

## Related Topic: Clustering

**Input:**  $n$  sites:  $S = \{s_1, s_2, \dots, s_n\}$

**Output:** Locations of  $k$  centers:  $C = \{c_1, c_2, \dots, c_k\}$

**Objective:** Minimize the total squared distance from each site to the nearest center.

$$\text{Minimize } \sum_{s \in S} \text{dist}(s, C)^2$$

$\text{dist}(s, C)$ : Distance from  $s$  to the nearest center.

$$\text{dist}(s, C) = \text{Min}_{c \in C} \{\text{dist}(s, c)\}$$

## Reformulation

■ site ( $n$  sites)

● center ( $k$  centers)

- (1) Divide the  $n$  sites into  $k$  clusters based on the nearest center.

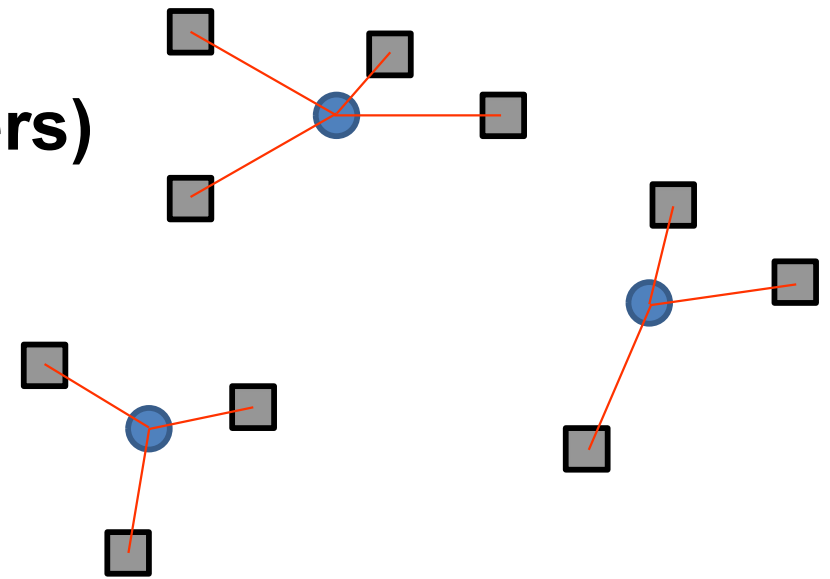
$$S = \{s_1, s_2, \dots, s_n\}$$




$$S = S_1 \cup S_2 \cup \dots \cup S_k$$

- (2) Reformulate the objective function as follows:

$$\text{Minimize } \sum_{j=1}^k \sum_{s \in S_j} \text{dist}(s, c_j)^2$$



$$\text{Minimize } \sum_{j=1}^k \sum_{s \in S_j} \text{dist}(s, c_j)^2$$

$$S = \{s_1, s_2, \dots, s_n\}$$


$$S = S_1 \cup S_2 \cup \dots \cup S_k$$

$s$  and  $c_j$ : Points in the 2D space.

**k-means Algorithm:** Iterate the following two steps from a random partition of  $S$  into  $k$  subsets:  $S_1, S_2, \dots, S_k$

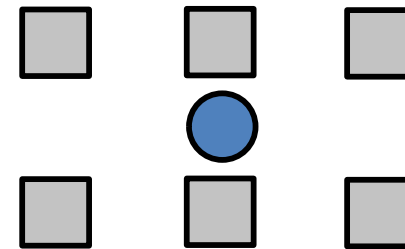
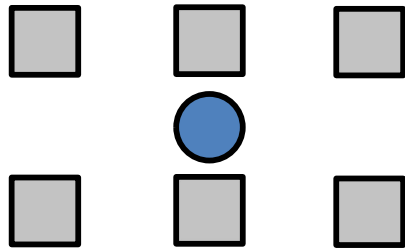
$$(i) \quad c_j = \frac{1}{|S_j|} \sum_{s \in S_j} s, \quad j = 1, 2, \dots, k.$$

$$(ii) \quad S_j = \{s \mid \text{dist}(s, c_j) = \min_{l=1, \dots, k} \text{dist}(s, c_l)\}, \quad j = 1, 2, \dots, k.$$

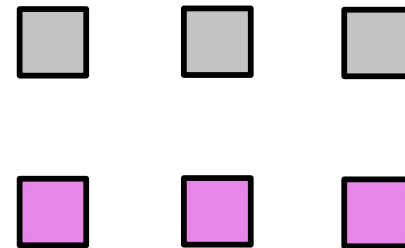
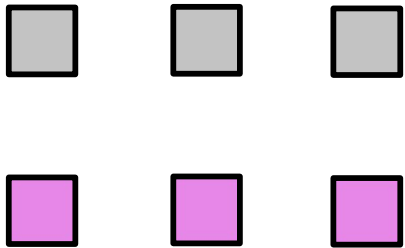
### Exercise 7-1:

In the k-means algorithm, we can start with (i) using an initial partition  $\{S_1, S_2, \dots, S_k\}$  or with (ii) using initial centers  $\{c_1, c_2, \dots, c_k\}$ . Design a good initialization method for the k-means algorithm with (i) or (ii).

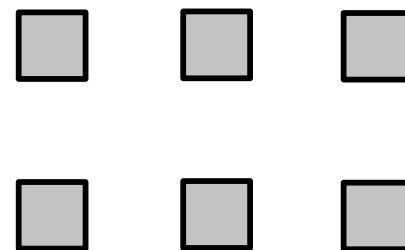
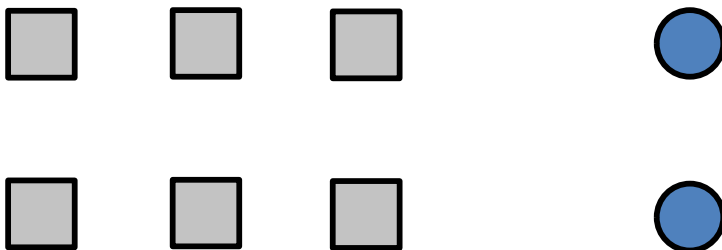
## Example (12 sites and 2 centers)



If we start with the following partitions:



If we start with the following centers:

