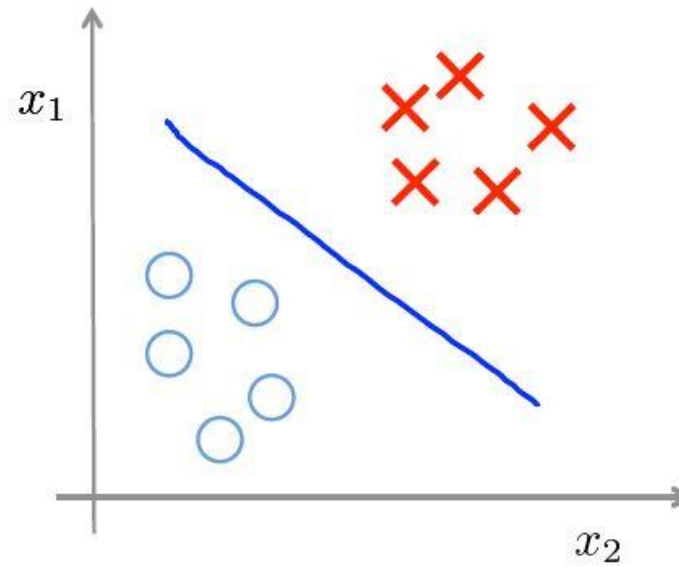# Clustering

Unsupervised learning introduction

# Supervised learning
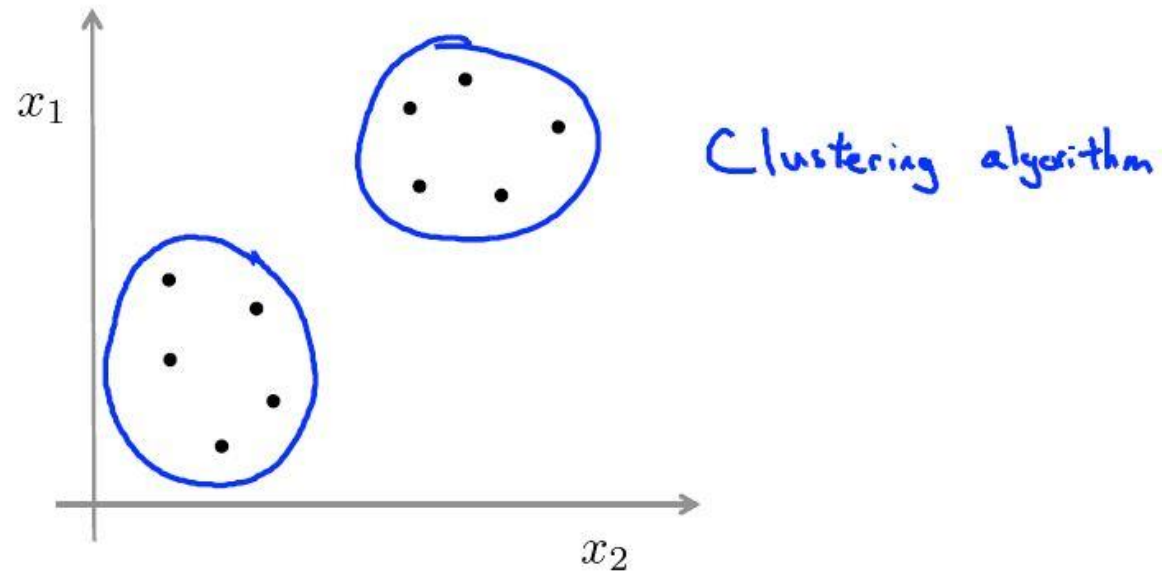


Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \ldots, (x^{(m)}, y^{(m)})\}$

