



Dr. George Karraz, Ph. D.

# Artificial Intelligence

## Lecture XI

# Clustering

**Dr. George Karraz, Ph.D.**

# Today's Topic: Clustering

---

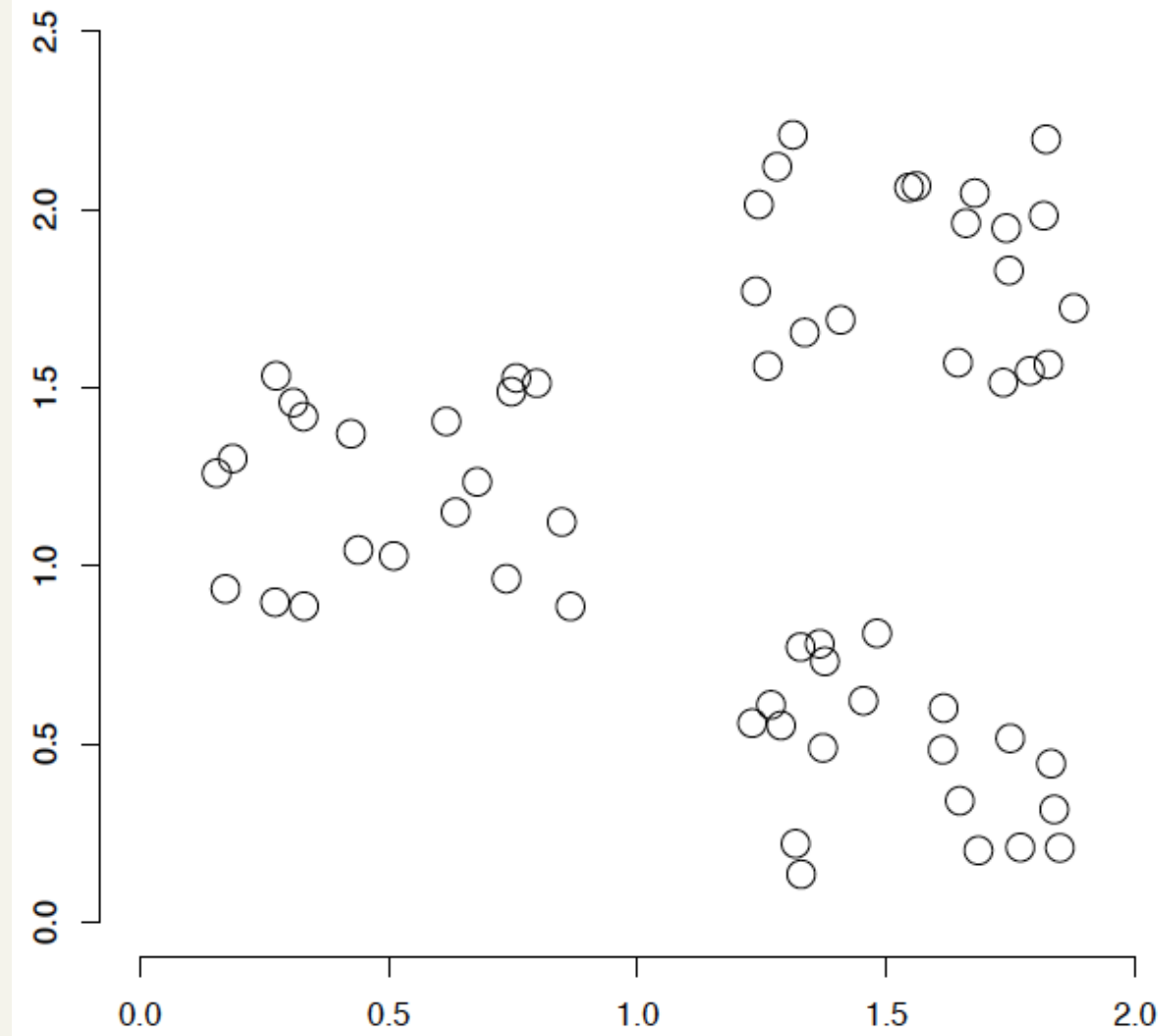
- Introduction
- Clustering algorithms
  - Partitional
  - Hierarchical

# What is clustering?

---

- **Clustering**: the process of grouping a set of objects into classes of similar objects
  - samples within a cluster should be similar.
  - samples from different clusters should be dissimilar.
- The commonest form of *unsupervised learning*
  - Unsupervised learning = learning from raw data, as opposed to supervised data where a classification of examples is given
  - A common and important task that finds many applications in IR and other places

# A data set with clear cluster structure



- How would you design an algorithm for finding the three clusters in this case?