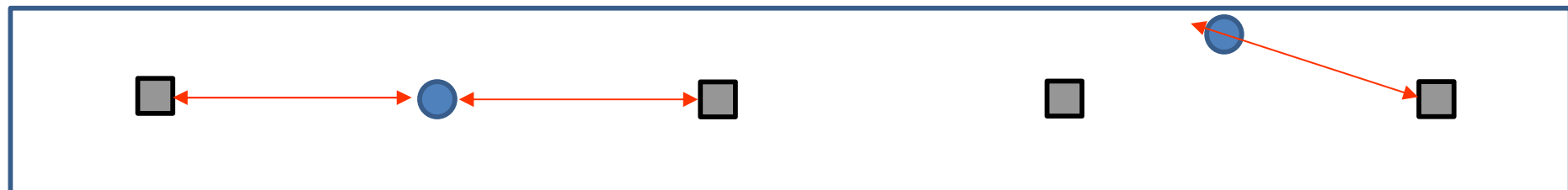
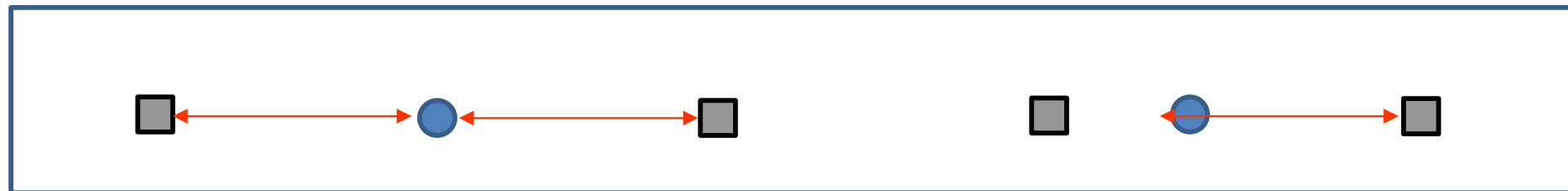
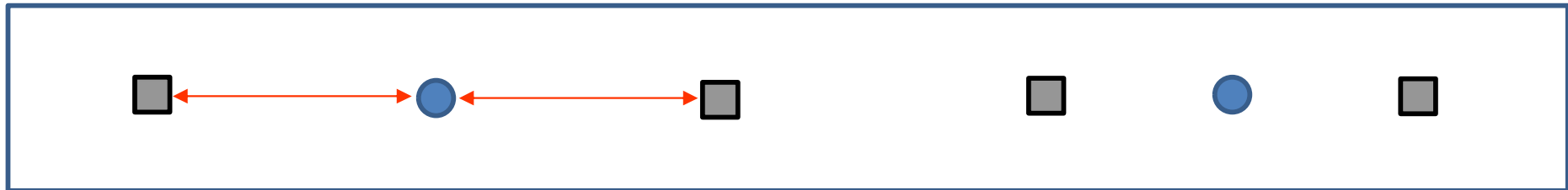


# Difficulty of “Min-Max” objective function: (“minimize the worst case” objective function)

Minimization of the maximum distance from each site to the nearest center.

$$\text{Minimize } \max_{s \in S} \text{dist}(s, C)$$

**All the following solutions are optimal (for  $k = 2$ ).**



# Comparison of Problems:

- (1) Minimization of the maximum distance from each site to the nearest center.

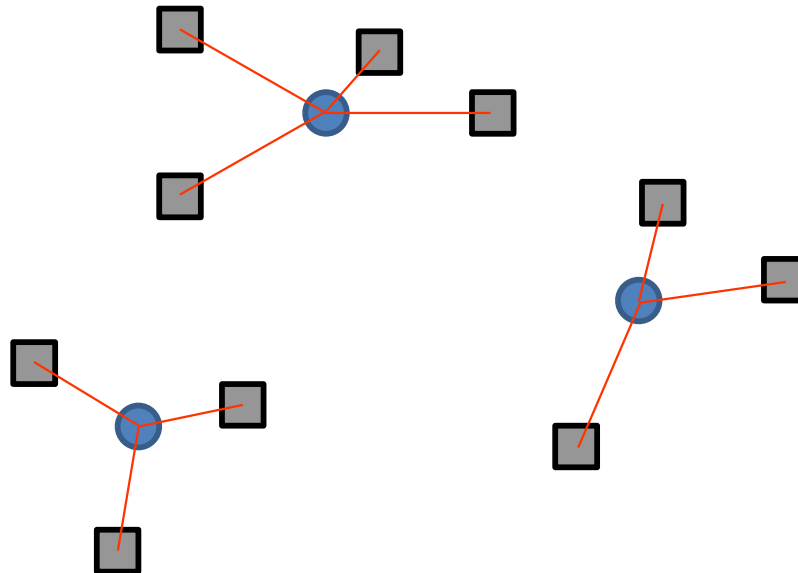
$$\text{Minimize } \max_{s \in S} \text{dist}(s, C)$$

- (2) Minimization of the total squared distance from each site to the nearest center

$$\text{Minimize } \sum_{s \in S} \text{dist}(s, C)^2$$

**Q. Which is a better problem formulation?**

■ site  
● center

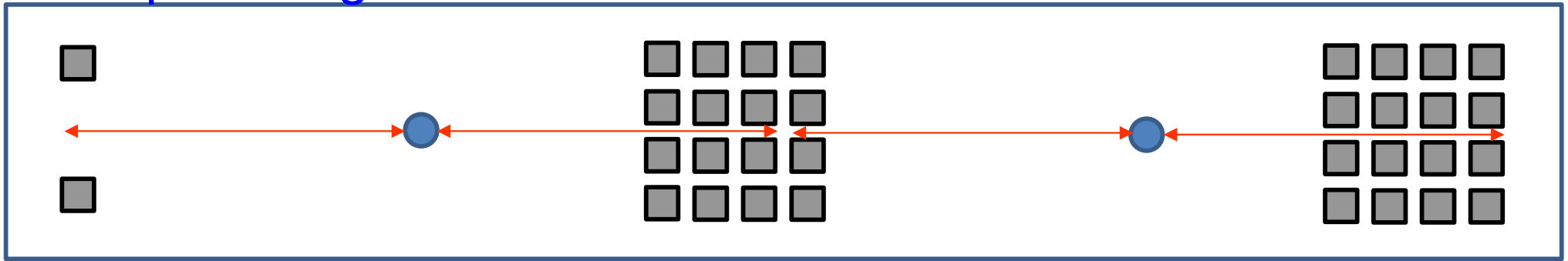


# Comparison of Problems:

(1) **Minimization of the maximum distance** from each site to the nearest center.