# Unsupervised learning



Clustering algorithm

Training set: $\{x^{(1)}, x^{(2)}, x^{(3)}, \ldots, x^{(m)}\}$
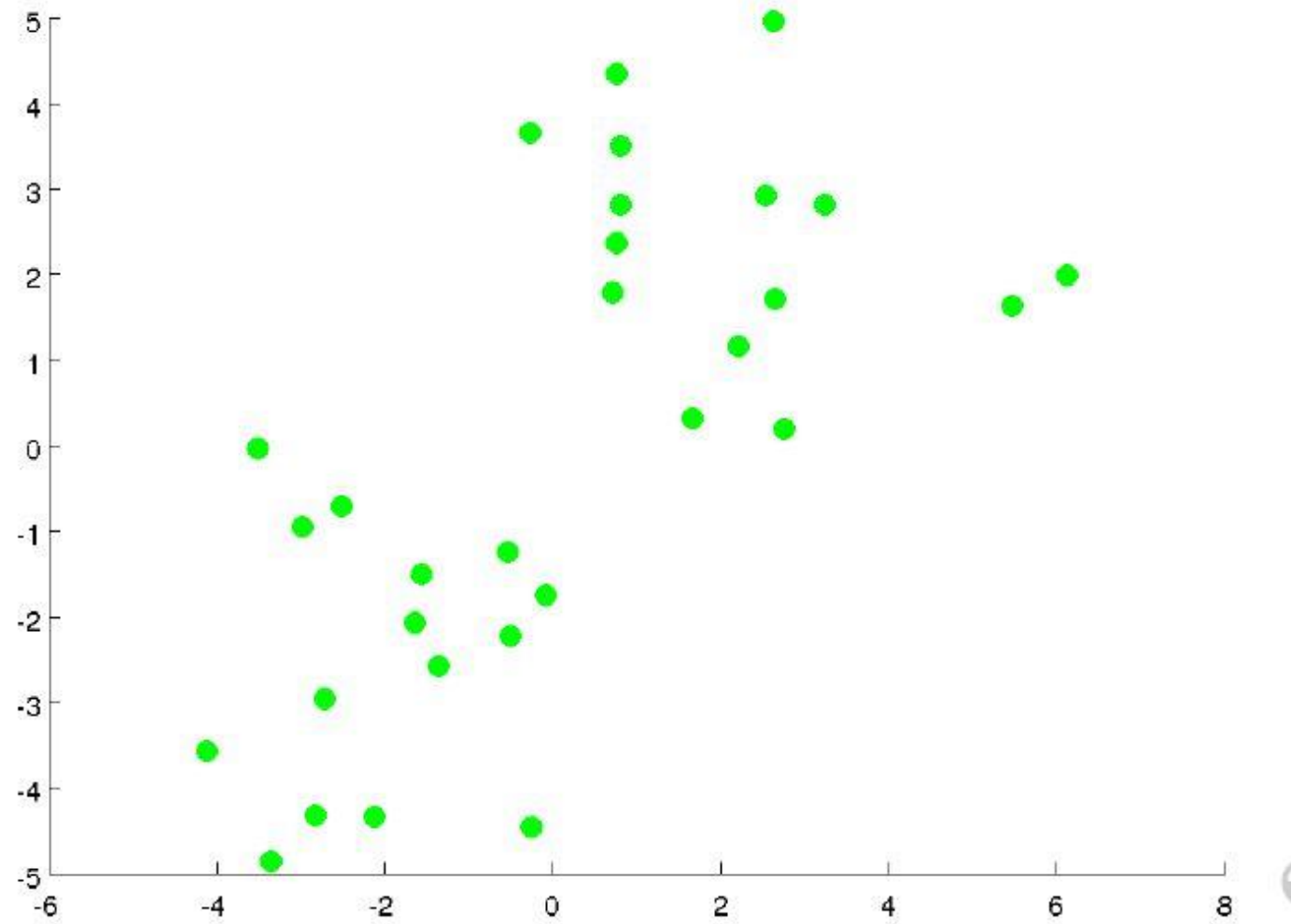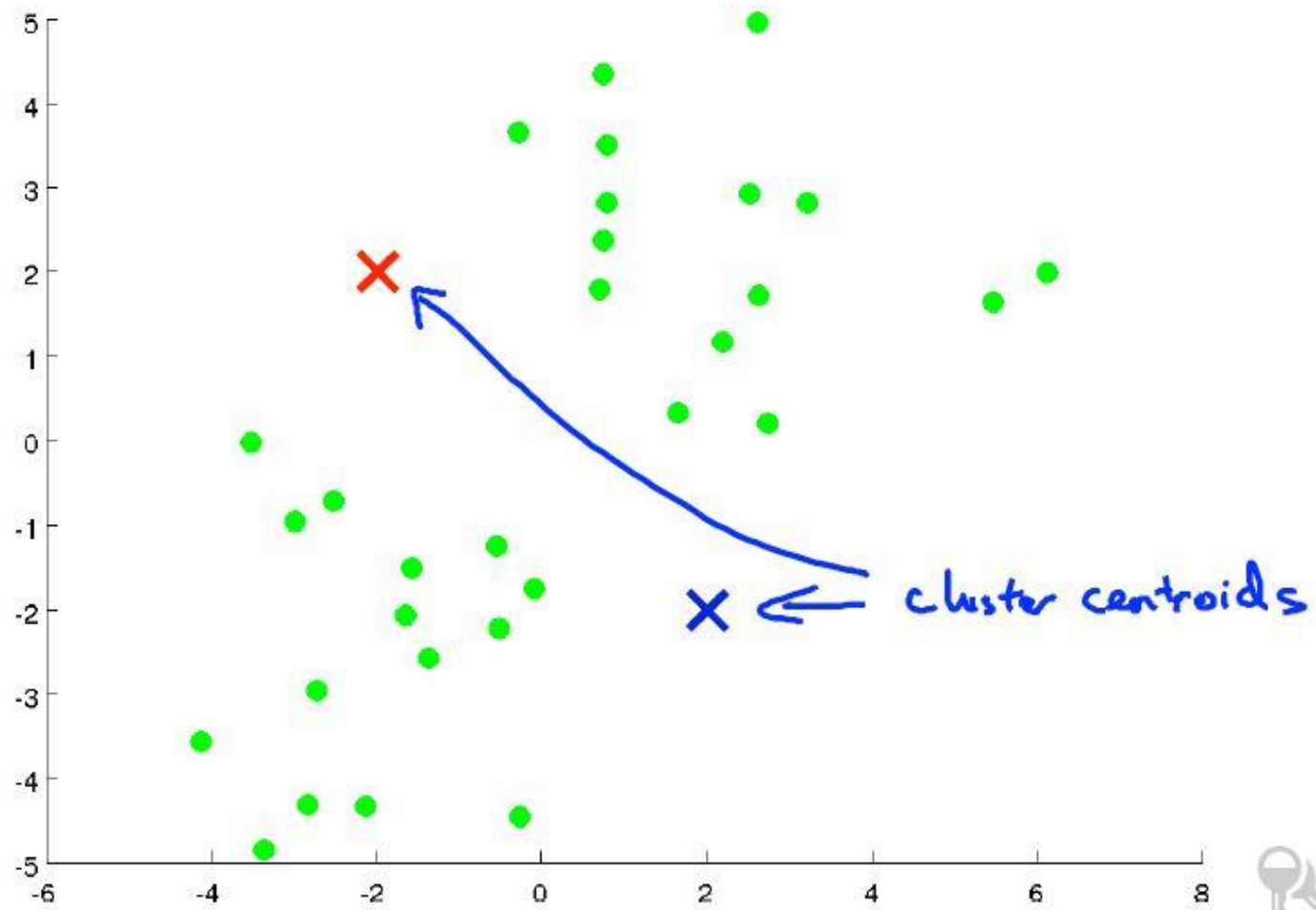