# Clustering Algorithms

---

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

# Hard vs. soft clustering

---

- Hard clustering: Each samples belongs to exactly one cluster
  - More common and easier to do
- Soft clustering: A samples can belong to more than one cluster.
  - Makes more sense for applications like creating browsable hierarchies
  - You may want to put a pair of sneakers in two clusters: (i) sports apparel and (ii) shoes
  - You can only do that with a soft clustering approach.
- We won't do soft clustering today. See IIR 16.5, 18

# Partitioning Algorithms

---

- Partitioning method: Construct a partition of  $n$  samples into a set of  $K$  clusters
- Given: a set of samples and the number  $K$
- Find: a partition of  $K$  clusters that optimizes the chosen partitioning criterion
  - Globally optimal: exhaustively enumerate all partitions
  - Effective heuristic methods:  $K$ -means and  $K$ -medoids algorithms



# K-Means

---

- Assumes samples 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.

# K-Means Algorithm

---

Select  $K$  random centers  $\{s_1, s_2, \dots, s_K\}$  as seeds.  
Until clustering converges or other stopping criterion:

For each  $s_i$ :

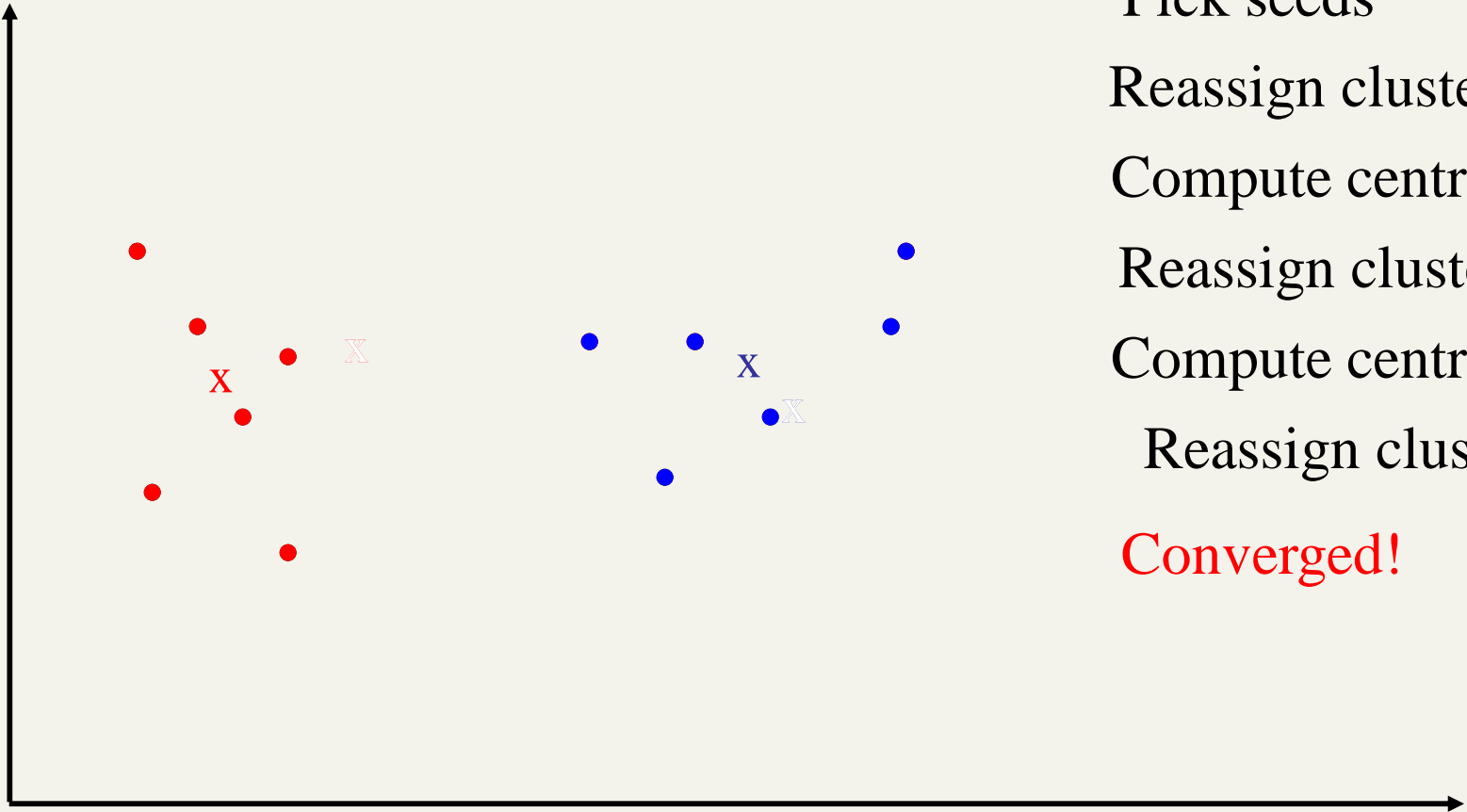
Assign  $d_i$  to the cluster  $c_j$  such that  $dist(x_i, s_j)$  is minimal.

