

Clustering

COMP462/561: Computational Biology Methods

Fall 2016

M & W: 10:00 am – 11:30 am

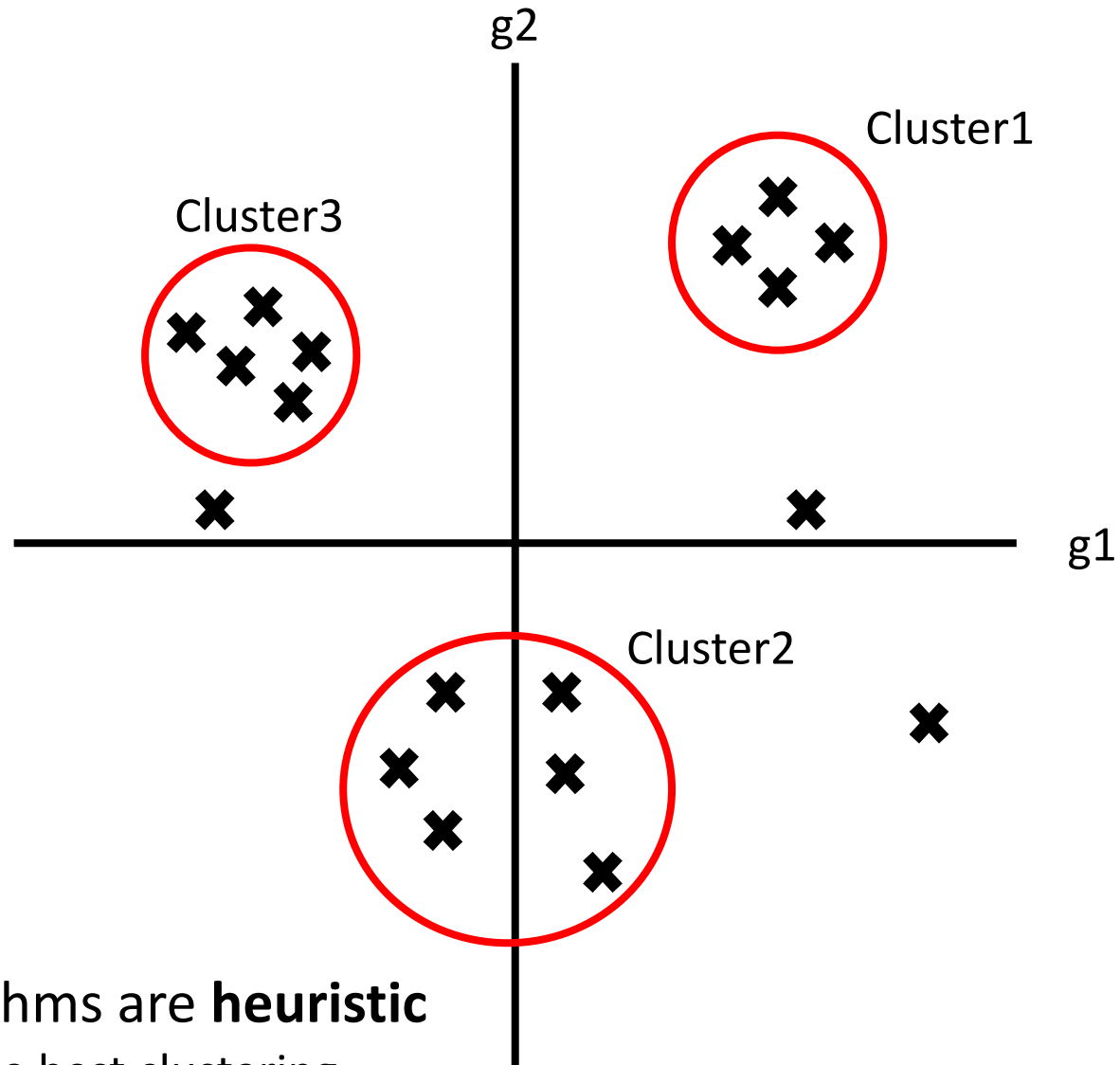
*Based on Course Notes by Dr. Mathieu Blanchette

Motivation

Given: A collection of unlabeled samples $X_1 \dots X_n$, where X_i represents the data for sample i

Goal: Partition samples into groups that are similar within themselves but dissimilar between

	X_1	...	X_n
gene1			
gene2			
gene3			
...			
gene _{$k-1$}			
gene _{k}			



- All the clustering algorithms are **heuristic**
 - They don't guarantee the best clustering

