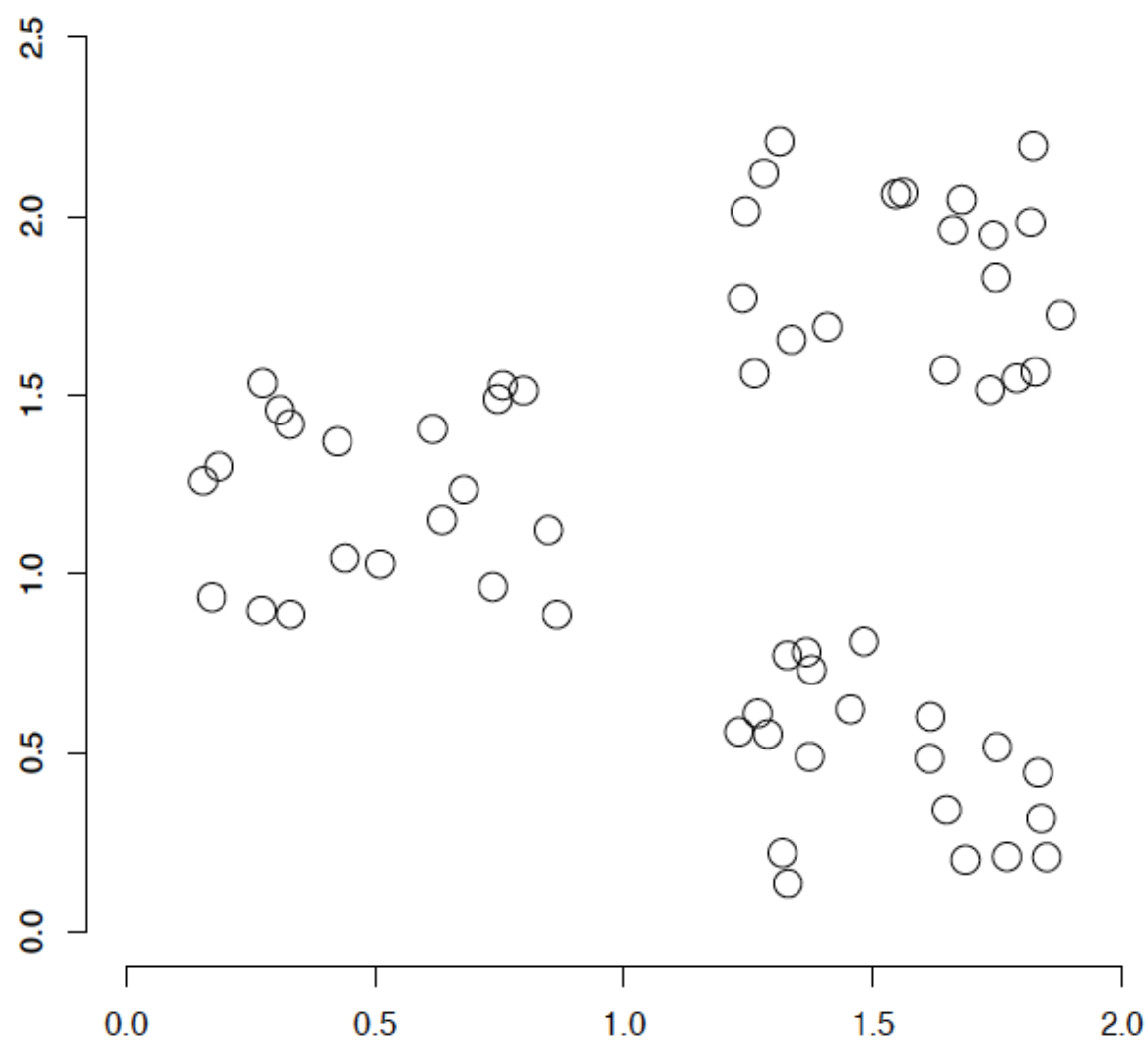


# مبانی بازیابی اطلاعات و جستجوی وب

۱۰- خوشه بندی

# What is clustering?

- **Clustering**: the process of grouping a set of objects into classes of similar objects
  - Documents within a cluster should be similar.
  - Documents from different clusters should be dissimilar.
- The commonest form of *unsupervised learning*



## Applications of clustering in IR

- Whole corpus analysis/navigation
  - Better user interface: search without typing
- For improving recall in search applications
  - Better search results
- For better navigation of search results
  - Effective “user recall” will be higher
- For speeding up vector space retrieval
  - Cluster-based retrieval gives faster search

# Issues for clustering

- Representation for clustering
  - Document representation
  - Need a notion of similarity/distance (cosine, Euclidean)
- How many clusters?
  - Fixed a priori?
  - Completely data driven?

