

The background features a complex network of thin, light-colored lines forming a mesh or Voronoi diagram. Scattered throughout are small, colored dots in shades of green, blue, and orange. A prominent, thicker red line forms a large, irregular loop in the center. On the left side, there is a vertical strip with a grid of small '+' symbols and a series of horizontal bars of varying lengths and colors (orange, red, brown).

# **Basic Concepts of Partitioning Algorithms**

# Partitioning Algorithms: Basic Concepts

---

- ❑ Partitioning method: Discovering the groupings in the data by optimizing a specific objective function and iteratively improving the quality of partitions
- ❑ *K*-partitioning method: Partitioning a dataset ***D*** of ***n*** objects into a set of ***K*** clusters so that an objective function is optimized (e.g., the sum of squared distances is minimized, where  $c_k$  is the centroid or medoid of cluster  $C_k$ )

❑ A typical objective function: **Sum of Squared Errors (SSE)**

$$SSE(C) = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - c_k\|^2$$

- ❑ Problem definition: Given *K*, find a partition of *K clusters* that optimizes the chosen partitioning criterion
  - ❑ Global optimal: Needs to exhaustively enumerate all partitions
  - ❑ Heuristic methods (i.e., greedy algorithms): *K-Means*, *K-Medians*, *K-Medoids*, etc.



The background features a complex collage of geometric and data-related elements. It includes a network of red lines connecting green dots, a grid of small purple arrows, and various abstract shapes in shades of brown, orange, and blue. A large white chevron shape points to the right, framing the title.

# The *K-Means* Clustering Method

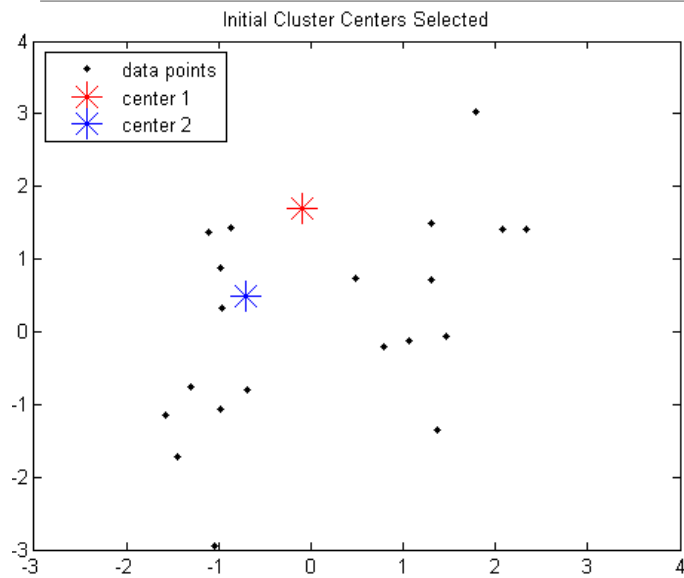


# The *K-Means* Clustering Method

---

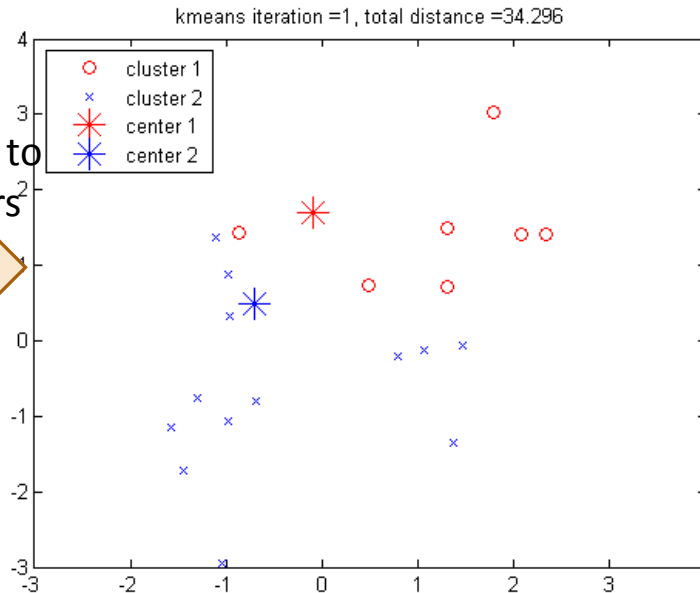
- ❑ *K-Means* (MacQueen'67, Lloyd'57/'82)
  - ❑ Each cluster is represented by the center of the cluster
- ❑ Given  $K$ , the number of clusters, the *K-Means* clustering algorithm is outlined as follows
  - ❑ Select  $K$  points as initial centroids
  - ❑ **Repeat**
    - ❑ Form  $K$  clusters by assigning 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