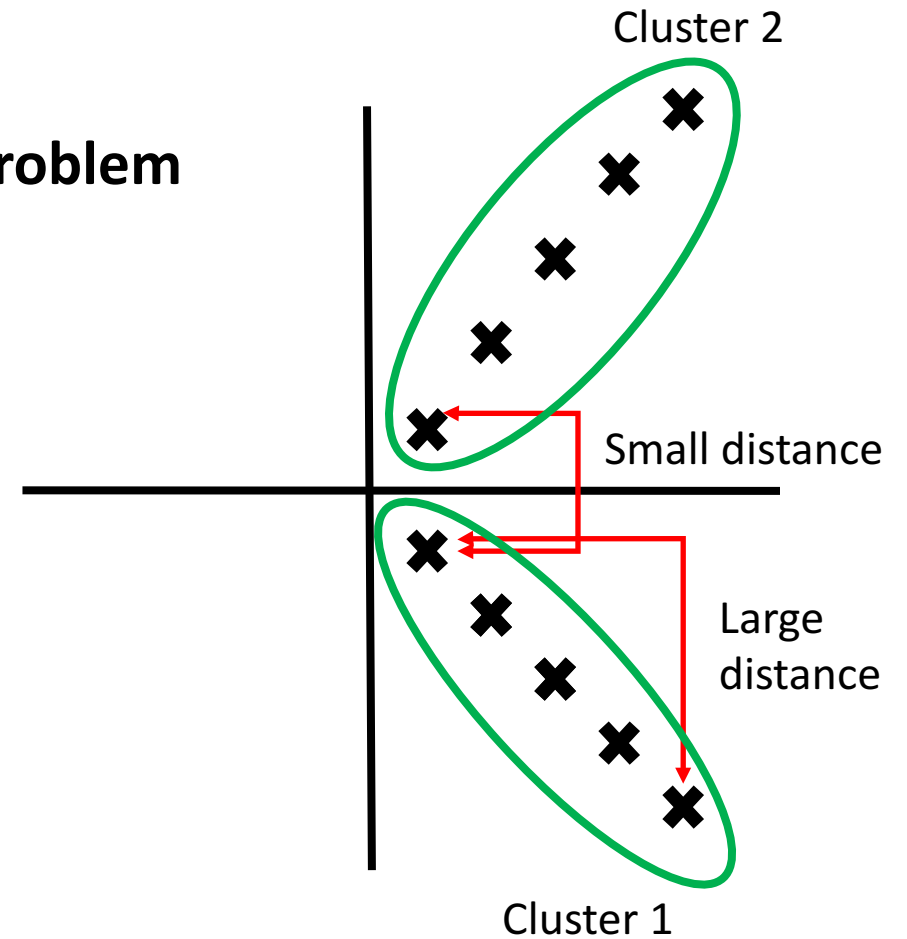
Similarity (or Distance) Measures

Given: Two expression profiles, X_i and X_j

Euclidean Distance

$$d_E(X_i, X_j) = \sqrt{\sum_{g=1 \dots k} (X_{i,g} - X_{j,g})^2}$$

Problem

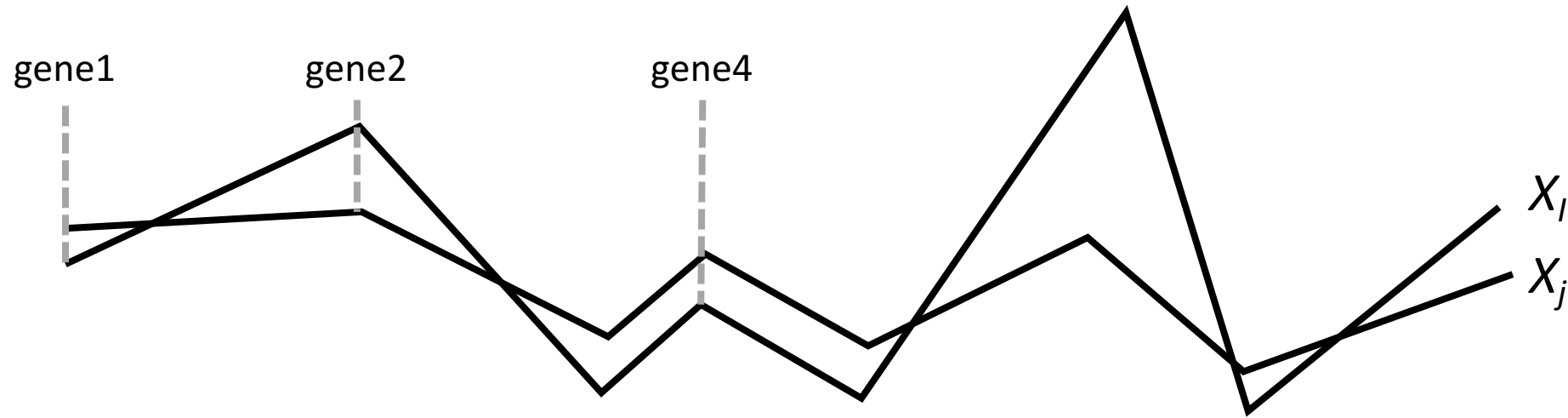


Pearson Correlation Coefficient

Similarity Measure

$$\begin{aligned} \text{Sim}(X_i, X_j) &= \frac{\text{Cov}(X_i, X_j)}{\sqrt{\text{Var}(X_i) \times \text{Var}(X_j)}} \\ &= \frac{\sum (X_i(g) - \bar{X}_i)(X_j(g) - \bar{X}_j)}{\sqrt{(\sum (X_i(g) - \bar{X}_i)^2) \times (\sum (X_j(g) - \bar{X}_j)^2)}} \end{aligned}$$

