

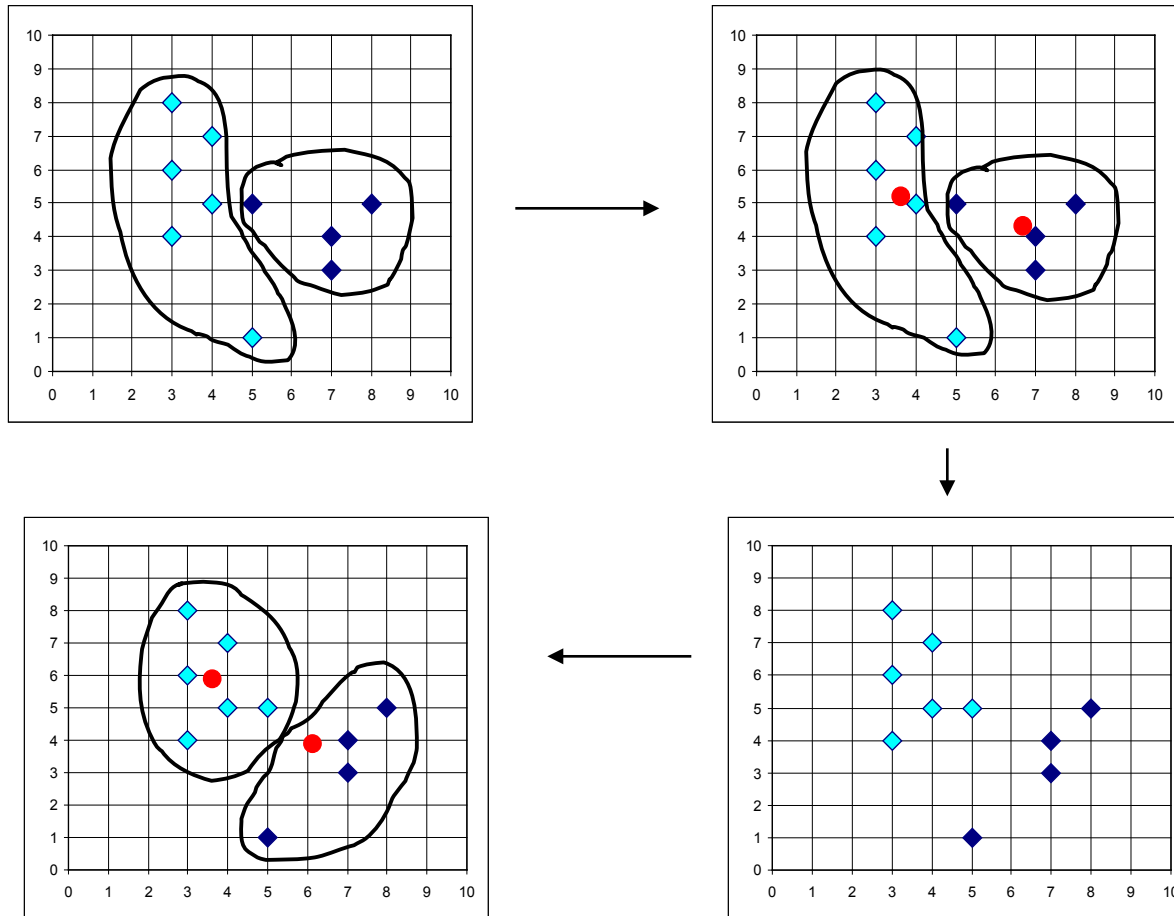
# Clustering

# *K-Means* Clustering

- Given  $k$ , the *k-means* algorithm consists of four steps:
  - Select initial centroids at random.
  - Assign each object to the cluster with the nearest centroid.
  - Compute each centroid as the mean of the objects assigned to it.
  - Repeat previous 2 steps until no change.

# *K-Means* Clustering (contd.)

- Example



# Comments on the *K-Means* Method

- Strengths

- *Relatively efficient*:  $O(tkn)$ , where  $n$  is # objects,  $k$  is # clusters, and  $t$  is # iterations. Normally,  $k, t \ll n$ .
- Often terminates at a *local optimum*. The *global optimum* may be found using techniques such as *simulated annealing* and *genetic algorithms*

- Weaknesses

- Applicable only when *mean* is defined (what about categorical data?)
- Need to specify  $k$ , the *number* of clusters, in advance
- Trouble with noisy data and *outliers*
- Not suitable to discover clusters with *non-convex shapes*