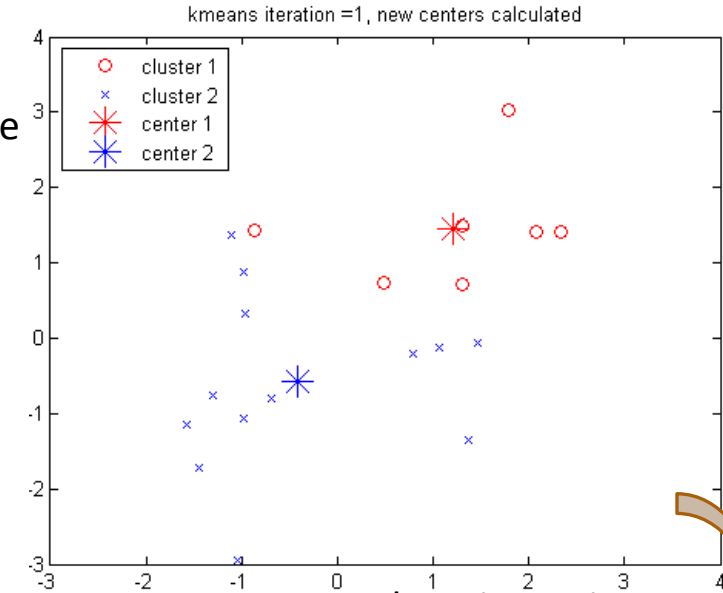
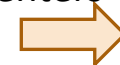


The original data points & randomly select  $K = 2$  centroids

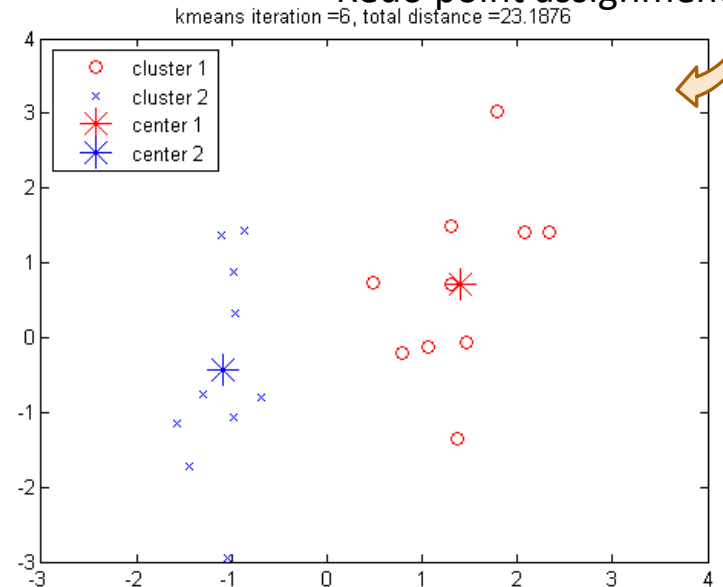
Assign  
points to  
clusters



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

# Discussion on the *K-Means* Method

---

- ❑ **Efficiency:**  $O(tKn)$  where  $n$ : # of objects,  $K$ : # of clusters, and  $t$ : # of iterations
  - ❑ Normally,  $K, t \ll n$ ; 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