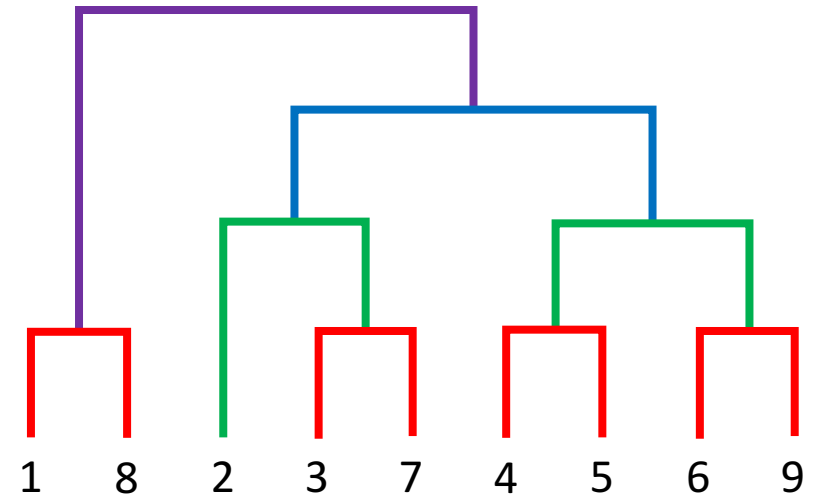
Pearson Correlation Coefficient Cont'd



- Different expression level
 - But always goes in the same direction

Hierarchical Clustering

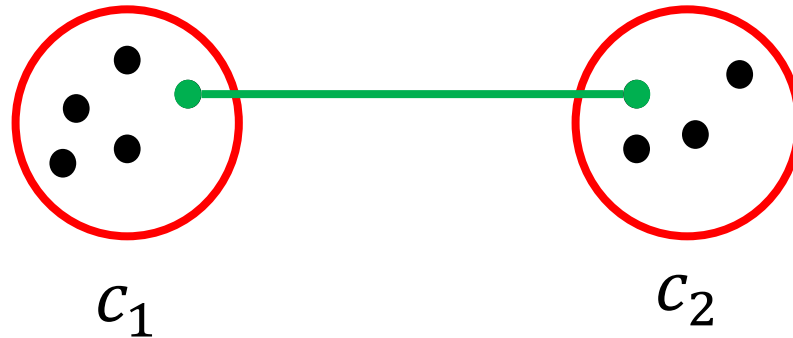
1. Start with each data point in its own cluster
2. Find the two clusters that are the closest and merge them
3. Repeat step two until all data points belong to a single cluster



Measuring Similarity Between Clusters

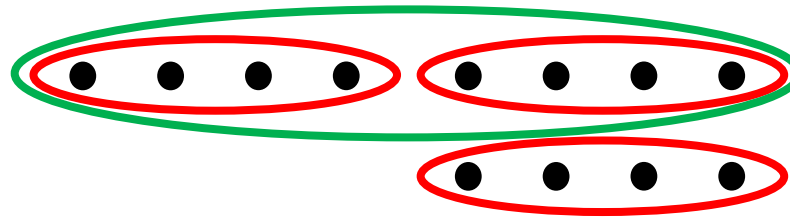
1) Single Linkage approach

$$Sim(c_1, c_2) = \max_{x \in c_1, y \in c_2} \{sim(x, y)\}$$



Problem

- Given the following data points:

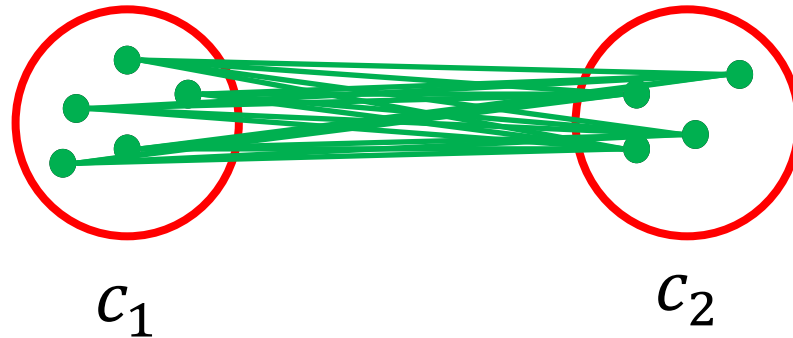


- Apply single linkage approach to clustering
- Get long and skinny clusters by having one point near the others
 - Shouldn't the two clusters on the right pair better together?

Measuring Similarity Between Clusters

2) Average linkage

$$Sim(c_1, c_2) = \frac{1}{|c_1| \cdot |c_2|} \sum_{x \in c_1, y \in c_2} Sim(x, y)$$



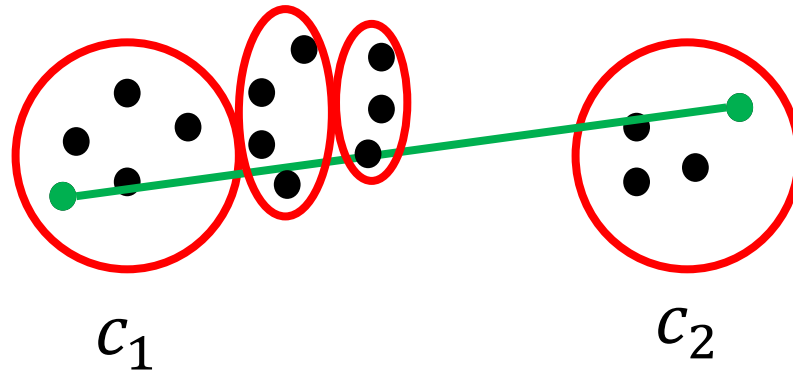
Take all pairs!