$$\text{Minimize } \max_{s \in S} \text{dist}(s, C)$$

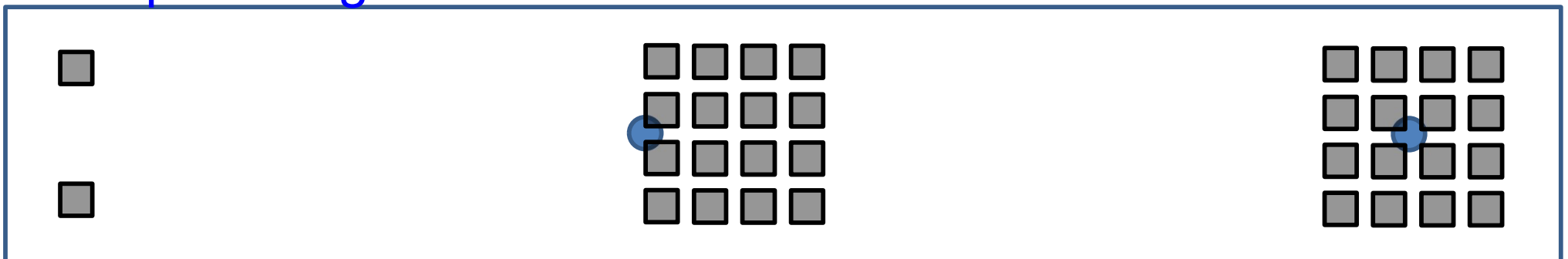
Example of a good solution



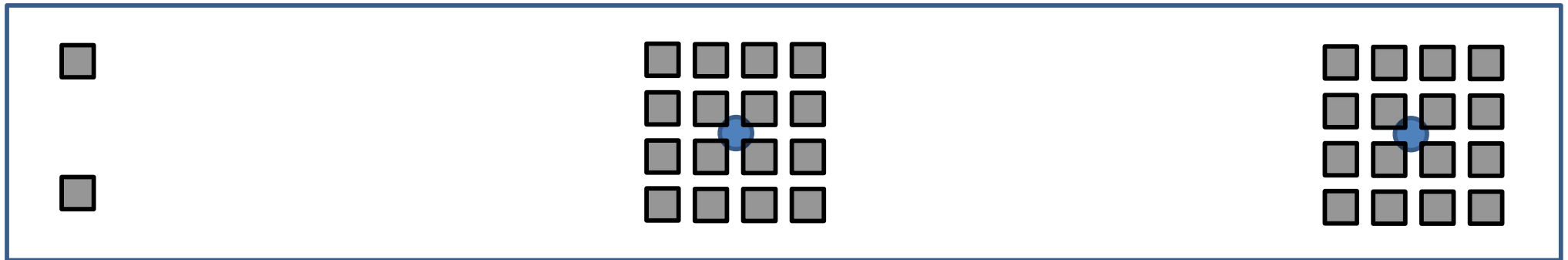
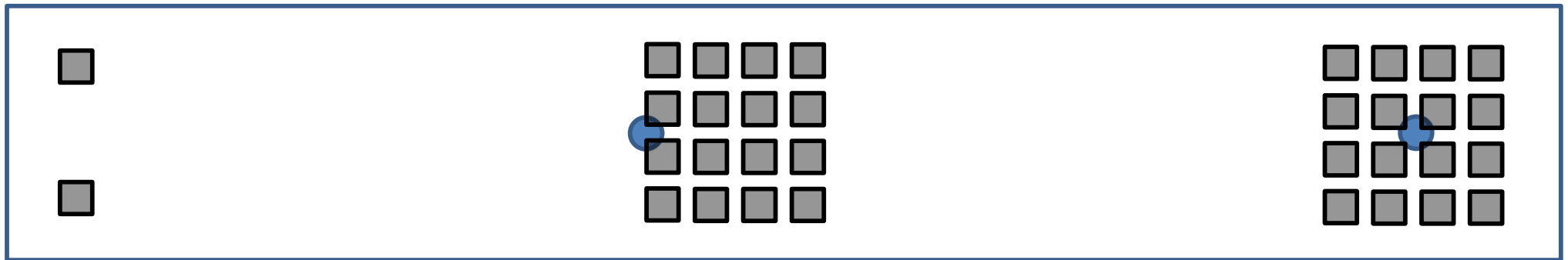
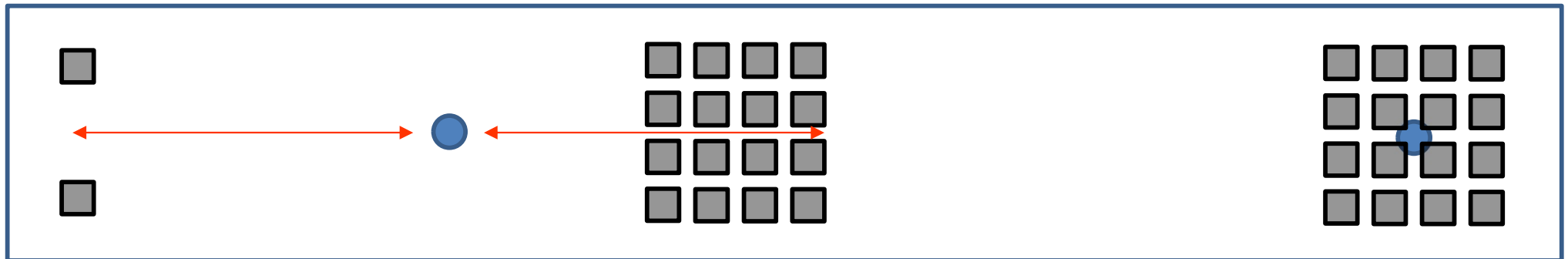
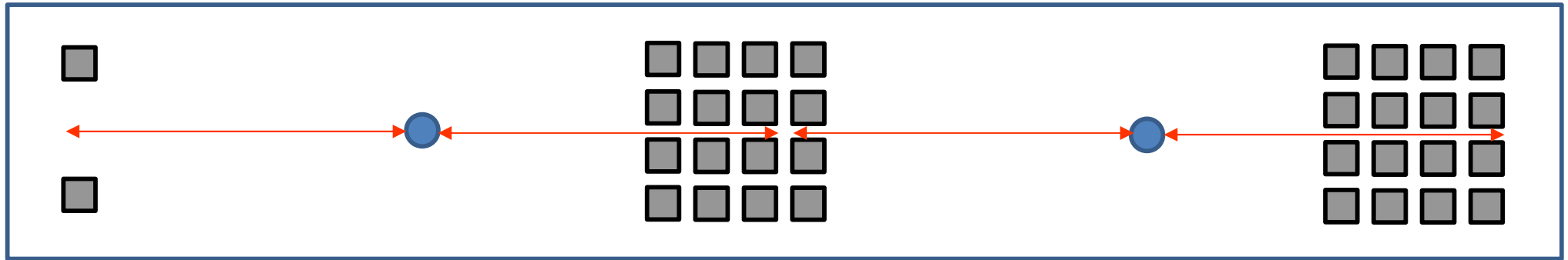
(2) **Minimization of the total squared distance** from each site to the nearest center