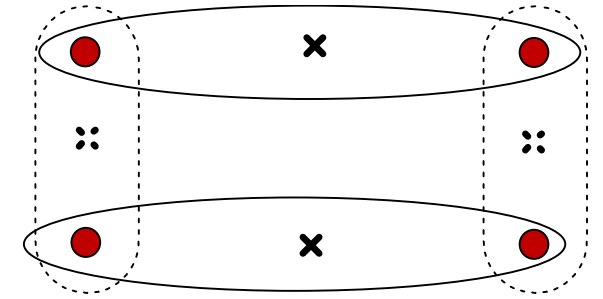


# Initialization of K-Means Clustering

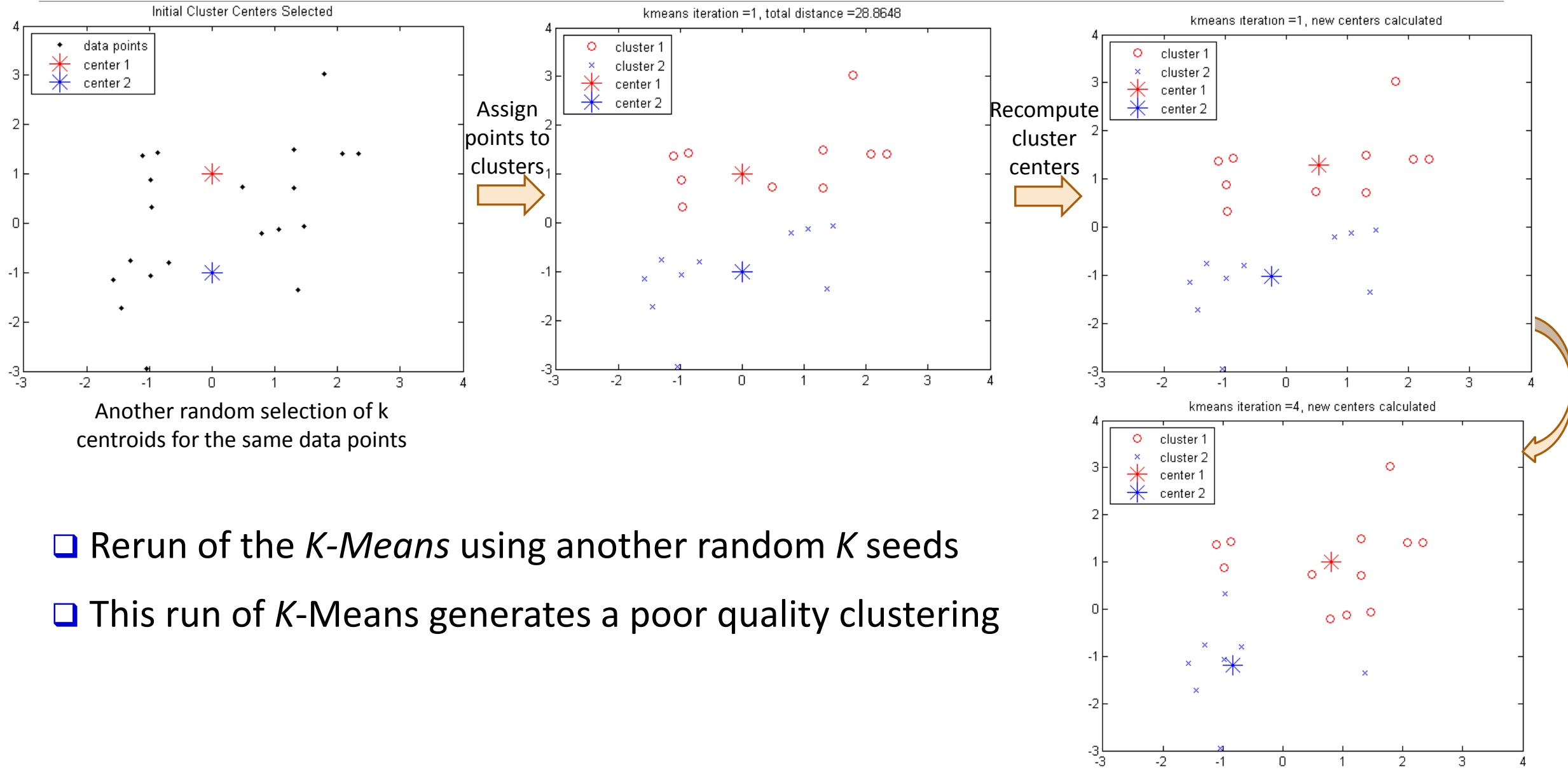


# Initialization of K-Means

- ❑ Different initializations may generate rather different clustering results (some could be far from optimal)
- ❑ Original proposal (MacQueen'67): Select  $K$  seeds randomly
  - ❑ Need to run the algorithm multiple times using different seeds
- ❑ There are many methods proposed for better initialization of  $k$  seeds
  - ❑ ***K-Means++*** (Arthur & Vassilvitskii'07):
    - ❑ The first centroid is selected at random
    - ❑ The next centroid selected is the one that is farthest from the currently selected (selection is based on a weighted probability score)
    - ❑ The selection continues until  $K$  centroids are obtained

