

Module 5 : Unsupervised Learning

Supervised learning vs. unsupervised learning

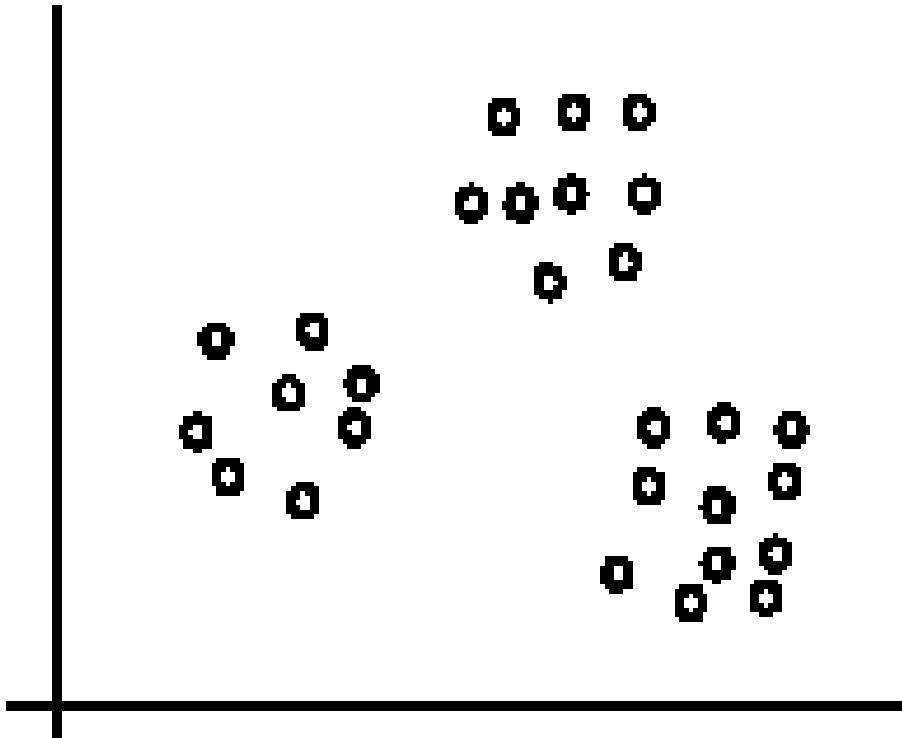
- ❏ **Supervised learning:** discover patterns in the data that relate data attributes with a target (class) attribute.
 - ❏ These patterns are then utilized to predict the values of the target attribute in future data instances.
- ❏ **Unsupervised learning:** The data have no target attribute.
 - ❏ We want to explore the data to find some intrinsic structures in them.

Clustering

- ❏ Clustering is a technique for finding **similarity groups** in data, called **clusters**. i.e.,
 - ❏ it groups data instances that are similar to (near) each other in one cluster and data instances that are very different (far away) from each other into different clusters.
- ❏ Clustering is often called an **unsupervised learning** task as no class values denoting an *a priori* grouping of the data instances are given, which is the case in supervised learning.
- ❏ Due to historical reasons, clustering is often considered synonymous with unsupervised learning.
 - ❏ In fact, association rule mining is also unsupervised

An illustration

- ❏ The data set has three natural groups of data points, i.e., 3 natural clusters.



What is clustering for?

- ❏ Let us see some real-life examples
- ❏ **Example 1:** groups people of similar sizes together to make “small”, “medium” and “large” T-Shirts.
 - ❏ Tailor-made for each person: too expensive
 - ❏ One-size-fits-all: does not fit all.
- ❏ **Example 2:** In marketing, segment customers according to their similarities
 - ❏ To do targeted marketing.

What is clustering for? (cont...)

- ❖ **Example 3:** Given a collection of text documents, we want to organize them according to their content similarities,
 - ❖ To produce a topic hierarchy
- ❖ **In fact, clustering is one of the most utilized data mining techniques.**
 - ❖ It has a long history, and used in almost every field, e.g., medicine, psychology, botany, sociology, biology, archeology, marketing, insurance, libraries, etc.
 - ❖ In recent years, due to the rapid increase of online documents, text clustering becomes important.

Types of Clustering



Hierarchical algorithms: these find successive clusters using previously established clusters.

1. Agglomerative ("bottom-up"): Agglomerative algorithms begin with each element as a separate cluster and merge them into successively larger clusters.

2. Divisive ("top-down"): Divisive algorithms begin with the whole set and proceed to divide it into successively smaller clusters.



Partitional clustering: Partitional algorithms determine all clusters at once. They include:



K-means and derivatives



Fuzzy c-means clustering



QT clustering algorithm

Aspects of clustering

- ❏ A distance (similarity, or dissimilarity) function
- ❏ Clustering quality
 - ❏ Inter-clusters distance \Rightarrow maximized
 - ❏ Intra-clusters distance \Rightarrow minimized
- ❏ The **quality** of a clustering result depends on the algorithm, the distance function, and the application.

K-means clustering

- ❏ K-means is a partitional clustering algorithm

- ❏ Let the set of data points (or instances) D be

$$\{x_1, x_2, \dots, x_n\},$$

where $x_i = (x_{i1}, x_{i2}, \dots, x_{ir})$ is a vector in a real-valued space $X \subseteq \mathbb{R}^r$, and r is the number of attributes (dimensions) in the data.

- ❏ The k-means algorithm partitions the given data into k clusters.

- ❏ Each cluster has a cluster center, called centroid.
- ❏ k is specified by the user

K-means algorithm

- ❖ Given k , the *k-means* algorithm works as follows:
 1. Randomly choose k data points (**seeds**) to be the initial **centroids**, cluster centers
 2. Assign each data point to the closest **centroid**
 3. Re-compute the **centroids** using the current cluster memberships.
 4. If a convergence criterion is not met, go to 2).