# Clustering Algorithms

- Flat algorithms
  - Usually start with a random (partial) partitioning
  - Refine it iteratively
    - $K$  means clustering
- Hierarchical algorithms
  - Bottom-up, agglomerative
  - (Top-down, divisive)

# K-Means

- Assumes documents are real-valued vectors.
- Clusters based on *centroids* (aka the *center of gravity* or mean) of points in a cluster,  $c$ :

$$\vec{\mu}(c) = \frac{1}{|c|} \sum_{\vec{x} \in c} \vec{x}$$

- Reassignment of instances to clusters is based on distance to the current cluster centroids.
- Minimizing the Cost function:

$$J = \sum_x \|x - \mu(c_x)\|^2$$

## K-Means Algorithm

Select  $K$  random docs  $\{s_1, s_2, \dots, s_K\}$  as seeds.

Until clustering *converges* (or other stopping criterion):

For each doc  $d_i$ :

Assign  $d_i$  to the cluster  $c_j$  such that  $\text{dist}(x_i, s_j)$  is minimal.

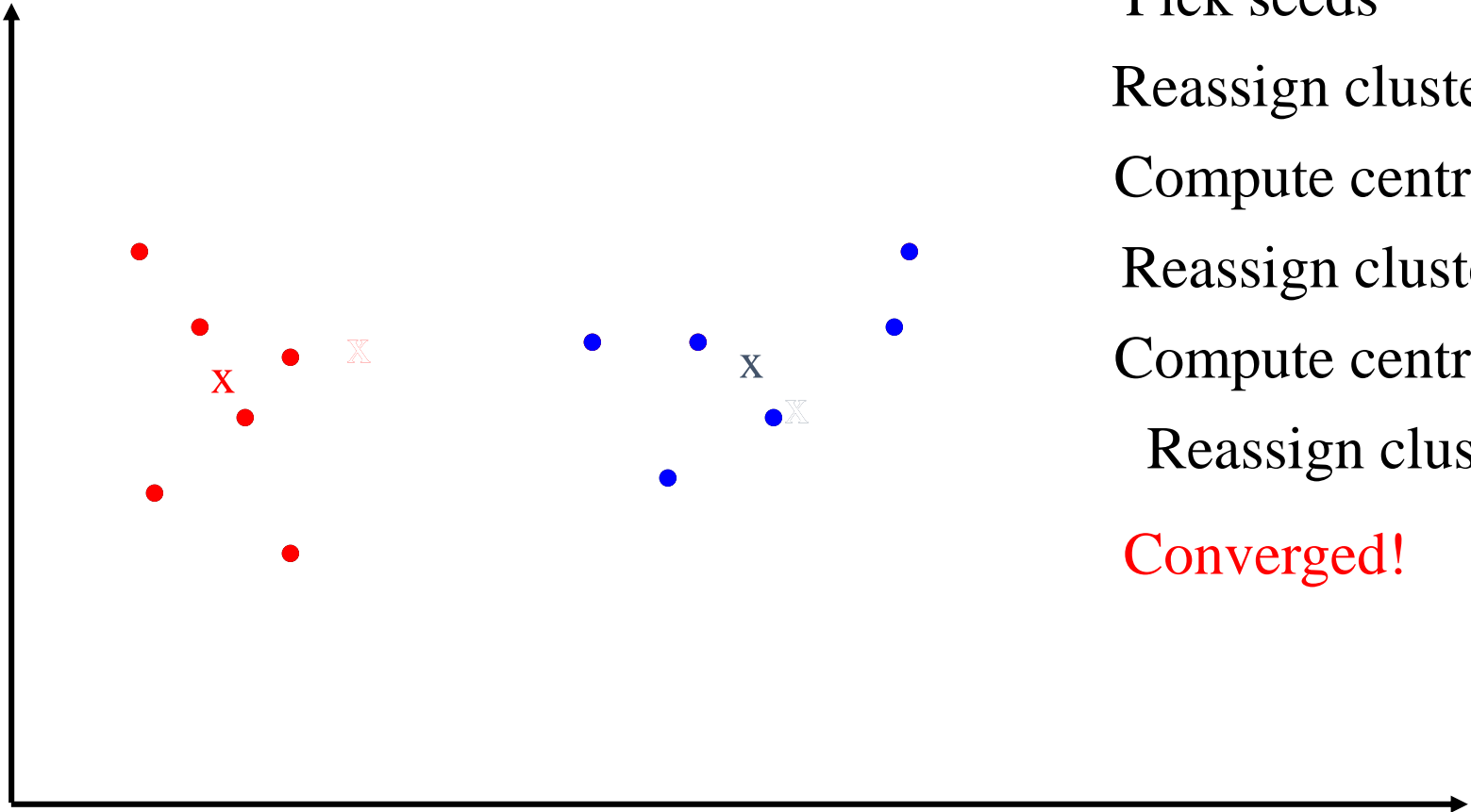
*(Next, update the seeds to the centroid of each cluster)*