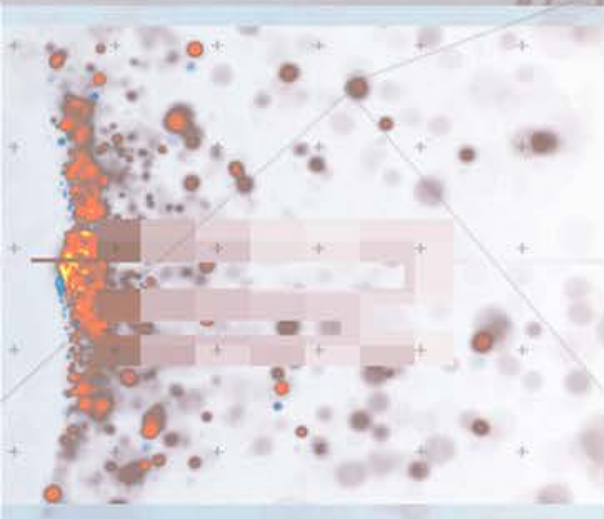


# Example: Poor Initialization May Lead to Poor Clustering





# The *K-Medoids* Clustering Method



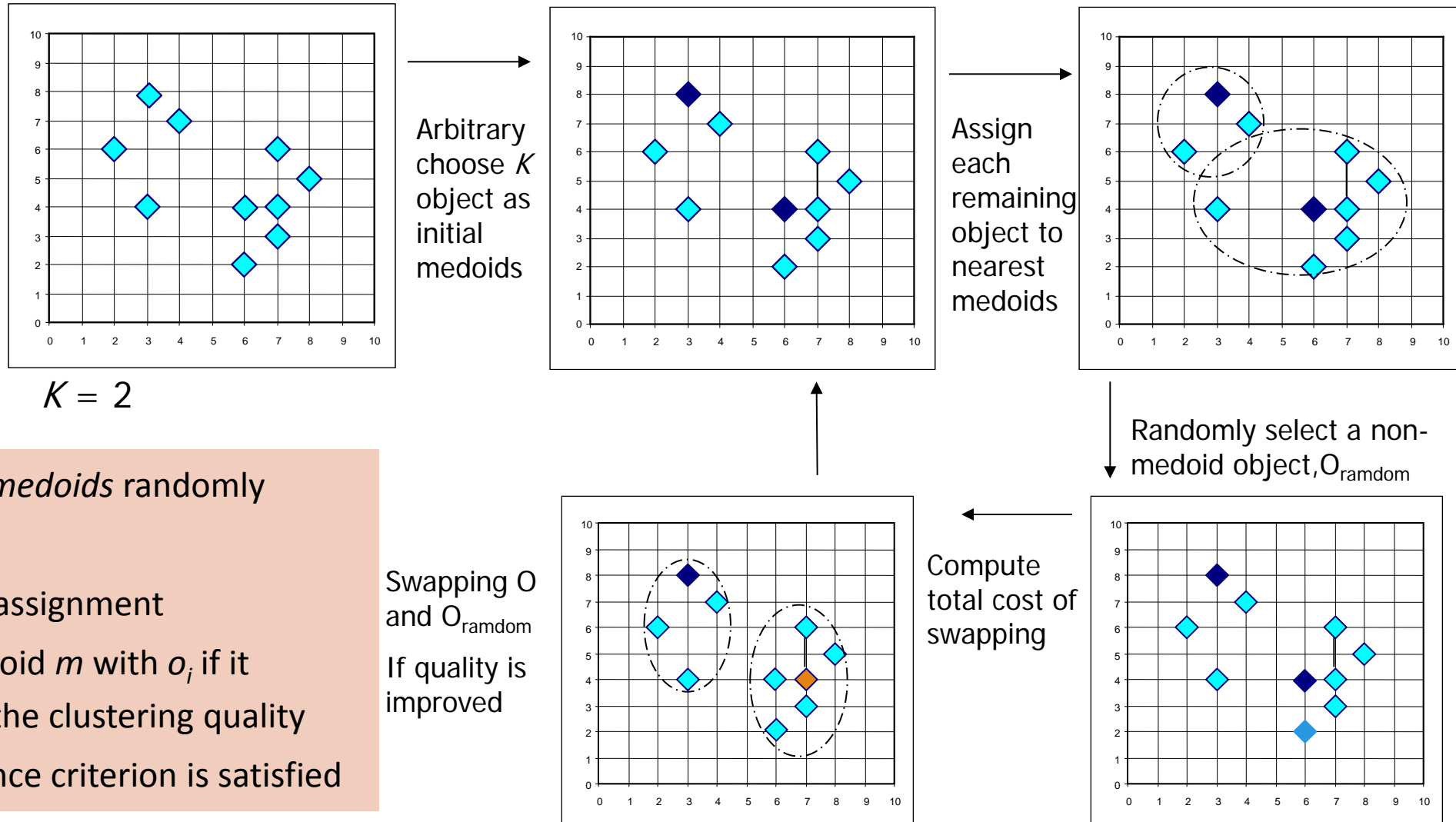
# Handling Outliers: From *K-Means* to *K-Medoids*

