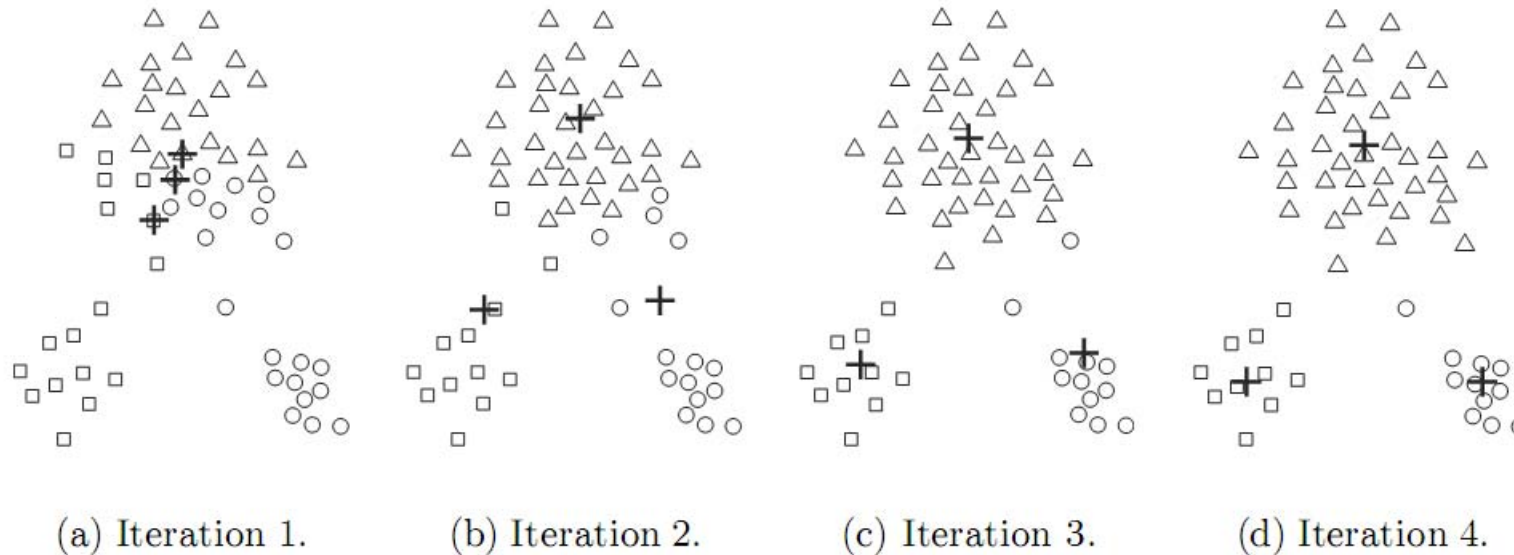


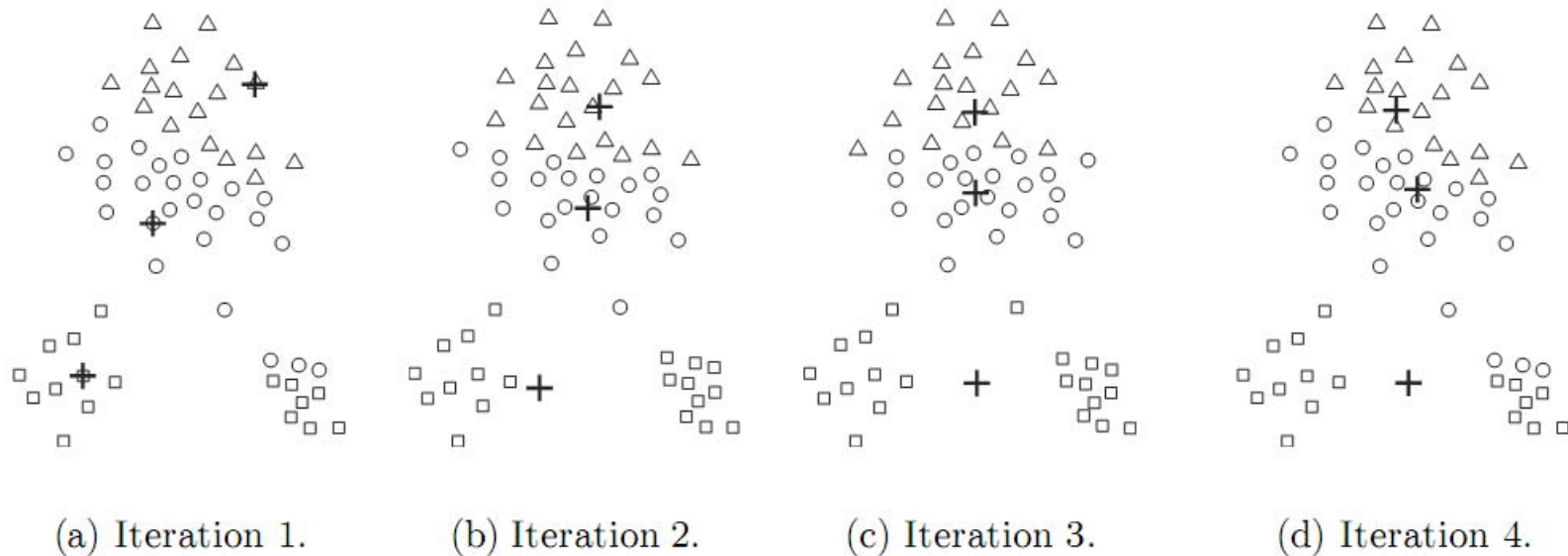
K-Means clustering



- 1: Select K points as initial centroids.
- 2: **repeat**
- 3: Form K clusters by assigning each point to its closest centroid.
- 4: Recompute the centroid of each cluster.
- 5: **until** Centroids do not change.

*I can use different
types of distances !!!*

How many centers? How to initialize them?



- Often the cluster centers are initialized to some of the points, picked at random
- Which “optimality criterion” is the k-means algorithm optimizing?

The algorithm attempts to minimize a sum of square criterion:

$$\sum_{j=1}^K \sum_{n \in S_j} d(x_n, \mu_j)^2$$

K number of clusters

x_n is a vector representing the n -th data point

$n \in S_j$ indicates the set of points belonging to cluster S_j

μ_j is the mean of the data points in S_j

$d(x_n, \mu_j)$ indicates the a **distance** between x_n and μ_j

Example of distances: Euclidean distance,
corrected Pearson correlation coefficient