Measuring Similarity Between Clusters

3) Complete linkage

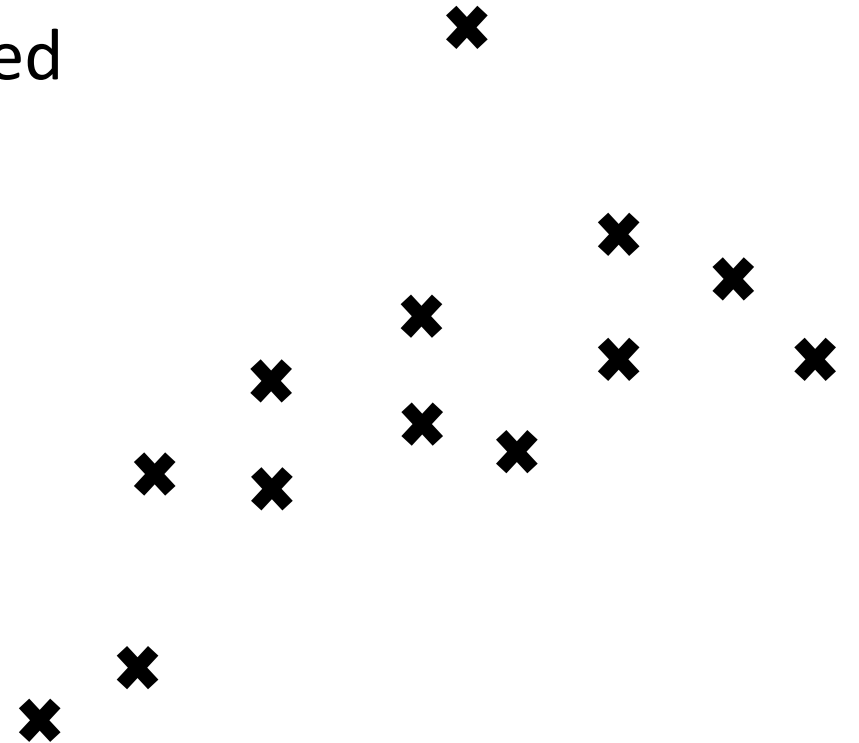
Makes very compact clusters

$$Sim(c_1, c_2) = \min_{x \in c_1, y \in c_2} sim(x, y)$$



K-Means Algorithm

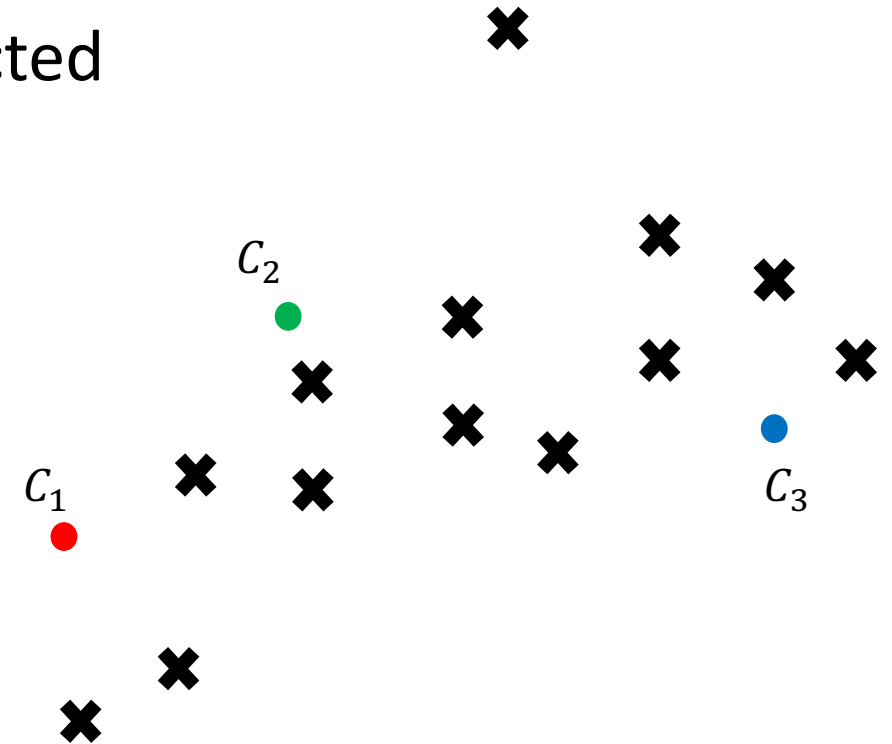
- 'k' is the number of clusters desired / expected
 - Each cluster has a centroid
1. Randomly choose k centroids
 2. Assign data points to nearest centroid
 3. Move centroid to center of cluster
 4. Repeat 2-4. Stop when no change to data point assignment



K-Means Algorithm

- 'k' is the number of clusters desired / expected
- Each cluster has a centroid

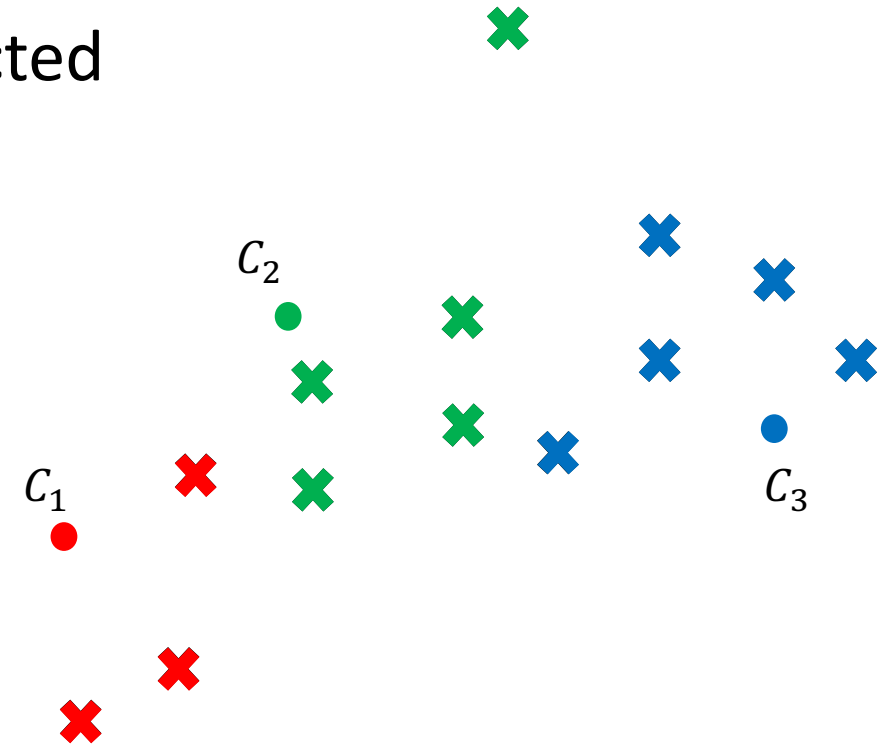
1. Randomly choose k centroids
2. Assign data points to nearest centroid
3. Move centroid to center of cluster
4. Repeat 2-4. Stop when no change to data point assignment