---

- ❑ 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



# 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



The background features a collage of data visualization elements. At the top and bottom, there are network graphs with nodes and edges in shades of red, orange, and green. On the left side, there is a vertical strip showing a heatmap or a series of small plots with orange and red data points. The central text is overlaid on a white, angular geometric shape.

# **The *K-Medians* and *K-Modes* Clustering Methods**

# ***K-Medians: Handling Outliers by Computing Medians***

---

- ❑ Medians are less sensitive to outliers than means
  - ❑ Think of the median salary vs. mean salary of a large firm when adding a few top executives!
- ❑ **K-Medians**: Instead of taking the **mean** value of the object in a cluster as a reference point, **medians** are used ( $L_1$ -norm as the distance measure)
- ❑ The criterion function for the *K-Medians* algorithm:
$$S = \sum_{k=1}^K \sum_{x_{ij} \in C_k} |x_{ij} - med_{kj}|$$
- ❑ The *K-Medians* clustering algorithm:
  - ❑ Select  $K$  points as the initial representative objects (i.e., as initial  $K$  medians)
  - ❑ **Repeat**
    - ❑ Assign every point to its nearest median
    - ❑ Re-compute the median using the median of each individual feature
  - ❑ **Until** convergence criterion is satisfied



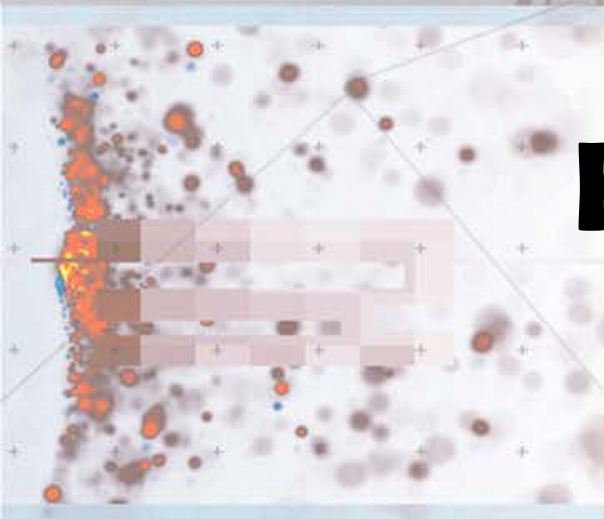
# K-Modes: Clustering Categorical Data

---

- ❑ *K-Means* cannot handle non-numerical (categorical) data
  - ❑ Mapping categorical value to 1/0 cannot generate quality clusters for high-dimensional data
- ❑ ***K-Modes***: An extension to *K-Means* by replacing means of clusters with ***modes***
- ❑ Dissimilarity measure between object  $X$  and the center of a cluster  $Z$ 
  - ❑  $\Phi(x_j, z_j) = 1 - n_j^r/n_l$  when  $x_j = z_j$  ; 1 when  $x_j \neq z_j$ 
    - ❑ where  $z_j$  is the categorical value of attribute  $j$  in  $Z_l$ ,  $n_l$  is the number of objects in cluster  $l$ , and  $n_j^r$  is the number of objects whose attribute value is  $r$
- ❑ This dissimilarity measure (distance function) is **frequency-based**
- ❑ Algorithm is still based on iterative *object cluster assignment* and *centroid update*
- ❑ A ***fuzzy K-Modes*** method is proposed to calculate a ***fuzzy cluster membership value*** for each object to each cluster
- ❑ A mixture of categorical and numerical data: Using a ***K-Prototype*** method