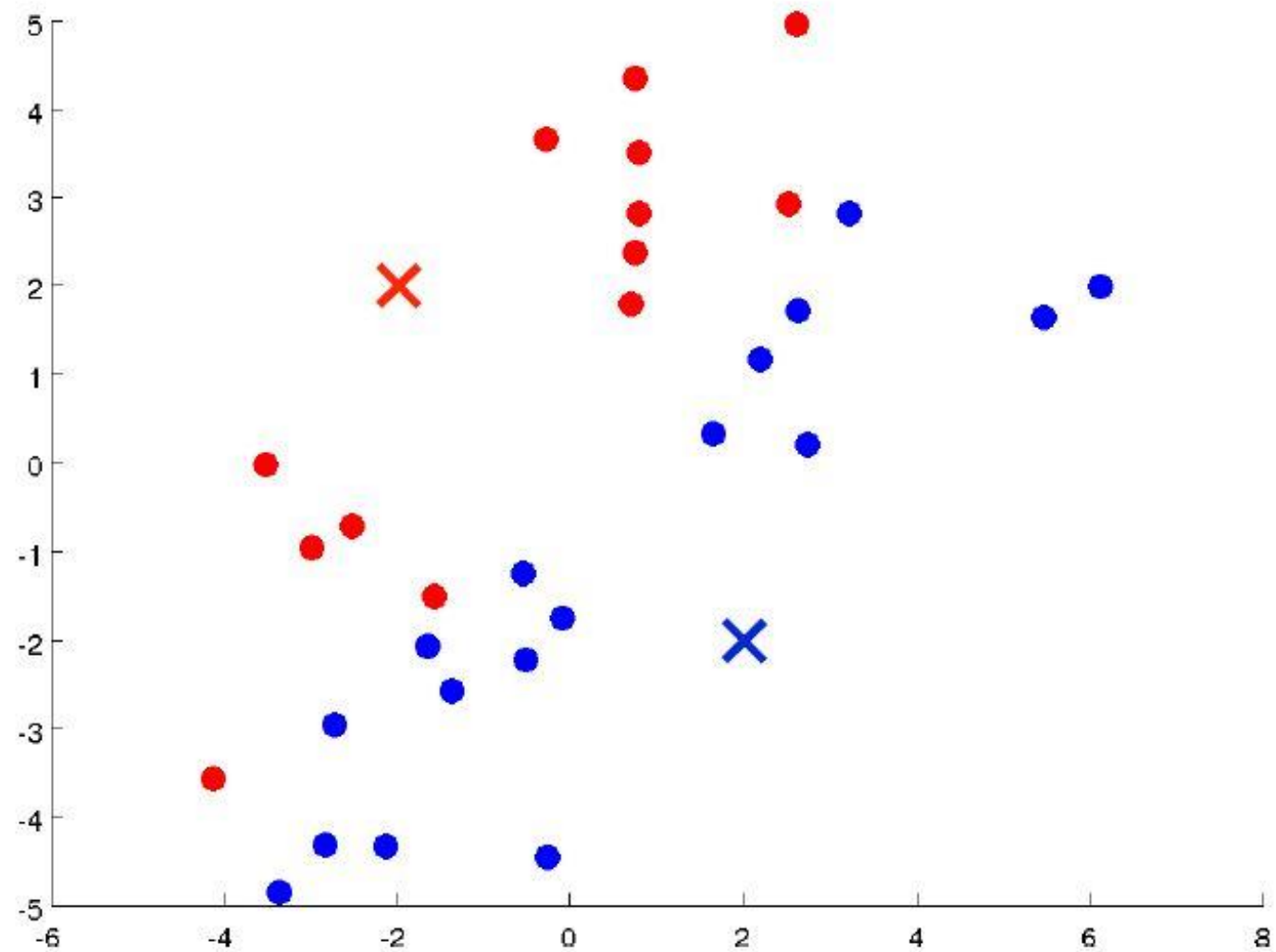
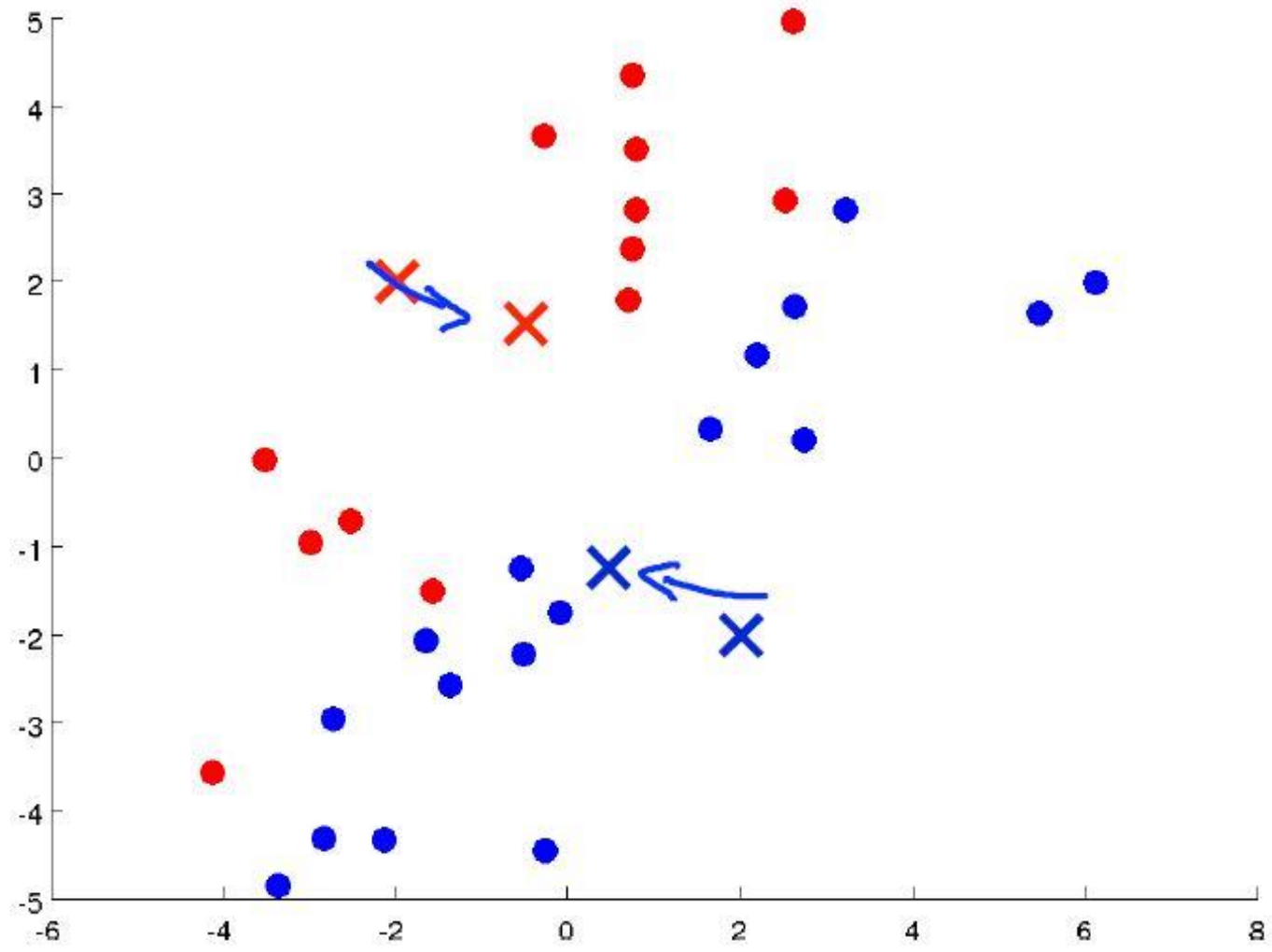# Applications of clustering



