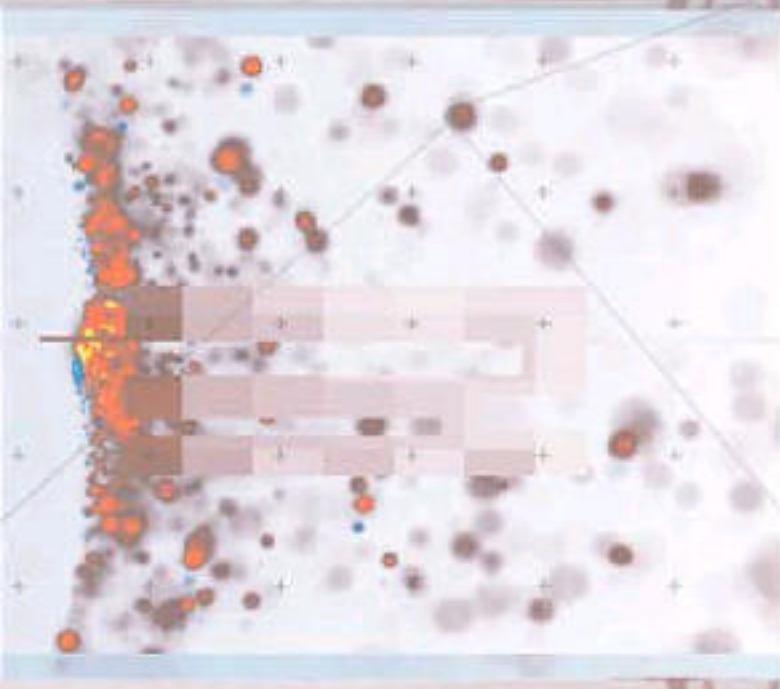




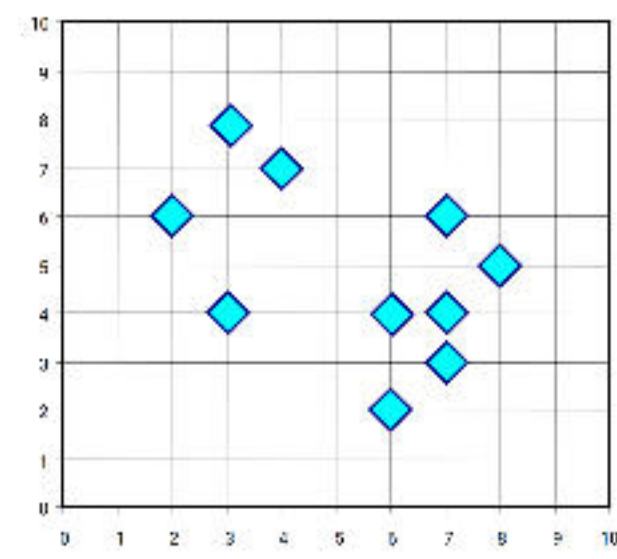
The *K-Medoids* Clustering Method



Handling Outliers: From *K-Means* to *K-Medoids*

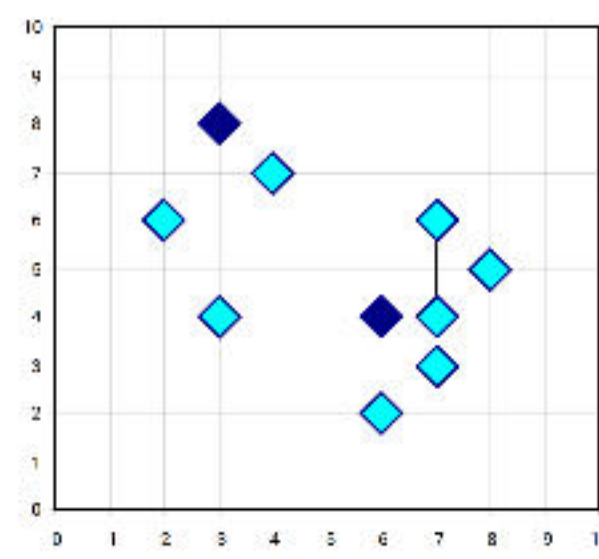
- ❑ The *K-Means* algorithm is sensitive to outliers!—since an object with an extremely large value may substantially distort the distribution of the data
- ❑ *K-Medoids*: Instead of taking the **mean** value of the object in a cluster as a reference point, **medoids** can be used, which is the **most centrally located** object in a cluster
- ❑ The *K-Medoids* clustering algorithm:
 - ❑ Select K points as the initial representative objects (i.e., as initial K medoids)
 - ❑ **Repeat**
 - ❑ Assigning each point to the cluster with the closest medoid
 - ❑ Randomly select a non-representative object o_i
 - ❑ Compute the total cost S of swapping the medoid m with o_i
 - ❑ If $S < 0$, then swap m with o_i to form the new set of medoids
 - ❑ **Until** convergence criterion is satisfied