$$\text{Minimize } \sum_{s \in S} \text{dist}(s, C)^2$$

Example of a good solution



# Which is the best solution?





# Comparison of Algorithms:

- (1) Minimization of the maximum distance from each site to the nearest center.

$$\text{Minimize } \max_{s \in S} \text{dist}(s, C)$$

## Center Selection Algorithm:

**Simple heuristics (a greedy algorithm)**  
**2-Approximation algorithm**

- (2) Minimization of the total squared distance from each site to the nearest center

$$\text{Minimize } \sum_{s \in S} \text{dist}(s, C)^2$$

## K-means Algorithm

**Iterative adjustment algorithm**  
**(iterations of two greedy algorithms)**  
**Not an exact optimization algorithm**

# Comparison of Algorithms:

- (2) Minimization of the total squared distance from each site to the nearest center

$$\text{Minimize } \sum_{s \in S} \text{dist}(s, C)^2$$

## k-means Algorithm

**Iterative adjustment algorithm**  
**(iterations of two greedy algorithms)**  
**Not an exact optimization algorithm**

- (3) Minimization of the total distance from each site to the nearest center

$$\text{Minimize } \sum_{s \in S} \text{dist}(s, C)$$

## k-medoids Algorithm

**Iterative adjustment algorithm**  
**(iterations of two greedy algorithms)**  
**Not an exact optimization algorithm**

## **k-means Algorithm (start with $k$ centers):**

Iterate the following two steps from  $k$  centers  $c_1, c_2, \dots, c_k$ :

$$(i) S_j = \{s \mid \text{dist}(s, c_j) = \min_{l=1,2,\dots,k} \text{dist}(s, c_l)\}, \quad j = 1, 2, \dots, k.$$

$$(ii) c_j = \frac{1}{|S_j|} \sum_{s \in S_j} s, \quad j = 1, 2, \dots, k.$$

## **k-medoids Algorithm (start with $k$ centers):**

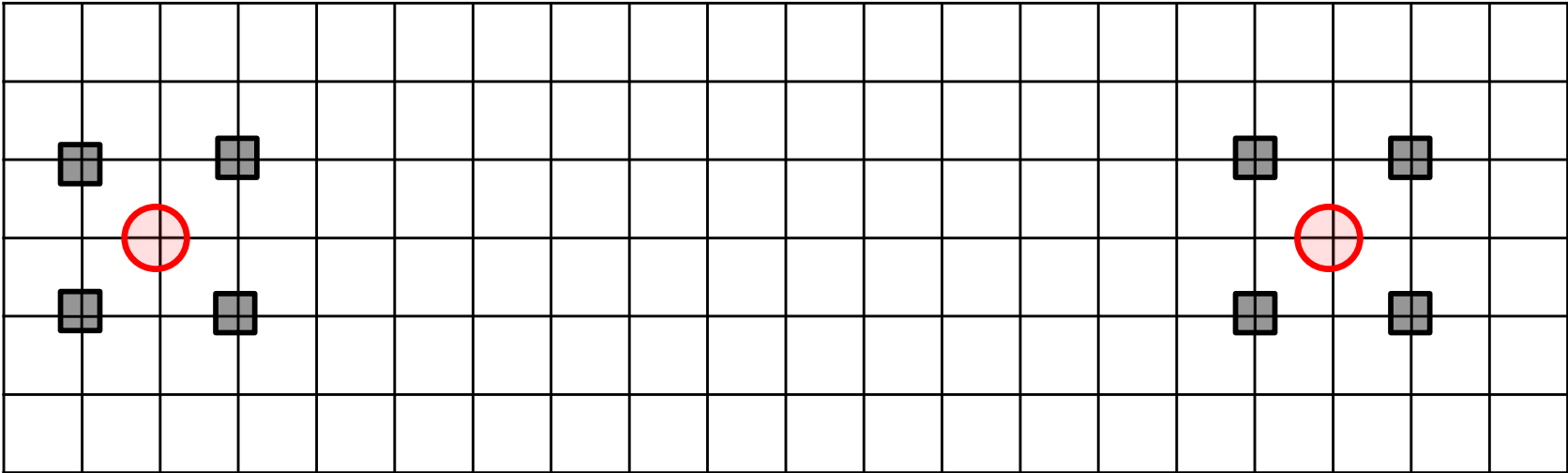
Iterate the following two steps from  $k$  centers  $c_1, c_2, \dots, c_k$ :

$$(i) S_j = \{s \mid \text{dist}(s, c_j) = \min_{l=1,2,\dots,k} \text{dist}(s, c_l)\}, \quad j = 1, 2, \dots, k.$$

