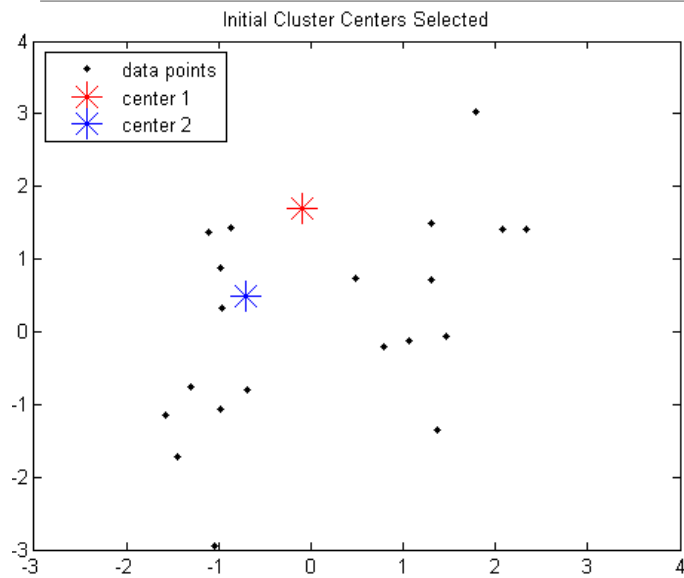# The *K-Means* Clustering Method
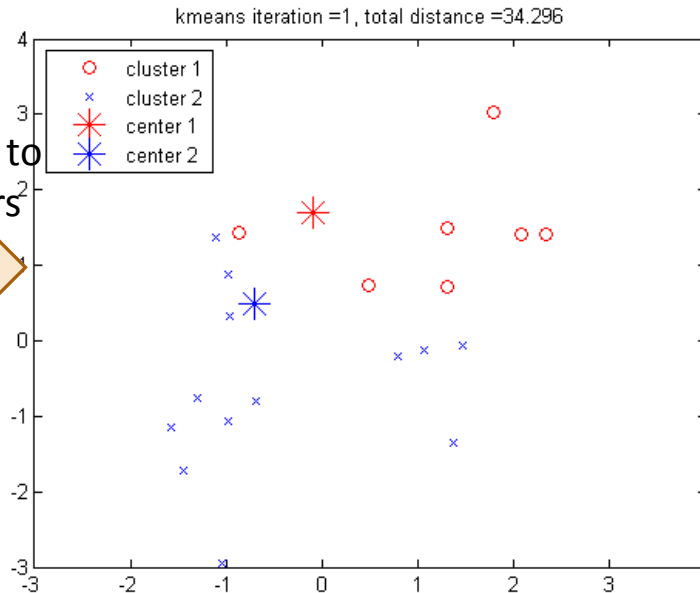
# The *K-Means* Clustering Method

- ❑ <u>*K-Means*</u> (MacQueen'67, Lloyd'57/'82)

  - ❑ Each cluster is represented by the center of the cluster

- ❑ Given K, the number of clusters, the *K-Means* clustering algorithm is outlined as follows