K-Means Algorithm

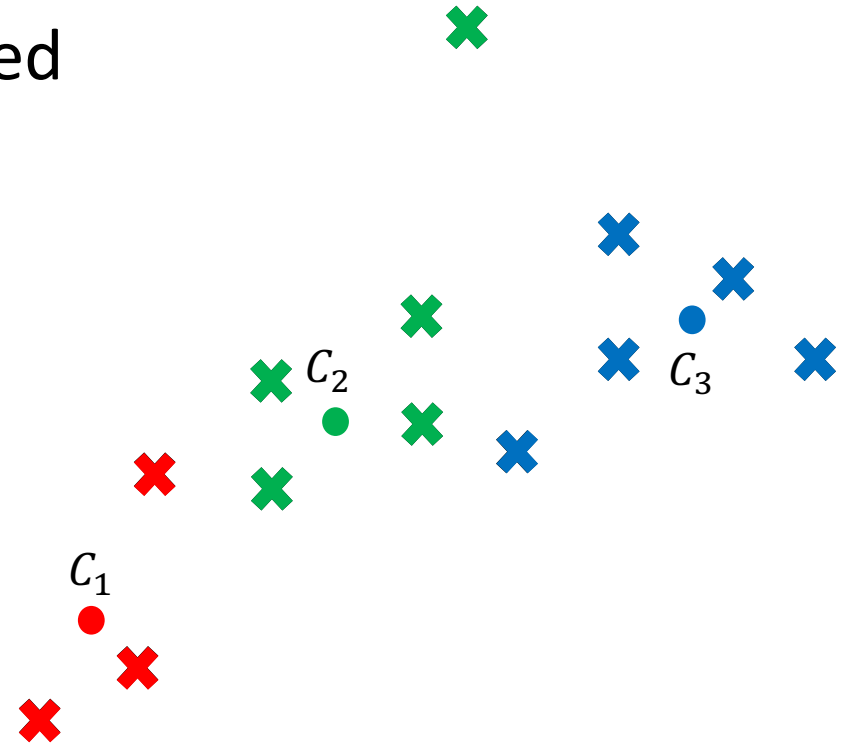
- 'k' is the number of clusters desired / expected
- Each cluster has a centroid

1. Randomly choose k centroids
2. Assign data points to nearest centroid
3. Move centroid to center of cluster
4. Repeat 2-4. Stop when no change to data point assignment