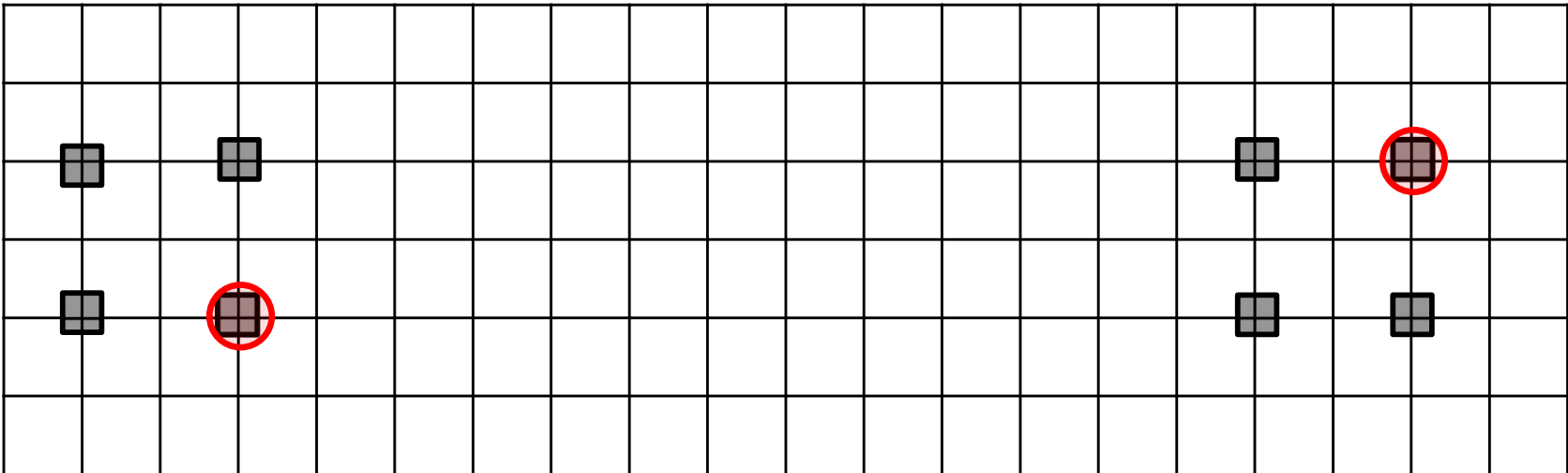
(ii) In each cluster  $S_j$  ( $j = 1, 2, \dots, k$ ), choose the site  $c_j$  which minimizes the total distance to all the other sites in cluster  $S_j$ .

**Medoid:** (mathematics) A mathematically representative object in a set of objects; it has the smallest average dissimilarity to all other objects in the set.