  - ❑ Select *K* points as initial centroids

  - ❑ **Repeat**

    - ❑ Form *K* clusters by assigning each point to its closest centroid

    - ❑ Re-compute the centroids (i.e., *mean point*) of each cluster

  - ❑ **Until** convergence criterion is satisfied

- ❑ Different kinds of measures can be used

  - ❑ Manhattan distance ($L_1$ norm), Euclidean distance ($L_2$ norm), Cosine similarity
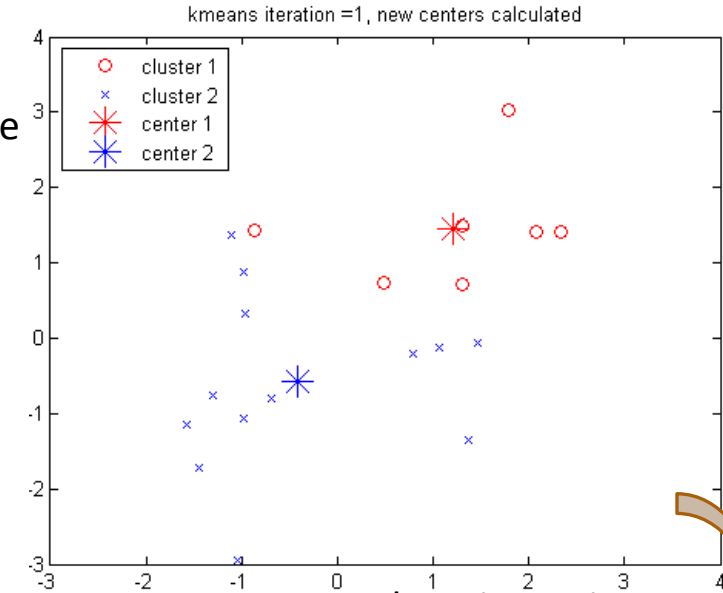
# Example: *K-Means* Clustering



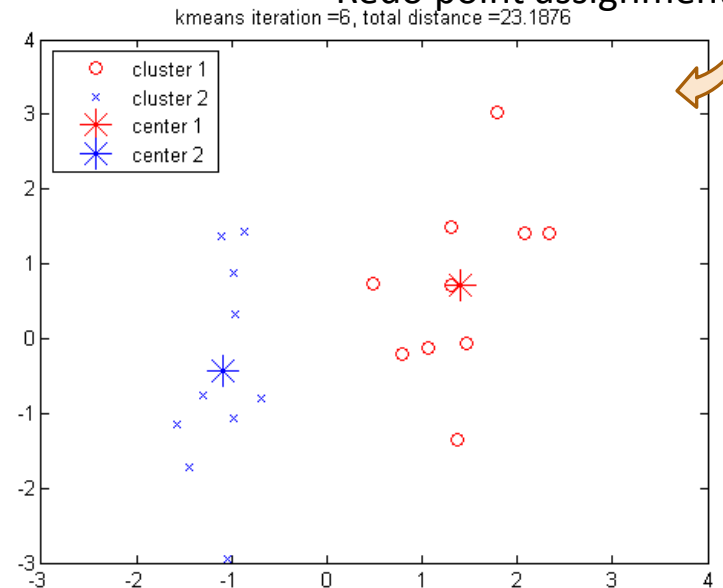The original data points & randomly select *K* = 2 centroids

Assign points to clusters

Recompute cluster centers

Redo point assignment

*Execution of the K-Means* Clustering Algorithm

Select *K* points as initial centroids

**Repeat**

- Form *K* clusters by assigning each point to its closest centroid
- Re-compute the centroids (i.e., *mean point*) of each cluster

**Until** convergence criterion is satisfied

3

# Discussion on the *K-Means* Method

- ❑ **Efficiency**: $O(tKn)$ where $n$: # of objects, $K$: # of clusters, and $t$: # of iterations
  - ❑ Normally, $K$, $t << n$; thus, an efficient method
- ❑ K-means clustering often ***terminates at*** **a** ***local optimal***
  - ❑ Initialization can be important to find high-quality clusters
- ❑ **Need to specify *K*,** the *number* of clusters, in advance
  - ❑ There are ways to automatically determine the "*best*" K
  - ❑ In practice, one often runs a range of values and selected the "*best*" K value
- ❑ **Sensitive to noisy data and *outliers***
  - ❑ Variations: Using K-medians, K-medoids, etc.
- ❑ K-means is applicable only to objects in a continuous n-dimensional space
  - ❑ Using the K-modes for ***categorical data***
- ❑ Not suitable to discover clusters with ***non-convex shapes***
  - ❑ Using density-based clustering, kernel *K*-means, etc.

# Variations of *K-Means*

❑ There are many variants of the *K-Means* method, varying in different aspects

    ❑ Choosing better initial centroid estimates

       ❑ *K-means++, Intelligent K-Means, Genetic K-Means*     To be discussed in this lecture

    ❑ Choosing different representative prototypes for the clusters

       ❑ *K-Medoids, K-Medians, K-Modes*     To be discussed in this lecture

    ❑ Applying feature transformation techniques

       ❑ *Weighted K-Means, Kernel K-Means*     To be discussed in this lecture