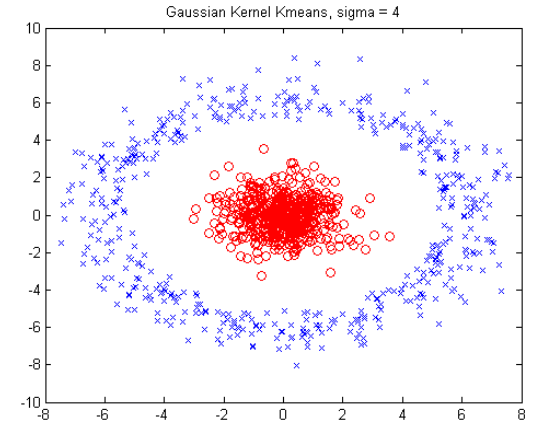


# Kernel K-Means Clustering



# ***Kernel K-Means Clustering***

- ❑ *Kernel K-Means* can be used to detect non-convex clusters
  - ❑ *K-Means* can only detect clusters that are linearly separable
- ❑ Idea: Project data onto the high-dimensional kernel space, and then perform *K-Means* clustering
  - ❑ Map data points in the input space onto a high-dimensional feature space using the kernel function
  - ❑ Perform *K-Means* on the mapped feature space
- ❑ Computational complexity is higher than K-Means
  - ❑ Need to compute and store  $n \times n$  kernel matrix generated from the kernel function on the original data
- ❑ The widely studied spectral clustering can be considered as a variant of Kernel K-Means clustering



# Kernel Functions and Kernel K-Means Clustering

## □ Typical kernel functions:

□ Polynomial kernel of degree  $h$ :  $K(\mathbf{X}_i, \mathbf{X}_j) = (\mathbf{X}_i \cdot \mathbf{X}_j + 1)^h$

□ Gaussian radial basis function (RBF) kernel:  $K(\mathbf{X}_i, \mathbf{X}_j) = e^{-\|\mathbf{X}_i - \mathbf{X}_j\|^2 / 2\sigma^2}$

□ Sigmoid kernel:  $K(\mathbf{X}_i, \mathbf{X}_j) = \tanh(\kappa \mathbf{X}_i \cdot \mathbf{X}_j - \delta)$

□ The formula for kernel matrix  $K$  for any two points  $x_i, x_j \in C_k$  is  $K_{x_i x_j} = \phi(x_i) \bullet \phi(x_j)$

□ The SSE criterion of *kernel K-means*: 
$$SSE(C) = \sum_{k=1}^K \sum_{x_i \in C_k} \|\phi(x_i) - c_k\|^2$$

□ The formula for the cluster centroid:

$$c_k = \frac{\sum_{x_i \in C_k} \phi(x_i)}{|C_k|}$$

□ Clustering can be performed without the actual individual projections  $\phi(x_i)$  and  $\phi(x_j)$  for the data points  $x_i, x_j \in C_k$