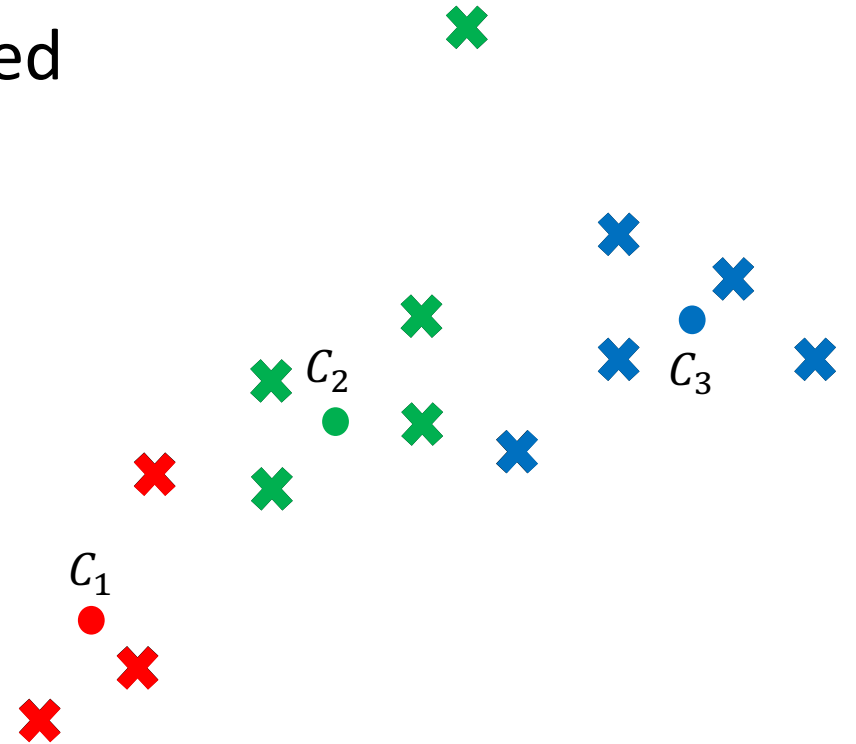
K-Means Algorithm

- 'k' is the number of clusters desired / expected
 - Each cluster has a centroid
1. Randomly choose k centroids
 2. Assign data points to nearest centroid
 3. Move centroid to center of cluster
 4. Repeat 2-4. Stop when no change to data point assignment



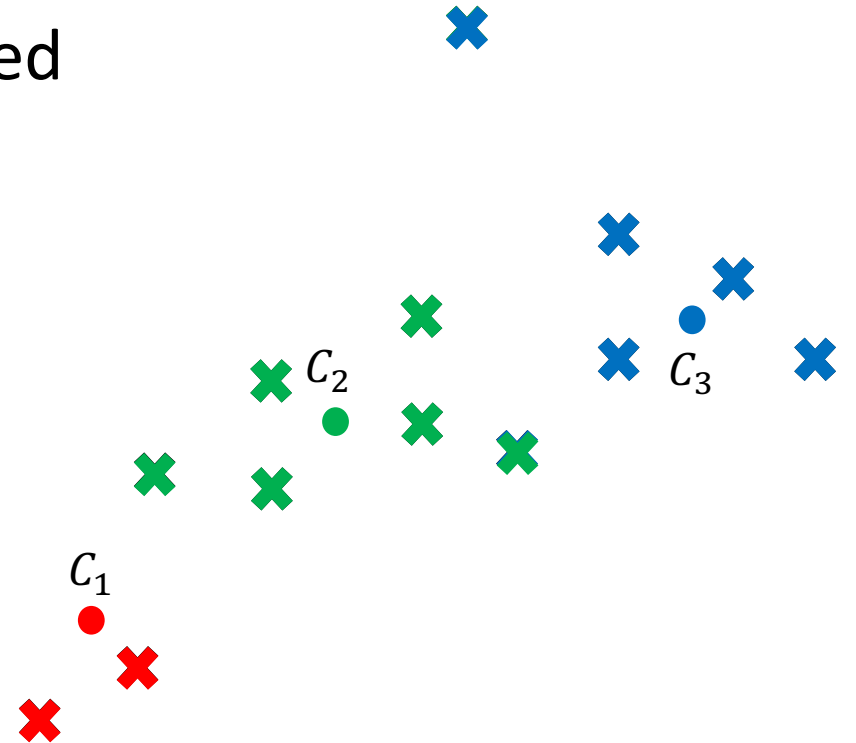
K-Means Algorithm

- 'k' is the number of clusters desired / expected
 - Each cluster has a centroid
1. Randomly choose k centroids
 2. Assign data points to nearest centroid
 3. Move centroid to center of cluster
 4. Repeat 2-4. Stop when no change to data point assignment