# k-means Algorithm



# k-medoids Algorithm

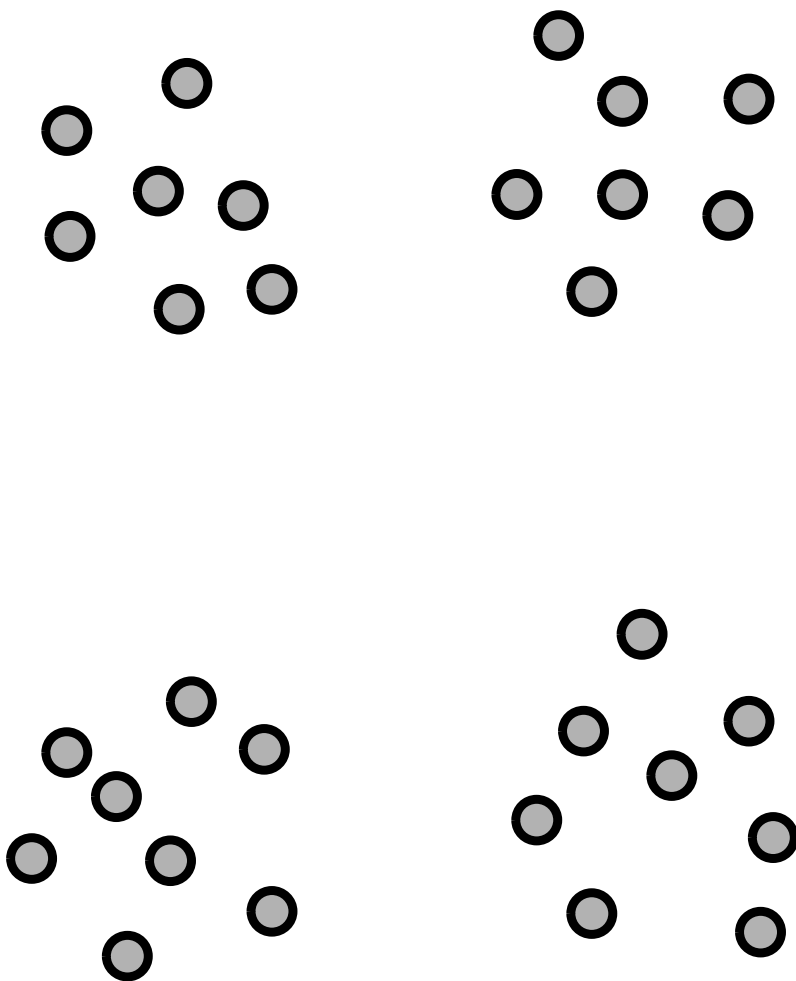


## **Exercise 7-2:**

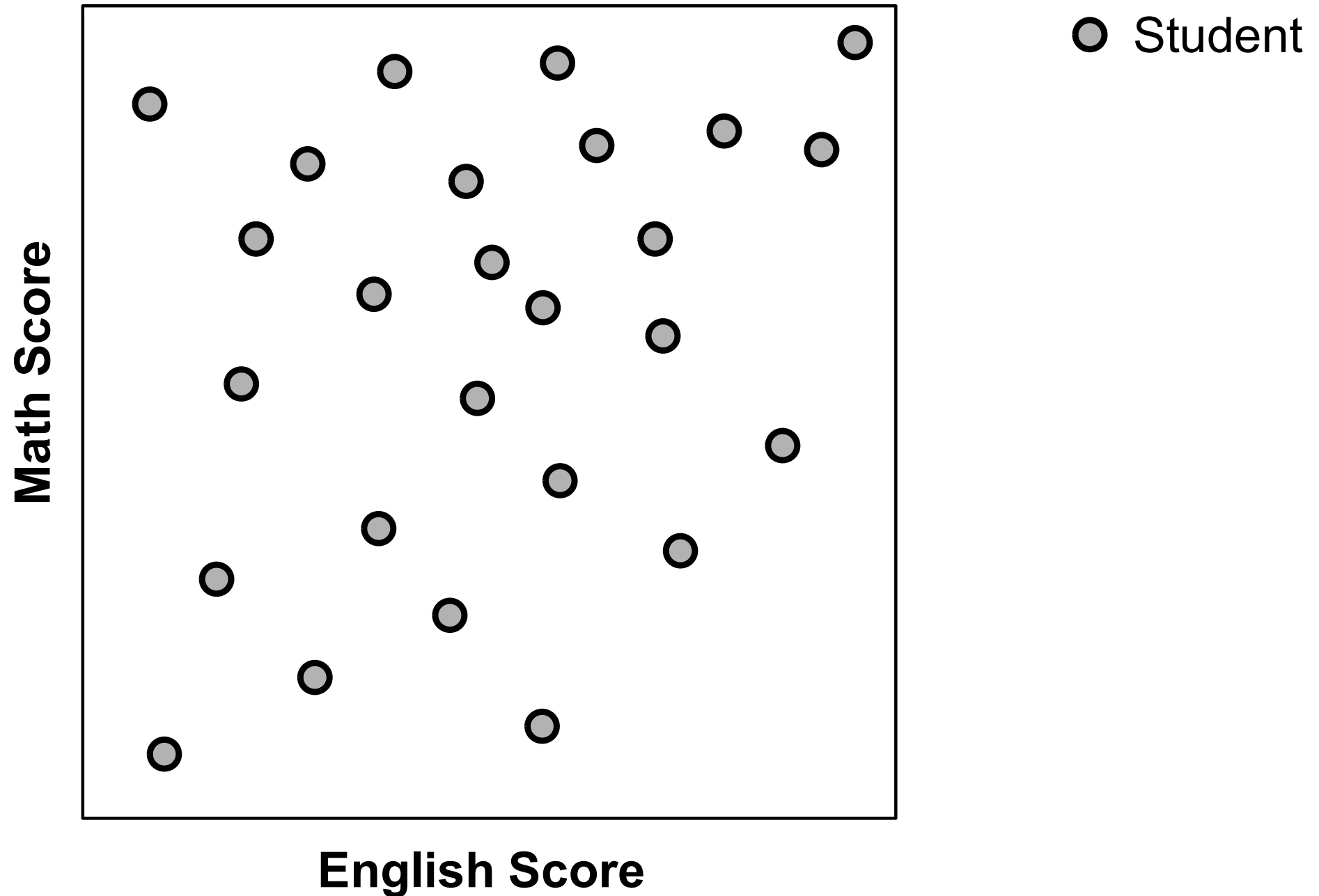
**Clearly demonstrate the difference between the *k*-means algorithm and the k-medoids algorithm using a test data set (i.e., create a test data set which can be used for clearly demonstrating the difference between the *k*-means algorithm and the k-medoids algorithm).**

# Recent Hot Topic: Fairness in Clustering

Problem: Find two clusters ( $k = 2$ )



# Classification Problem: Accept 10 Students



## Related Topic: Fuzzy Clustering

**Input:**  $n$  sites:  $S = \{s_1, s_2, \dots, s_n\}$ , and a constant  $m$  ( $m > 1$ )

**Output:** Locations of  $k$  centers:  $C = \{c_1, c_2, \dots, c_k\}$

Membership of  $s_i$  to  $c_j$ :  $\mu_{ij}$  ( $i = 1, \dots, n$ ;  $j = 1, \dots, k$ )

**Objective:** Minimize the weighted total squared distance from each site to each center.

$$\text{Minimize } \sum_{i=1}^n \sum_{j=1}^k (\mu_{ij})^m \text{dist}(s_i, c_j)^2$$

where

$$0 \leq \mu_{ij} \leq 1, \quad i = 1, 2, \dots, n; \quad j = 1, 2, \dots, k$$

$$\sum_{j=1}^k \mu_{ij} = 1, \quad i = 1, 2, \dots, n$$



## Related Topic: Fuzzy Clustering

**Input:**  $n$  sites:  $S = \{s_1, s_2, \dots, s_n\}$ , and a constant  $m$  ( $m > 1$ )

**Output:** Locations of  $k$  centers:  $C = \{c_1, c_2, \dots, c_k\}$

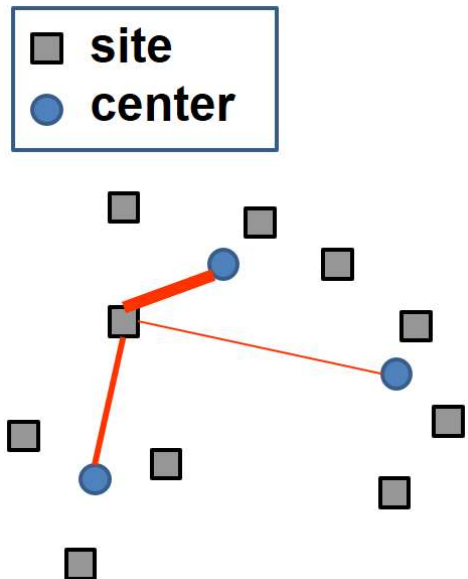
Membership of  $s_i$  to  $c_j$ :  $\mu_{ij}$  ( $i = 1, \dots, n$ ;  $j = 1, \dots, k$ )