*(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!**

# K Means Examples

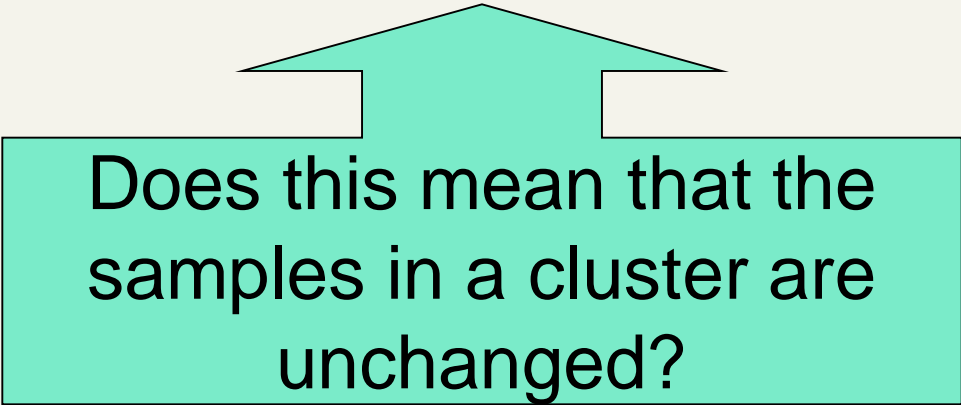
---

- Find the suitable clusters to classify the following data sets:
  - 1- { 1, 2, 1.5, 1.2, 2.5, 3.8}
  - 2- {1 , 2, 5, 7, 10, 15, 19}
- Find the R value that make the following data set separable into three clusters, then find the clusters:
  - 1- {10, 20, 7, 2, R, 15,100}

# Termination conditions

---

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



Does this mean that the samples in a cluster are unchanged?

# Convergence

---

- Why should the  $K$ -means algorithm ever reach a *fixed point*?
  - A state in which clusters don't change.
- $K$ -means is a special case of a general procedure known as the *Expectation Maximization (EM) algorithm*.
  - EM is known to converge.
  - Number of iterations could be large.
    - But in practice usually isn't

