

MACHINE LEARNING WITH PYTHON

# K-MEANS CLUSTERING

---

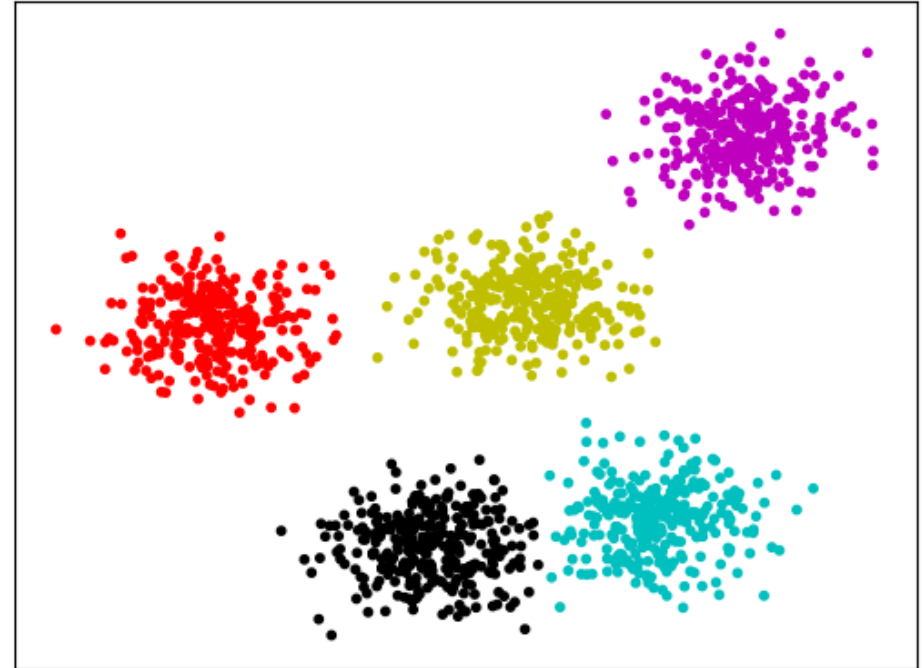
Themistoklis Diamantopoulos

# Clustering

- Split data into clusters according to features
- K-Means
  - Given the number of clusters  $k$
  - Split the data into clusters
  - Each cluster has a centroid
  - Minimize sum-of-squared-errors:

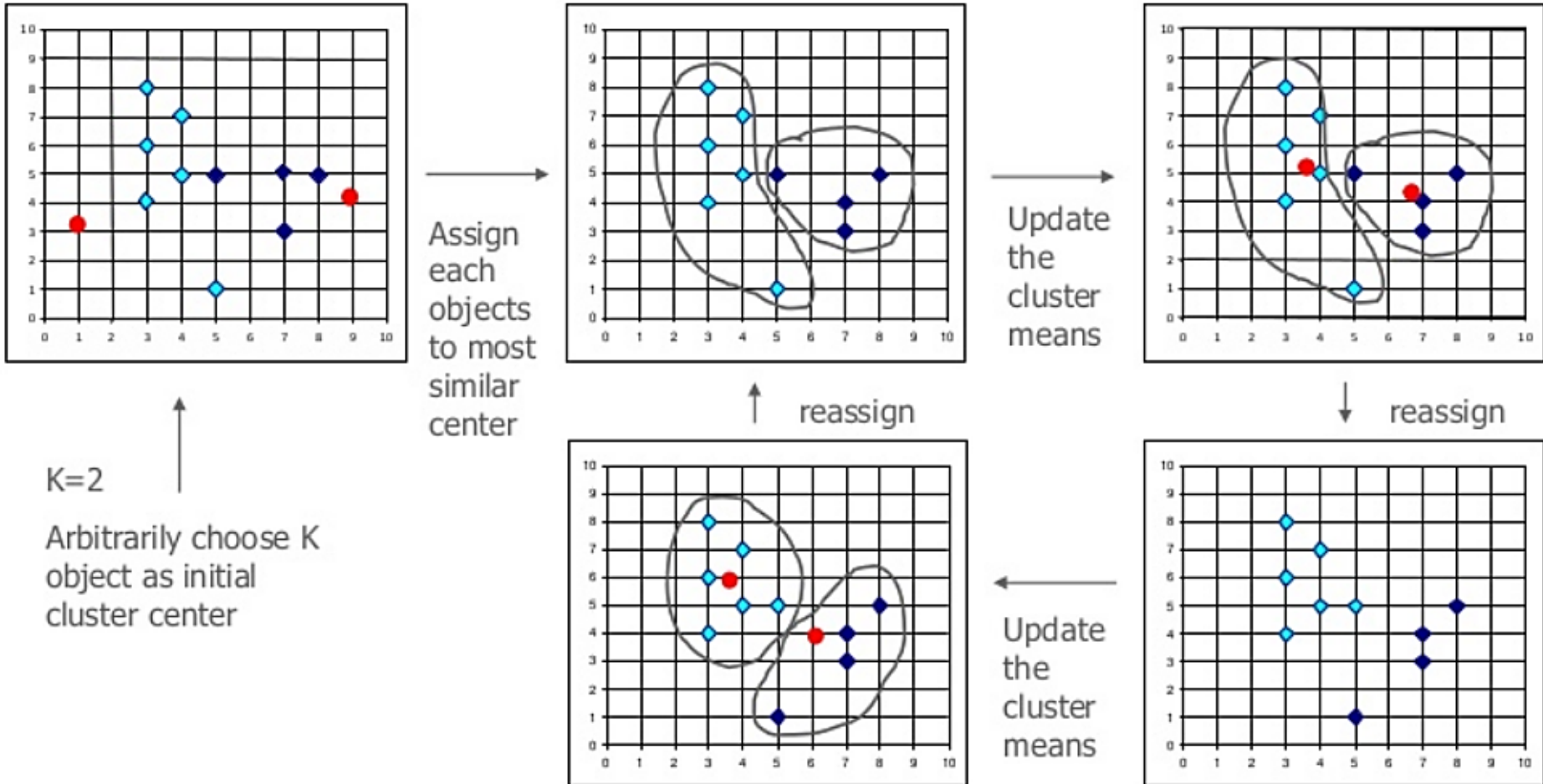
$$J = \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2$$