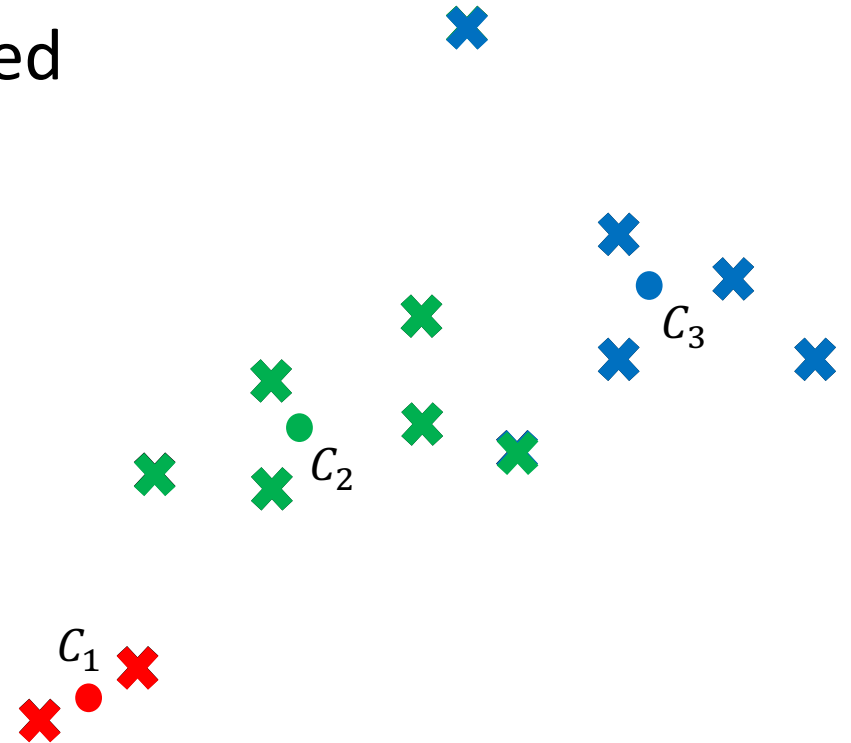
K-Means Algorithm

- 'k' is the number of clusters desired / expected
 - Each cluster has a centroid
1. Randomly choose k centroids
 2. Assign data points to nearest centroid
 3. Move centroid to center of cluster
 4. Repeat 2-4. Stop when no change to data point assignment



K-Means Algorithm

- 'k' is the number of clusters desired / expected
 - Each cluster has a centroid
1. Randomly choose k centroids
 2. Assign data points to nearest centroid
 3. Move centroid to center of cluster
 4. Repeat 2-4. Stop when no change to data point assignment



Cluster Validation

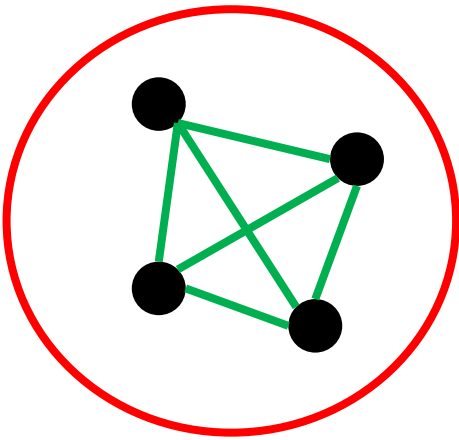
- **Cohesion:** measures how closely related data points in a cluster are (i.e., within cluster Sum of Squares [WSS])

$$WSS = \sum_i \sum_{x \in c_i} \|x - m_i\|^2$$

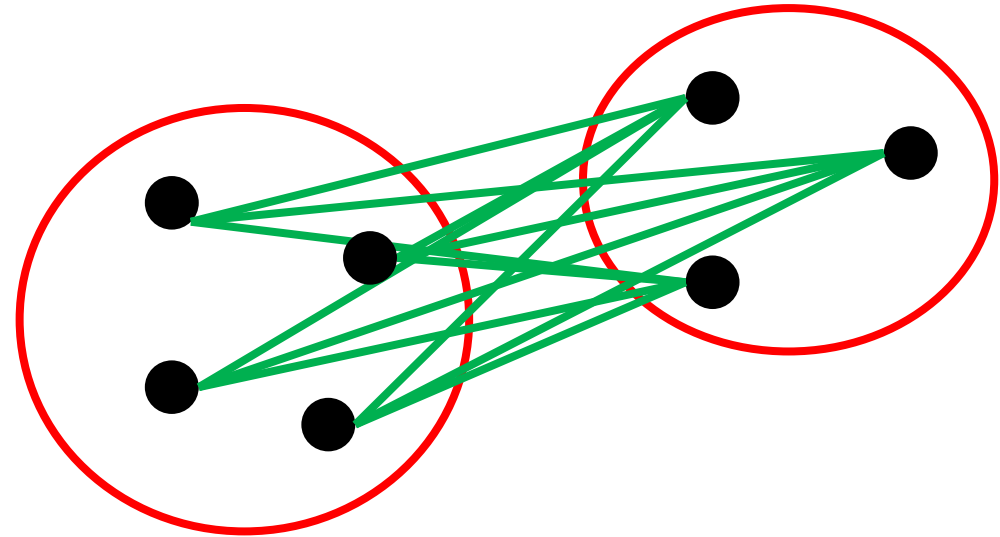
- **Separation:** measures how distinct or well-separated a cluster is from others (i.e., between cluster Sum of Squares [BSS])

$$BSS = \sum_i \sum_j |c_i| \cdot |c_j| \cdot \|m_i - m_j\|^2$$

Cohesion and Separation



Cohesion



Separation