→ Market segmentation

→ Social network analysis

# Clustering
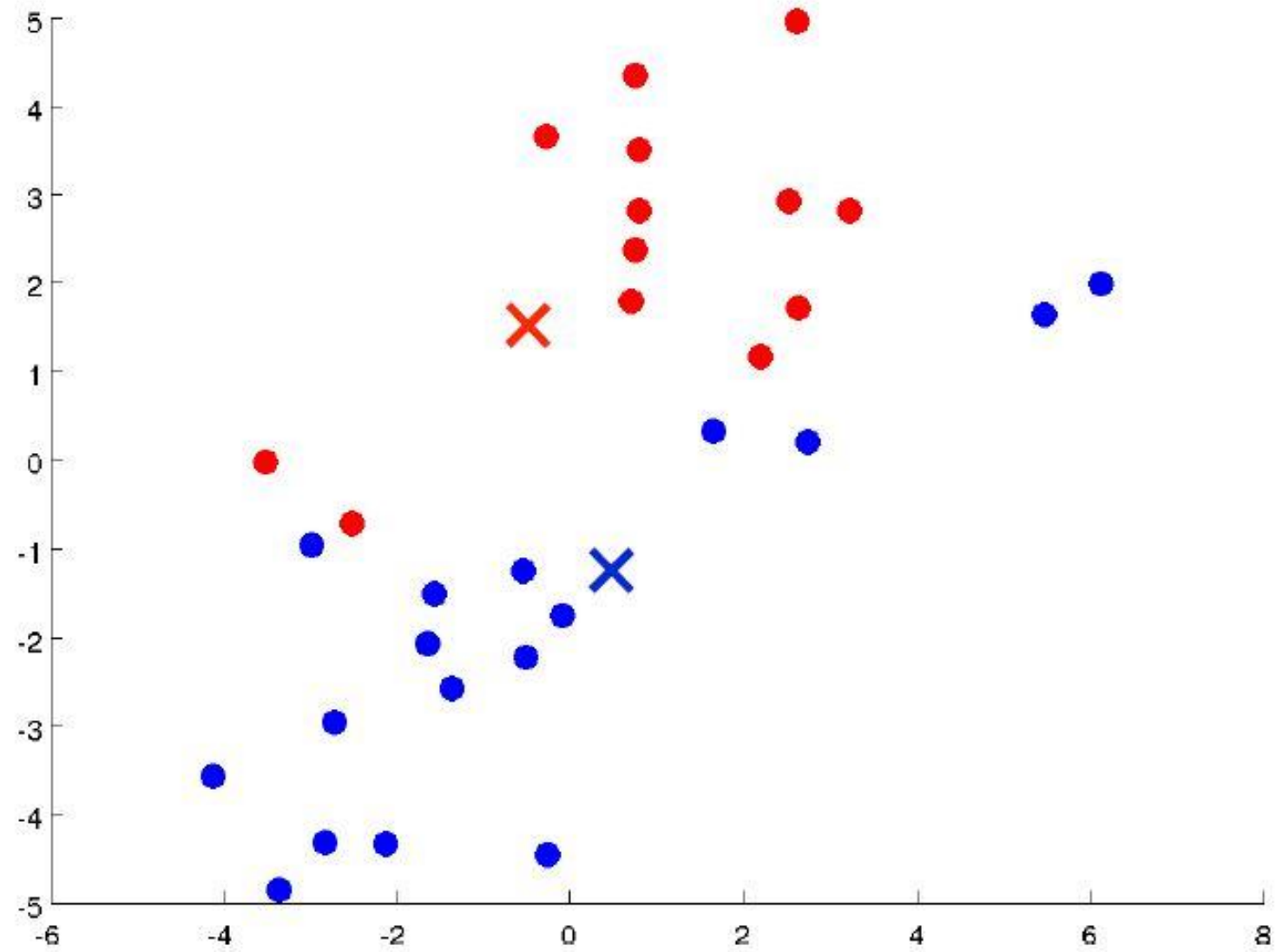