where  $x$  are datapoints and  $i=1 \dots k$  refer to clusters  $S_i$  with centroids  $\mu_i$



Euclidean distance:  $\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$

# K-Means Clustering



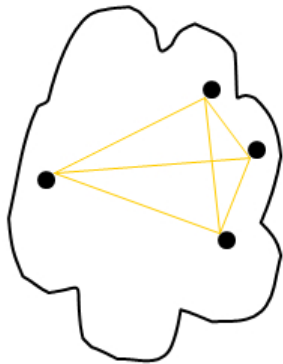
# Evaluation

- Cohesion (Within Cluster Sum of Squares)

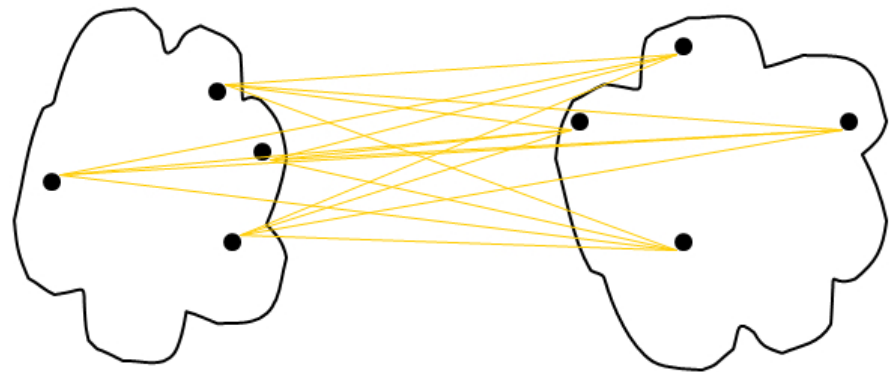
$$SSE = WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

- Separation (Between Cluster Sum of Squares)

$$BSS = \sum_i |C_i| (m - m_i)^2$$



cohesion



separation

# Silhouette

- Silhouette value of  $x$ :

$$s = (b - a) / \max(a, b)$$

where

- $a$  = average distance of  $x$  to the points in its cluster
- $b$  =  $\min(\text{average distance of } x \text{ to points in another cluster})$
- Value:
  - Positive/close to 1: possibly assigned to proper cluster
  - Negative: possibly assigned to wrong cluster
  - Close to 0: on border of 2 clusters