# Convergence of $K$ -Means

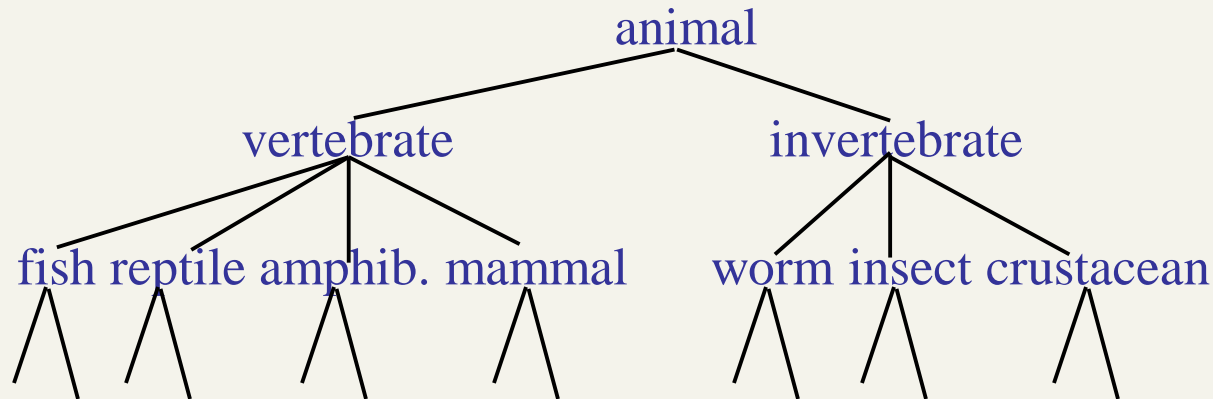
---

- Define goodness measure of cluster  $k$  as sum of squared distances from cluster centroid:
  - $G_k = \sum_i (s_i - c_k)^2$  (sum over all  $s_i$  in cluster  $k$ )
- $G = \sum_k G_k$
- Reassignment monotonically decreases  $G$  since each vector is assigned to the closest centroid.

# Hierarchical Clustering

---

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

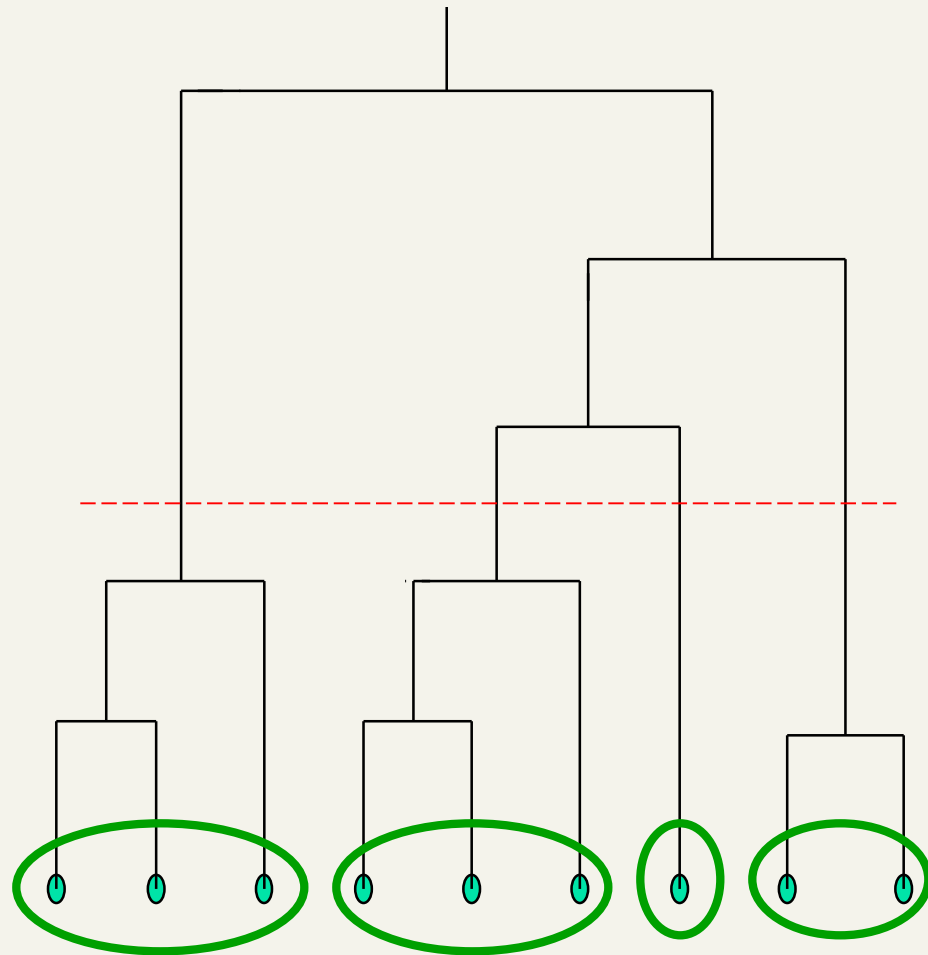


- One approach: recursive application of a partitional clustering algorithm.