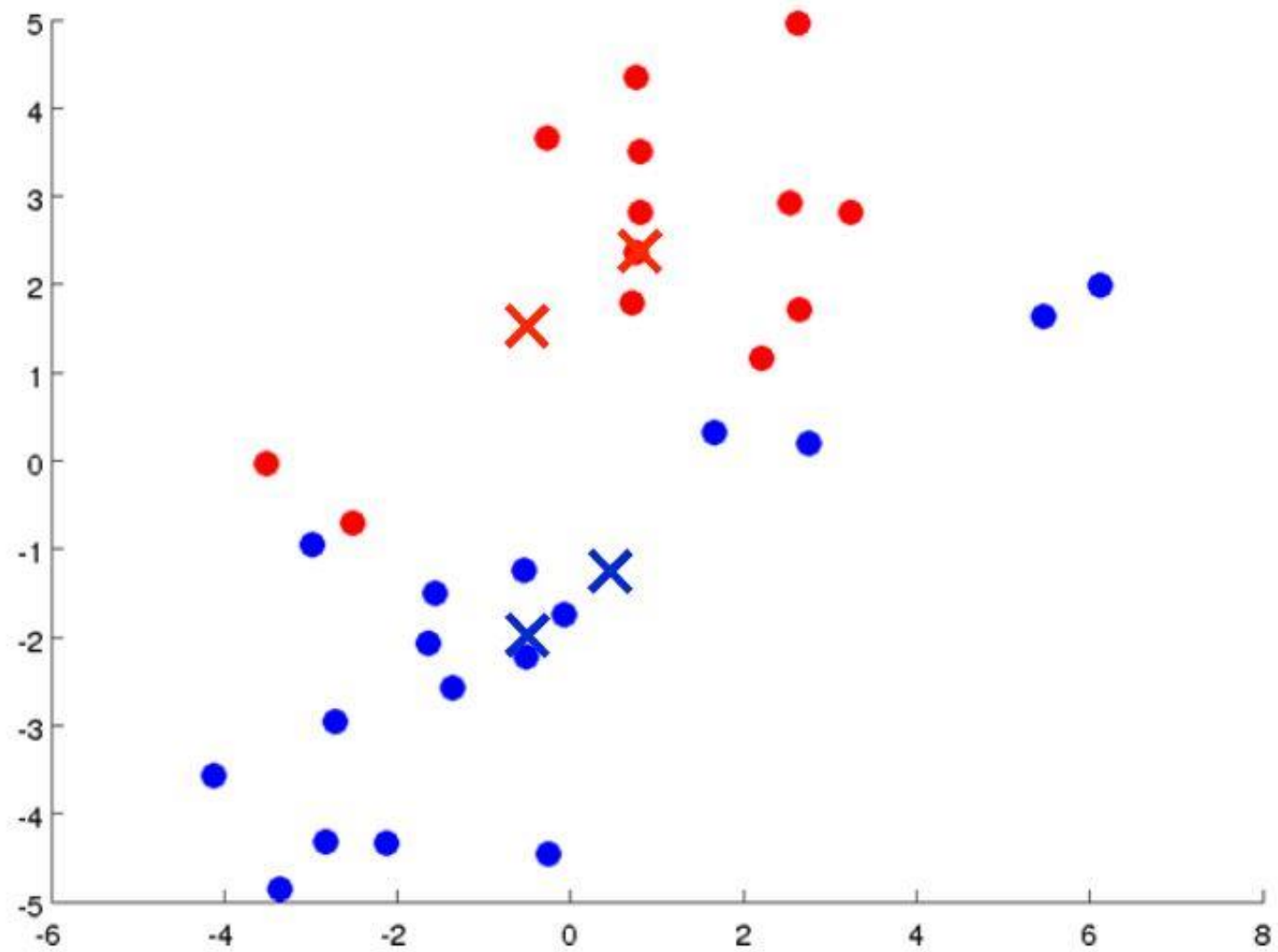
K-means algorithm

cluster centroids

**K-means algorithm**

Input:

- $K$ (number of clusters) ←
- Training set $\{x^{(1)}, x^{(2)}, \ldots, x^{(m)}\}$ ←

**K-means algorithm**

$\mu_1$ ×    $\mu_2$ ×

Randomly initialize $K$ cluster centroids $\mu_1, \mu_2, \ldots, \mu_K \in \mathbb{R}^n$

Repeat {

Cluster assignment step

for $i = 1$ to $m$

$\quad c^{(i)} :=$ index (from 1 to $K$) of cluster centroid

$\quad\quad$ closest to $x^{(i)}$

$\quad\quad\quad \min_k \|x^{(i)} - \mu_k\|^2$

$\quad\quad\quad\quad \hookrightarrow c^{(i)}$

Move centroid

for $k = 1$ to $K$

$\quad \rightarrow \mu_k :=$ average (mean) of points assigned to cluster $k$

$\quad\quad x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)} \rightarrow c^{(1)} = 2, \; c^{(5)} = 2, \; c^{(6)} = 2,$

$\quad\quad\quad\quad c^{(10)} = 2$

$$\mu_2 = \frac{1}{4}\left[ x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)} \right] \in \mathbb{R}^n$$

}

Activate Win
Go to PC settings

# Unsupervised Learning

➢ K-means Algorithm [1]

**Algorithm:** *k*-means. The *k*-means algorithm for partitioning, where each cluster's center is represented by the mean value of the objects in the cluster.

**Input:**

- ■ *k*: the number of clusters,

- ■ *D*: a data set containing *n* objects.

**Output:** A set of *k* clusters.

**Method:**

(1)  arbitrarily choose *k* objects from *D* as the initial cluster centers;
(2)  **repeat**
(3)      (re)assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;
(4)      update the cluster means, i.e., calculate the mean value of the objects for each cluster;
(5)  **until** no change;

**Figure 7.2**  The *k*-means partitioning algorithm.

# Unsupervised Learning

➢ K-means Algorithm [1]

Suppose that the data mining task is to cluster the following eight points (with $(x, y)$ representing location) into three clusters:

$$A_1(2, 10), A_2(2, 5), A_3(8, 4), B_1(5, 8), B_2(7, 5), B_3(6, 4), C_1(1, 2), C_2(4, 9).$$

The distance function is Euclidean distance. Suppose initially we assign $A_1$, $B_1$, and $C_1$ as the center of each cluster, respectively. Use the *k-means* algorithm to show *only*

(a) The three cluster centers after the first round execution

(b) The final three clusters

# K-Means Clustering – Solved Example

**Initial Centroids:**
A1: (2, 10)
B1: (5, 8)
C1: (1, 2)

**New Centroids:**
A1: (2, 10) ╱
B1: (6, 6) ─
C1: (1.5, 3.5) ╱

| Data Points | | | Distance to | | | Cluster |
|---|---|---|---|---|---|---|
| | | | 2  10 | 5  8 | 1  2 | |
| A1 | 2 | 10 | 0.00 | 3.61 | 8.06 | 1 |
| A2 | 2 | 5 | 5.00 | 4.24 | 3.16 | 3 |
| A3 | 8 | 4 | 8.49 | 5.00 | 7.28 | 2 |
| B1 | 5 | 8 | 3.61 | 0.00 | 7.21 | 2 |
| B2 | 7 | 5 | 7.07 | 3.61 | 6.71 | 2 |
| B3 | 6 | 4 | 7.21 | 4.12 | 5.39 | 2 |
| C1 | 1 | 2 | 8.06 | 7.21 | 0.00 | 3 |
| C2 | 4 | 9 | 2.24 | 1.41 | 7.62 | 2 |

$$d(p_1, p_2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

# K-Means Clustering – Solved Example

Current Centroids:

A1: (2, 10)

B1: (6, 6)