For each cluster  $c_j$

$$s_j = \mu(c_j)$$



# K Means Example ( $K=2$ )



Pick seeds

Reassign clusters

Compute centroids

Reassign clusters

Compute centroids

Reassign clusters

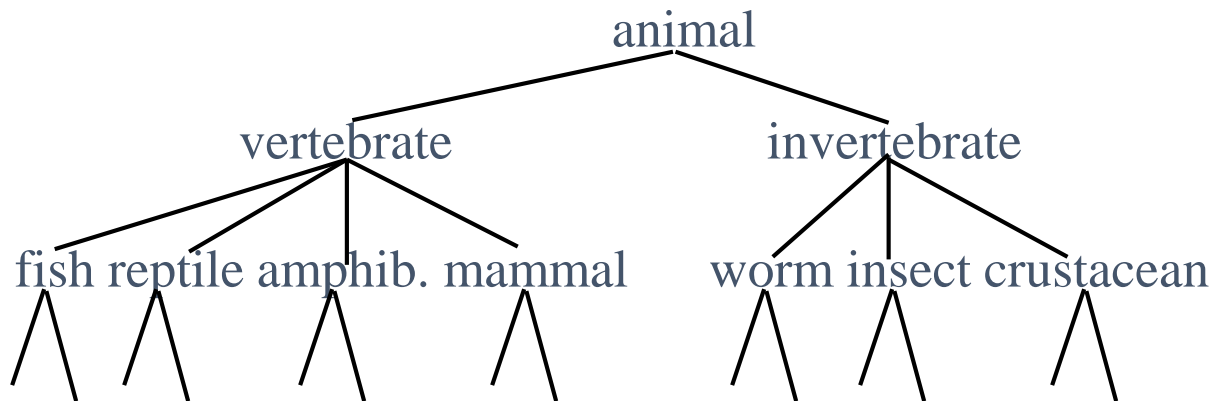
**Converged!**

## Termination conditions

- Several possibilities, e.g.,
  - A fixed number of iterations.
  - Doc partition unchanged.
  - Centroid positions don't change.

# Hierarchical Clustering

- Build a tree-based hierarchical taxonomy (*dendrogram*) from a set of documents.



- Clustering obtained by cutting the dendrogram at a desired level: each connected component forms a cluster.

## Hierarchical Agglomerative Clustering (HAC)

- Starts with each doc in a separate cluster
  - then repeatedly joins the closest pair of clusters, until there is only one cluster.
- The history of merging forms a binary tree or hierarchy.

## *Closest pair* of clusters

- Many variants to defining closest pair of clusters
- **Single-link**
  - Similarity of the *most* cosine-similar (single-link)
- **Complete-link**
  - Similarity of the “furthest” points, the *least* cosine-similar
- **Centroid**
  - Clusters whose centroids (centers of gravity) are the most cosine-similar
- **Average-link**
  - Average cosine between pairs of elements

# What Is A Good Clustering?

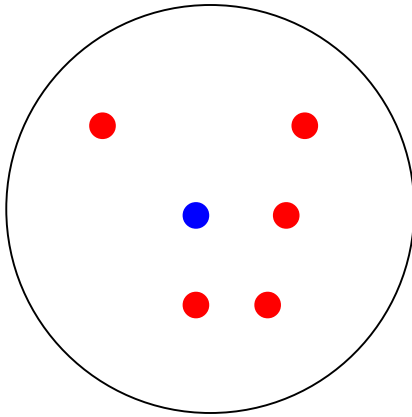
- Internal criterion
- External criterion

# External Evaluation of Cluster Quality

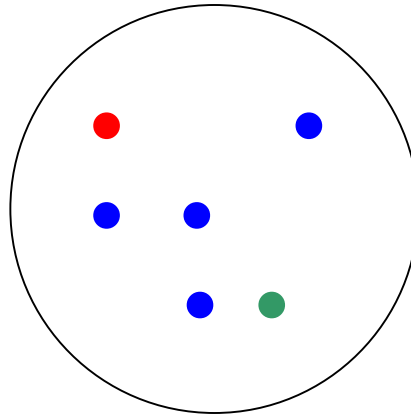
- Simple measure: purity, the ratio between the dominant class in the cluster  $\pi_i$  and the size of cluster  $\omega_i$

$$Purity(\omega_i) = \frac{1}{n_i} \max_j (n_{ij}) \quad j \in C$$

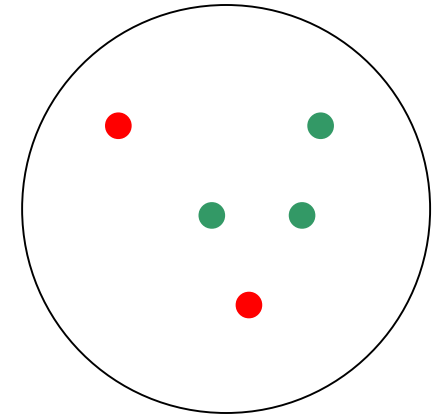
# Purity example



Cluster I



Cluster II



Cluster III

Cluster I: Purity =  $1/6 (\max(5, 1, 0)) = 5/6$

Cluster II: Purity =  $1/6 (\max(1, 4, 1)) = 4/6$

Cluster III: Purity =  $1/5 (\max(2, 0, 3)) = 3/5$



Rand Index measures between pair decisions. Here  $RI = 0.68$

Number of points	Same Cluster in clustering	Different Clusters in clustering
Same class in ground truth	20	24
Different classes in ground truth	20	72

# Rand index and Cluster F-measure

$$RI = \frac{A + D}{A + B + C + D}$$

■ فصل ۱۶ و ۱۷ کتاب An introduction to information retrieval