# Dendrogram: Hierarchical Clustering

- 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**
- **Complete-link**
- **Centroid**
- **Average-link**

# Single Link Agglomerative Clustering

---

- Use maximum similarity of pairs:

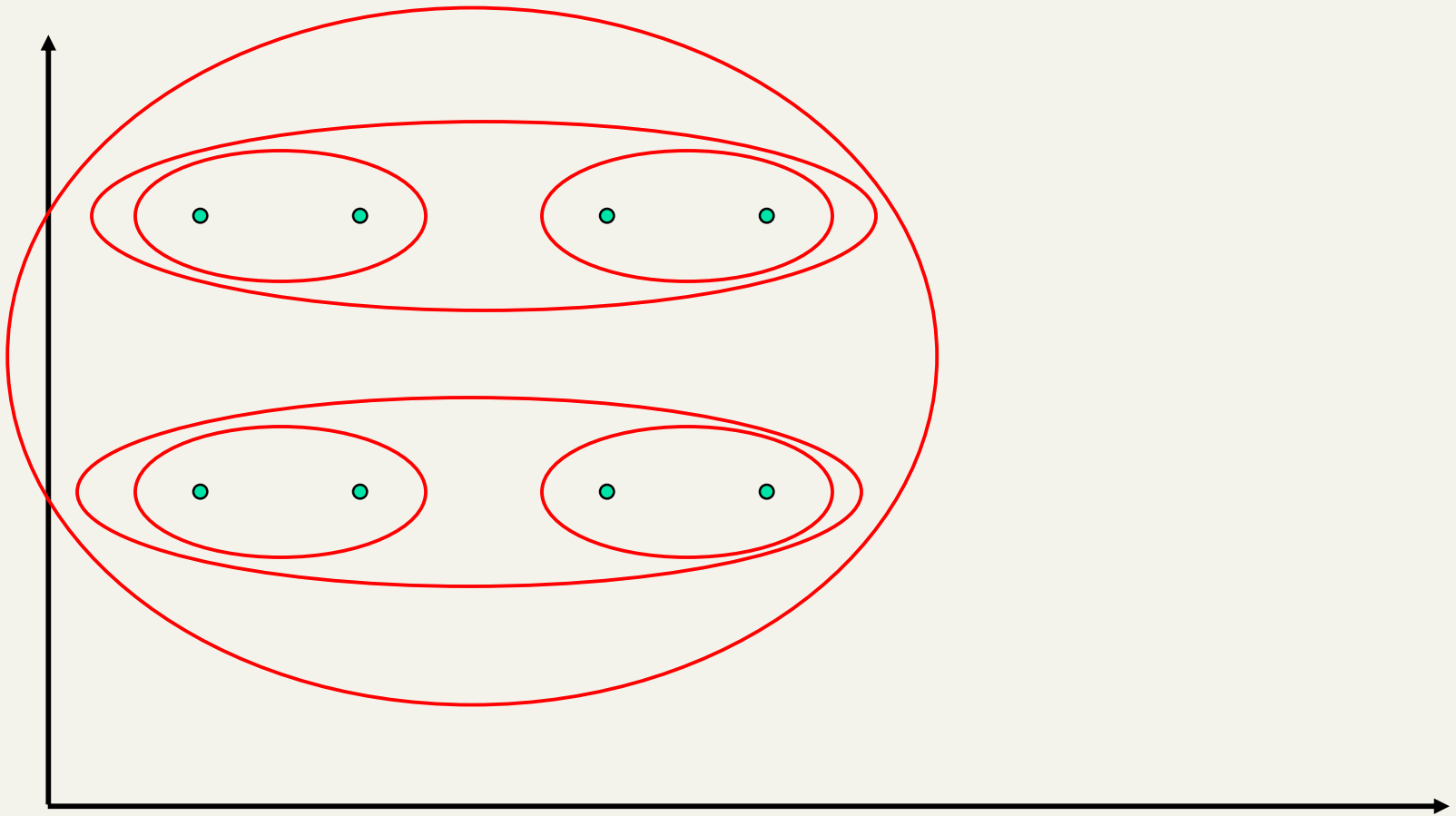
$$\mathit{sim}(c_i, c_j) = \max_{x \in c_i, y \in c_j} \mathit{sim}(x, y)$$

- Can result in “straggly” (long and thin) clusters due to chaining effect.
- After merging  $c_i$  and  $c_j$ , the similarity of the resulting cluster to another cluster,  $c_k$ , is:

$$\mathit{sim}((c_i \cup c_j), c_k) = \max(\mathit{sim}(c_i, c_k), \mathit{sim}(c_j, c_k))$$