**Objective:** Minimize the weighted total squared distance from each site to each center.

$$\text{Minimize } \sum_{i=1}^n \sum_{j=1}^k (\mu_{ij})^m \text{dist}(s_i, c_j)^2$$

where

Membership grade



$$0 \leq \mu_{ij} \leq 1, i = 1, 2, \dots, n; j = 1, 2, \dots, k$$

$$\sum_{j=1}^k \mu_{ij} = 1, i = 1, 2, \dots, n$$

## fuzzy c-means

$$\text{Minimize } \sum_{i=1}^n \sum_{j=1}^k (\mu_{ij})^m \text{dist}(s_i, c_j)^2$$

$$0 \leq \mu_{ij} \leq 1, i = 1, 2, \dots, n; j = 1, 2, \dots, k$$

$$\sum_{j=1}^k \mu_{ij} = 1, i = 1, 2, \dots, n$$

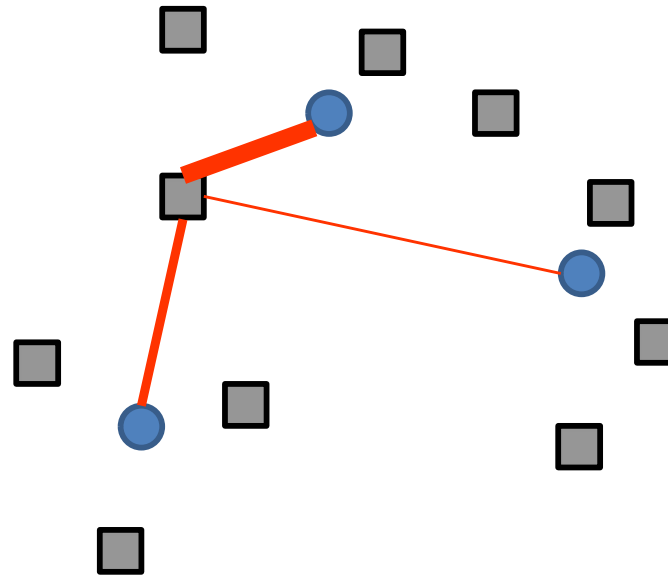
## k-means

$$\text{Minimize } \sum_{i=1}^n \sum_{j=1}^k (\mu_{ij})^1 \text{dist}(s_i, c_j)^2$$

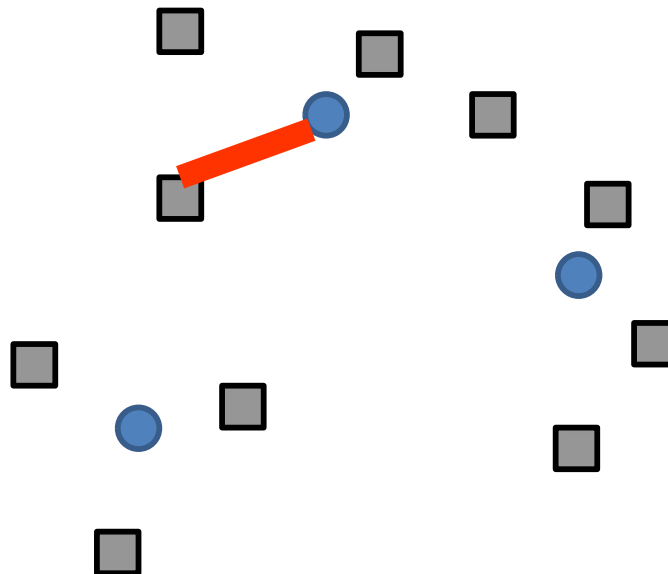
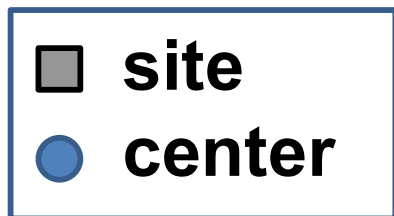
$$\mu_{ij} = 0 \text{ or } 1, i = 1, 2, \dots, n; j = 1, 2, \dots, k$$

$$\sum_{j=1}^k \mu_{ij} = 1, i = 1, 2, \dots, n$$

## fuzzy c-means



## *k*-means



**Fuzzy *c*-means Algorithm:** Iterate the following two steps from randomly specified values of  $\mu_{ij}$

$$\sum_{j=1}^k \mu_{ij} = 1, \quad i = 1, 2, \dots, n$$

$$(i) \quad c_j = \frac{\sum_{i=1}^n (\mu_{ij})^m s_i}{\sum_{i=1}^n (\mu_{ij})^m}, \quad j = 1, 2, \dots, k$$

$$(ii) \quad \mu_{ij} = \left[ \sum_{h=1}^k \left( \frac{\text{dist}(s_i, c_j)}{\text{dist}(s_i, c_h)} \right)^{\frac{2}{m-1}} \right]^{-1} \quad \text{for all } i \text{ and } j$$

$$\mu_{ij} = \left[ \sum_{h=1}^k \left( \frac{\text{dist}(s_i, c_j)}{\text{dist}(s_i, c_h)} \right)^{\frac{2}{m-1}} \right]^{-1} \quad \text{for all } i \text{ and } j$$

When  $m \rightarrow \infty$

$$\mu_{ij} \rightarrow \left[ \sum_{h=1}^k \left( \frac{\text{dist}(s_i, c_j)}{\text{dist}(s_i, c_h)} \right)^0 \right]^{-1} = 1/k$$

When  $m \rightarrow 1$

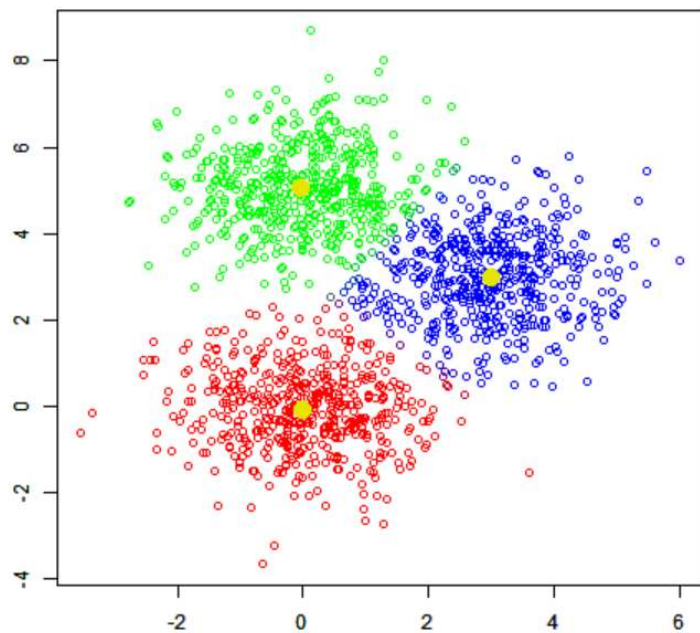
$$\mu_{ij} \rightarrow \left[ \sum_{h=1}^k \left( \frac{\text{dist}(s_i, c_j)}{\text{dist}(s_i, c_h)} \right)^{\infty} \right]^{-1} = 0 \text{ or } 1$$