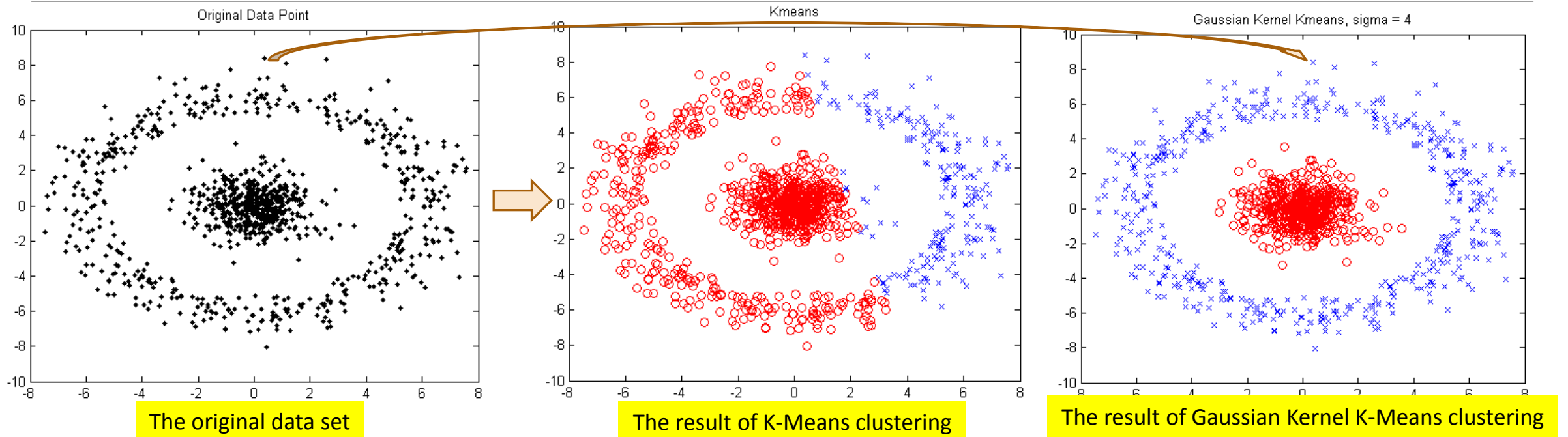


# Example: Kernel Functions and Kernel K-Means Clustering

- Gaussian radial basis function (RBF) kernel:  $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma^2}$
- Suppose there are 5 original 2-dimensional points:
  - $x_1(0, 0), x_2(4, 4), x_3(-4, 4), x_4(-4, -4), x_5(4, -4)$
- If we set  $\sigma$  to 4, we will have the following points in the kernel space
  - E.g.,  $\|x_1 - x_2\|^2 = (0 - 4)^2 + (0 - 4)^2 = 32$ , therefore,  $K(x_1, x_2) = e^{-\frac{32}{2 \cdot 4^2}} = e^{-1}$

Original Space			RBF Kernel Space ( $\sigma = 4$ )				
	$x$	$y$	$K(x_i, x_1)$	$K(x_i, x_2)$	$K(x_i, x_3)$	$K(x_i, x_4)$	$K(x_i, x_5)$
$x_1$	0	0	0	$e^{-\frac{4^2+4^2}{2 \cdot 4^2}} = e^{-1}$	$e^{-1}$	$e^{-1}$	$e^{-1}$
$x_2$	4	4	$e^{-1}$	0	$e^{-2}$	$e^{-4}$	$e^{-2}$
$x_3$	-4	4	$e^{-1}$	$e^{-2}$	0	$e^{-2}$	$e^{-4}$
$x_4$	-4	-4	$e^{-1}$	$e^{-4}$	$e^{-2}$	0	$e^{-2}$
$x_5$	4	-4	$e^{-1}$	$e^{-2}$	$e^{-4}$	$e^{-2}$	0

# Example: Kernel K-Means Clustering



- ❑ The above data set cannot generate quality clusters by K-Means since it contains non-convex clusters
- ❑ Gaussian RBF Kernel transformation maps data to a kernel matrix  $K$  for any two points  $x_i, x_j$ :  $K_{x_i x_j} = \phi(x_i) \bullet \phi(x_j)$  and Gaussian kernel:  $K(x_i, x_j) = e^{-\|x_i - x_j\|^2 / 2\sigma^2}$
- ❑ K-Means clustering is conducted on the mapped data, generating quality clusters

# Recommended Readings

---

- ❑ J. MacQueen. Some Methods for Classification and Analysis of Multivariate Observations. In *Proc. of the 5th Berkeley Symp. on Mathematical Statistics and Probability*, 1967
- ❑ S. Lloyd. Least Squares Quantization in PCM. *IEEE Trans. on Information Theory*, 28(2), 1982
- ❑ A. K. Jain and R. C. Dubes. Algorithms for Clustering Data. Prentice Hall, 1988
- ❑ L. Kaufman and P. J. Rousseeuw. Finding Groups in Data: An Introduction to Cluster Analysis. John Wiley & Sons, 1990
- ❑ R. Ng and J. Han. Efficient and Effective Clustering Method for Spatial Data Mining. VLDB'94
- ❑ B. Schölkopf, A. Smola, and K. R. Müller. Nonlinear Component Analysis as a Kernel Eigenvalue Problem. *Neural computation*, 10(5):1299–1319, 1998
- ❑ I. S. Dhillon, Y. Guan, and B. Kulis. Kernel K-Means: Spectral Clustering and Normalized Cuts. *KDD'04*
- ❑ D. Arthur and S. Vassilvitskii. K-means++: The Advantages of Careful Seeding. *SODA'07*
- ❑ C. K. Reddy and B. Vinzamuri. A Survey of Partitional and Hierarchical Clustering Algorithms, in (Chap. 4) Aggarwal and Reddy (eds.), *Data Clustering: Algorithms and Applications*. CRC Press, 2014
- ❑ M. J. Zaki and W. Meira, Jr.. *Data Mining and Analysis: Fundamental Concepts and Algorithms*. Cambridge Univ. Press, 2014