# Single Link Example

---



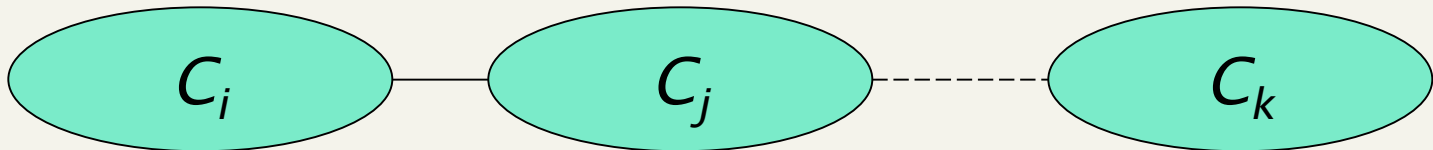
# Complete Link Agglomerative Clustering

- Use minimum similarity of pairs:

$$\text{sim}(c_i, c_j) = \min_{x \in c_i, y \in c_j} \text{sim}(x, y)$$

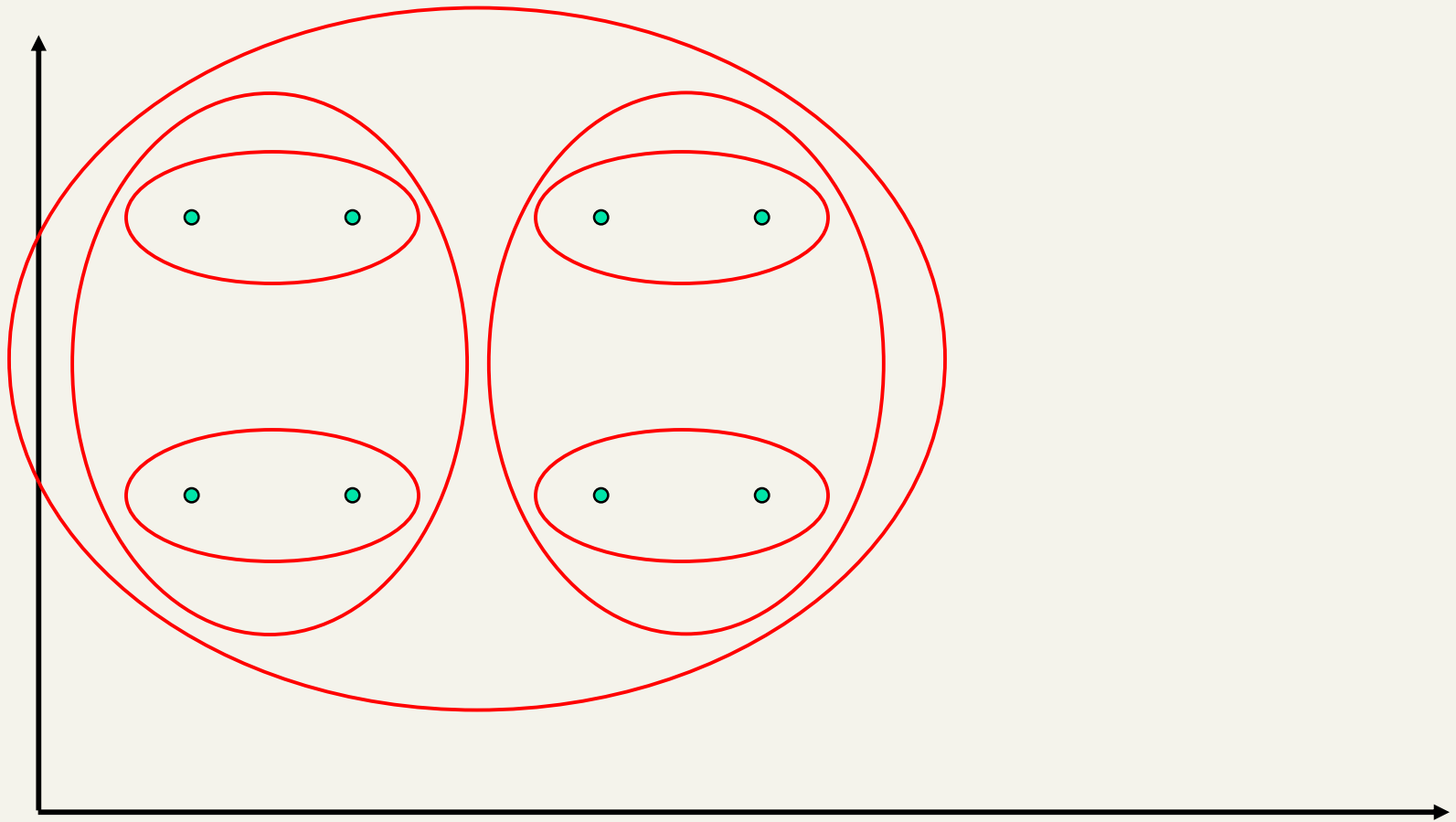
- Makes “tighter,” spherical clusters that are typically preferable.
- After merging  $c_i$  and  $c_j$ , the similarity of the resulting cluster to another cluster,  $c_k$ , is:

$$\text{sim}((c_i \cup c_j), c_k) = \min(\text{sim}(c_i, c_k), \text{sim}(c_j, c_k))$$



# Complete Link Example

---



# Group Average Agglomerative Clustering

- Similarity of two clusters = average similarity of all pairs within merged cluster.

$$sim(c_i, c_j) = \frac{1}{|c_i \cup c_j|(|c_i \cup c_j| - 1)} \sum_{\vec{x} \in (c_i \cup c_j)} \sum_{\vec{y} \in (c_i \cup c_j): \vec{y} \neq \vec{x}} sim(\vec{x}, \vec{y})$$

- Compromise between single and complete link.
- Two options:
  - Averaged across all ordered pairs in the merged cluster
  - Averaged over all pairs *between* the two original clusters
- No clear difference in efficacy



# Computing Group Average Similarity

---

- Always maintain sum of vectors in each cluster.

$$\vec{s}(c_j) = \sum_{\vec{x} \in c_j} \vec{x}$$

- Compute similarity of clusters in constant time:

$$\text{sim}(c_i, c_j) = \frac{(\vec{s}(c_i) + \vec{s}(c_j)) \bullet (\vec{s}(c_i) + \vec{s}(c_j)) - (|c_i| + |c_j|)}{(|c_i| + |c_j|)(|c_i| + |c_j| - 1)}$$

# What Is A Good Clustering?

---

- Internal criterion: A good clustering will produce high quality clusters in which:
  - the intra-class (that is, intra-cluster) similarity is high
  - the inter-class similarity is low
  - The measured quality of a clustering depends on both the samples representation and the similarity measure used

# External Evaluation of Cluster Quality

---

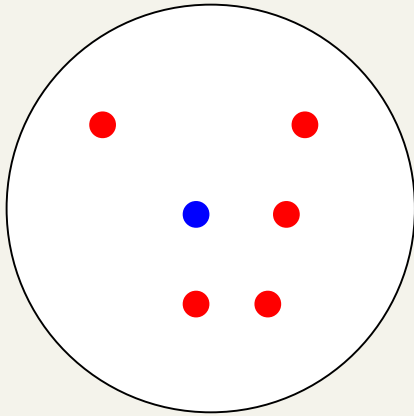
- Simple measure: purity, the ratio between the dominant class in the cluster  $\pi_i$  and the size of cluster  $\omega_i$

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

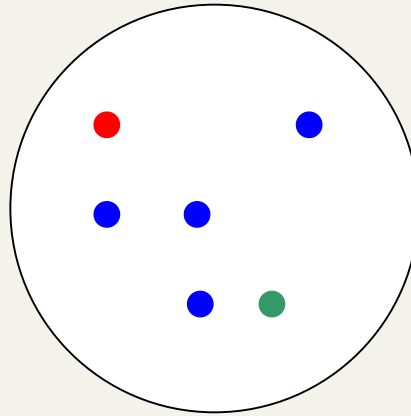
- Biased because having  $n$  clusters maximizes purity
- Others are entropy of classes in clusters (or mutual information between classes and clusters)

# Purity example

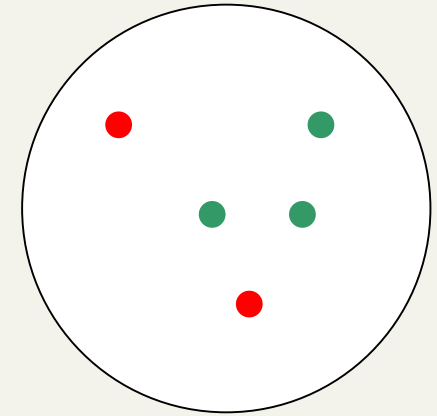
---



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 =  $1/5 (\max(2, 0, 3)) = 3/5$