C1: (1.5, 3.5)

| Data Points | | | Distance to | | | | | | Cluster | New Cluster |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 2 | 10 | 6 | 6 | 1.5 | 1.5 | | |
| A1 | 2 | 10 | 0.00 | | 5.66 | | 6.52 | | 1 | 1 |
| A2 | 2 | 5 | 5.00 | | 4.12 | | 1.58 | | 3 | 3 |
| A3 | 8 | 4 | 8.49 | | 2.83 | | 6.52 | | 2 | 2 |
| B1 | 5 | 8 | 3.61 | | 2.24 | | 5.70 | | 2 | 2 |
| B2 | 7 | 5 | 7.07 | | 1.41 | | 5.70 | | 2 | 2 |
| B3 | 6 | 4 | 7.21 | | 2.00 | | 4.53 | | 2 | 2 |
| C1 | 1 | 2 | 8.06 | | 6.40 | | 1.58 | | 3 | 3 |
| C2 | 4 | 9 | 2.24 | | 3.61 | | 6.04 | | 2 | 1 |

# K-Means Clustering – Solved Example

| Current Centroids: |
| --- |
| A1: (2, 10) |
| B1: (6, 6) |
| C1: (1.5, 3.5) |

| | Data Points | | Distance to | | | | | | Cluster | New Cluster |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | 2 | 10 | 6 | 6 | 1.5 | 1.5 | | |
| A1 | 2 | 10 | 0.00 | | 5.66 | | 6.52 | | 1 | 1 |
| A2 | 2 | 5 | 5.00 | | 4.12 | | 1.58 | | 3 | 3 |
| A3 | 8 | 4 | 8.49 | | 2.83 | | 6.52 | | 2 | 2 |
| B1 | 5 | 8 | 3.61 | | 2.24 | | 5.70 | | 2 | 2 |
| B2 | 7 | 5 | 7.07 | | 1.41 | | 5.70 | | 2 | 2 |
| B3 | 6 | 4 | 7.21 | | 2.00 | | 4.53 | | 2 | 2 |
| C1 | 1 | 2 | 8.06 | | 6.40 | | 1.58 | | 3 | 3 |
| C2 | 4 | 9 | 2.24 | | 3.61 | | 6.04 | | 2 | 1 |

New Centroids:

A1: (3, 9.5) ✓

B1: (6.5, 5.25) ✓

C1: (1.5, 3.5) ✓

# K-Means Clustering – Solved Example

Current Centroids:
A1: (3.67, 9)
B1: (7, 4.33)
C1: (1.5, 3.5)

| Data Points | | | Distance to | | | | | | Cluster | New Cluster |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 3.67 | 9 | 7 | 4.33 | 1.5 | 3.5 | | |
| A1 | 2 | 10 | 1.94 | | 7.56 | | 6.52 | | 1 | 1 |
| A2 | 2 | 5 | 4.33 | | 5.04 | | 1.58 | | 3 | 3 |
| A3 | 8 | 4 | 6.62 | | 1.05 | | 6.52 | | 2 | 2 |
| B1 | 5 | 8 | 1.67 | | 4.18 | | 5.70 | | 1 | 1 |
| B2 | 7 | 5 | 5.21 | | 0.67 | | 5.70 | | 2 | 2 |
| B3 | 6 | 4 | 5.52 | | 1.05 | | 4.53 | | 2 | 2 |
| C1 | 1 | 2 | 7.49 | | 6.44 | | 1.58 | | 3 | 3 |
| C2 | 4 | 9 | 0.33 | | 5.55 | | 6.04 | | 1 | 1 |

# K-Means Clustering – Solved Example

Current Centroids:

A1: (3.67, 9)
B1: (7, 4.33)
C1: (1.5, 3.5)

| Data Points | | | Distance to | | | | | | Cluster | New Cluster |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 3.67 | 9 | 7 | 4.33 | 1.5 | 3.5 | | |
| A1 | 2 | 10 | 1.94 | | 7.56 | | 6.52 | | 1 | 1 |
| A2 | 2 | 5 | 4.33 | | 5.04 | | 1.58 | | 3 | 3 |
| A3 | 8 | 4 | 6.62 | | 1.05 | | 6.52 | | 2 | 2 |
| B1 | 5 | 8 | 1.67 | | 4.18 | | 5.70 | | 1 | 1 |
| B2 | 7 | 5 | 5.21 | | 0.67 | | 5.70 | | 2 | 2 |
| B3 | 6 | 4 | 5.52 | | 1.05 | | 4.53 | | 2 | 2 |
| C1 | 1 | 2 | 7.49 | | 6.44 | | 1.58 | | 3 | 3 |
| C2 | 4 | 9 | 0.33 | | 5.55 | | 6.04 | | 1 | 1 |