PAM: A Typical *K-Medoids* Algorithm

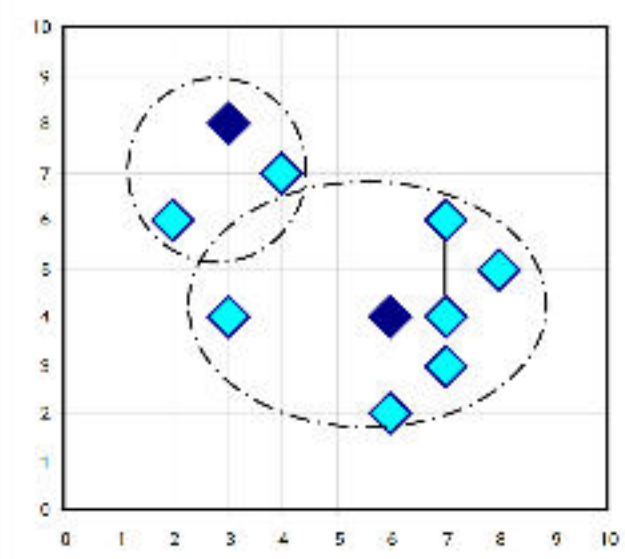


$K = 2$

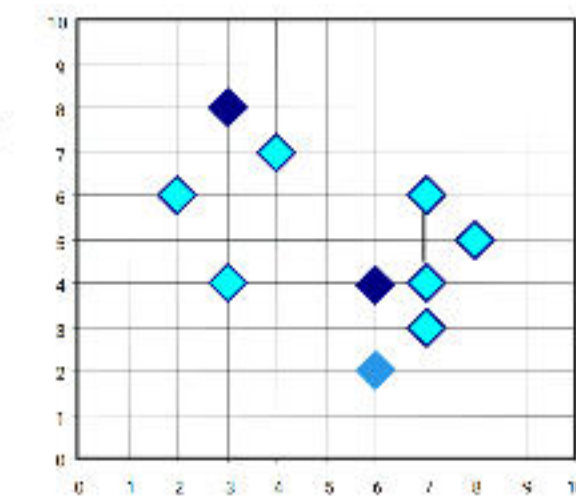
Arbitrary
choose K
object as
initial
medoids



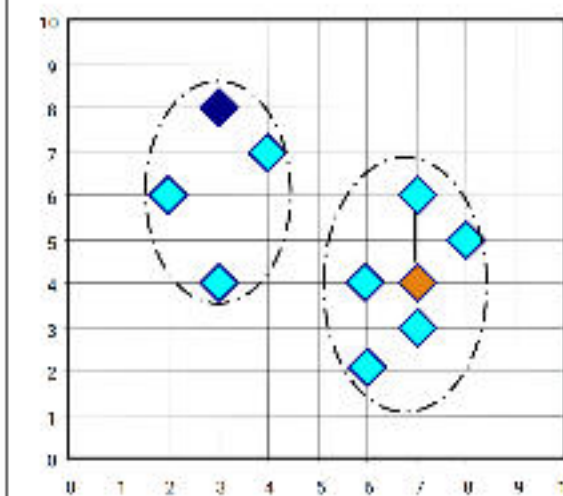
Assign
each
remaining
object to
nearest
medoids



Randomly select a non-
medoid object, O_{random}



Compute
total cost of
swapping



Swapping O
and O_{random}
If quality is
improved

Select initial K medoids randomly

Repeat

Object re-assignment

Swap medoid m with o_i if it
improves the clustering quality

Until convergence criterion is satisfied

Discussion on *K-Medoids* Clustering

- *K-Medoids* Clustering: Find *representative* objects (medoids) in clusters
- *PAM* (Partitioning Around Medoids: Kaufmann & Rousseeuw 1987)
 - Starts from an initial set of medoids, and
 - Iteratively replaces one of the medoids by one of the non-medoids if it improves the total sum of the squared errors (SSE) of the resulting clustering
 - *PAM* works effectively for small data sets but does not scale well for large data sets (due to the computational complexity)
 - Computational complexity: *PAM*: $O(K(n - K)^2)$ (quite expensive!)
- Efficiency improvements on *PAM*
 - *CLARA* (Kaufmann & Rousseeuw, 1990):
 - *PAM* on samples; $O(Ks^2 + K(n - K))$, s is the sample size
 - *CLARANS* (Ng & Han, 1994): Randomized re-sampling, ensuring efficiency + quality