### Exercise 7-3:

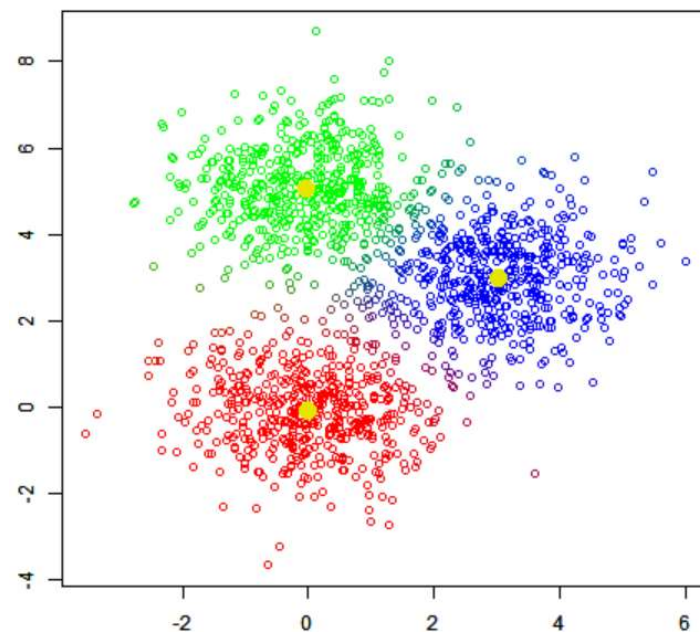
Clearly demonstrate the effects of  $m$  on the clustering results by the fuzzy c-means algorithm through computational experiments on a test data set (i.e., create a test data set which can be used for clearly demonstrating the effects of  $m$ ). **Try to create some beautiful figures.**

### Exercise 7-4:

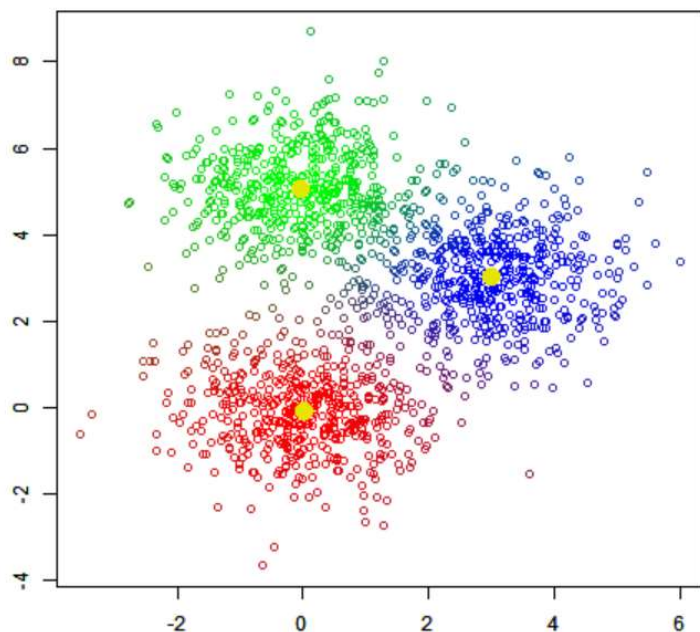
Clearly demonstrate the difference between the k-means algorithm and the fuzzy c-means algorithm through computational experiments on a test data set (i.e., create a test data set which can be used for clearly demonstrating the difference between the k-means algorithm and the fuzzy c-means algorithm). **Try to create some beautiful figures.**



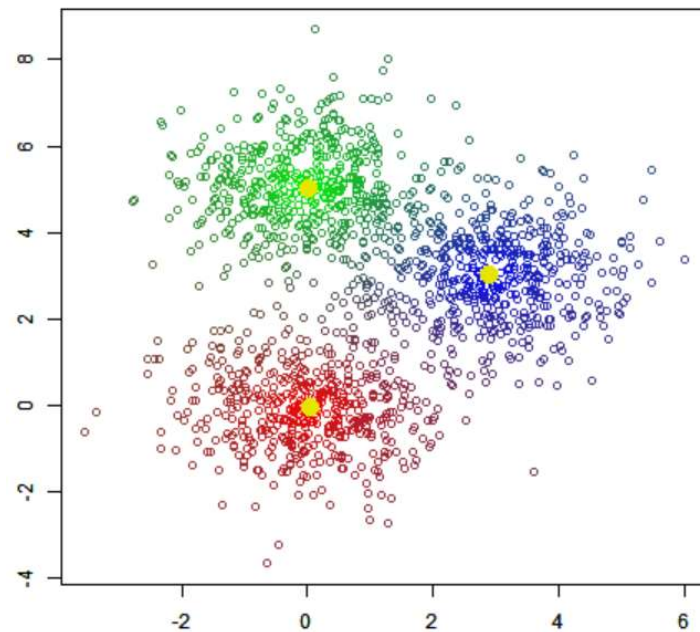
$m = 1.1$



$m = 1.5$



$m = 2$



$m = 3$

# Use of a user-defined hyper parameter ( $m$ in fuzzy c-means)

## Positive Aspects:

A more desirable result can be obtained by appropriately specifying the value of  $m$  (than the case of the fixed value of  $m$ )

Different results can be obtained by examining different values of  $m$  (we can choose one of them based on our preference).

## Negative Aspects:

It is not always easy to appropriately specify the value of  $m$ .

Undesirable results can be obtained when the value of  $m$  is inappropriate.



**Example:** If  $m$  is specified as  $m = 1$ , the algorithm does not work.

$$(i) \quad c_j = \frac{\sum_{i=1}^n (\mu_{ij})^m s_i}{\sum_{i=1}^n (\mu_{ij})^m}, \quad j = 1, 2, \dots, k$$

$$(ii) \quad \mu_{ij} = \left[ \sum_{h=1}^k \left( \frac{\text{dist}(s_i, c_j)}{\text{dist}(s_i, c_h)} \right)^{\frac{2}{m-1}} \right]^{-1} \quad \text{for all } i \text{ and } j$$