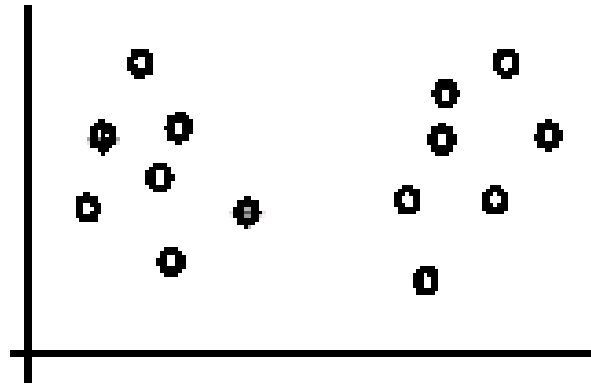
Stopping/convergence criterion

1. no (or minimum) re-assignments of data points to different clusters,
2. no (or minimum) change of centroids, or
3. minimum decrease in the **sum of squared error** (SSE),

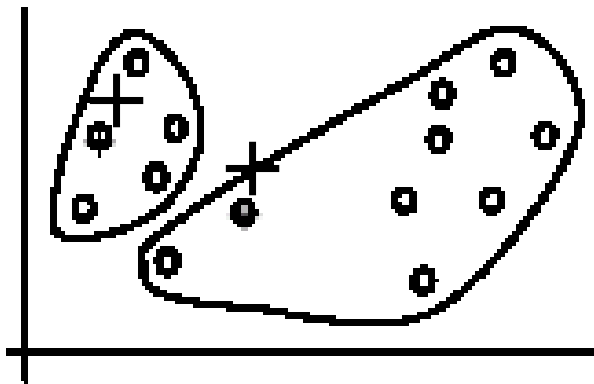
$$SSE = \sum_{j=1}^k \sum dist(x, m_j)^2 \quad (1)$$

- ❏ C_i is the j th cluster, \mathbf{m}_j is the centroid of cluster C_j (the mean vector of all the data points in C_j), and $dist(\mathbf{x}, \mathbf{m}_j)$ is the distance between data point \mathbf{x} and centroid \mathbf{m}_j .

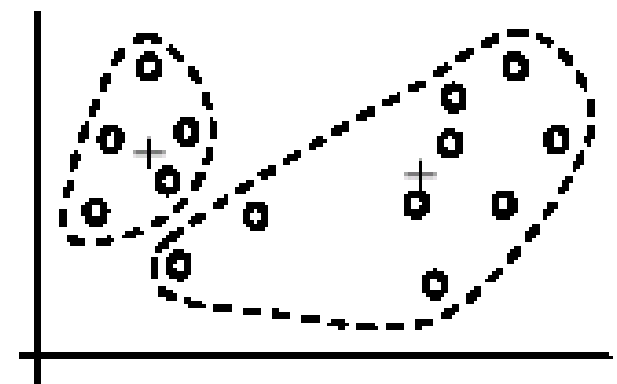
An example



(A). Random selection of k centers

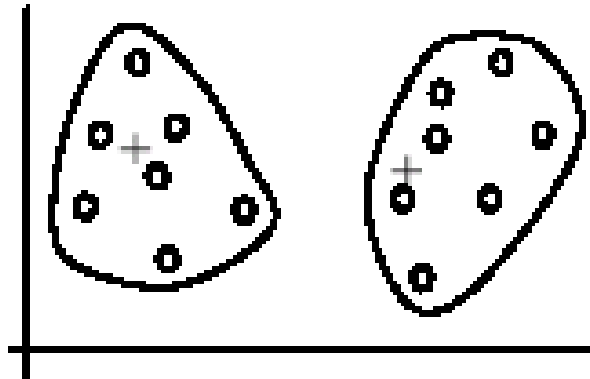


Iteration 1: (B). Cluster assignment

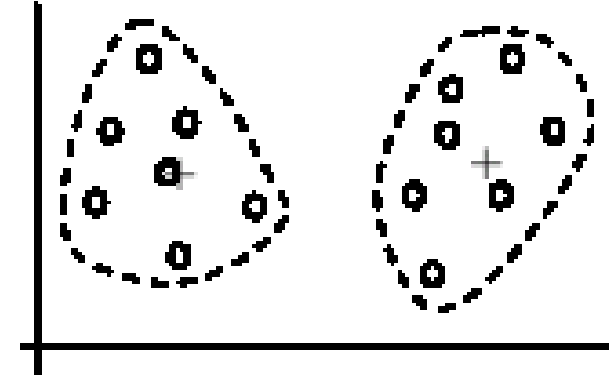


(C). Re-compute centroids

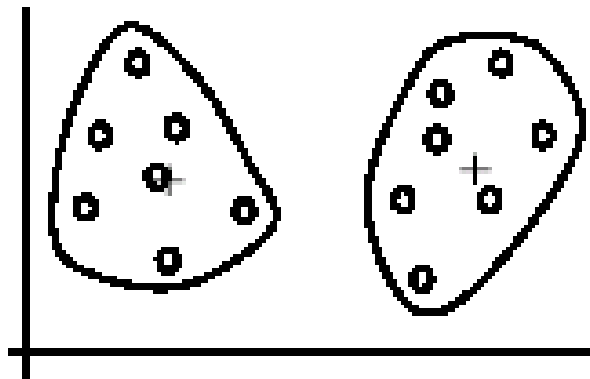
An example (cont ...)



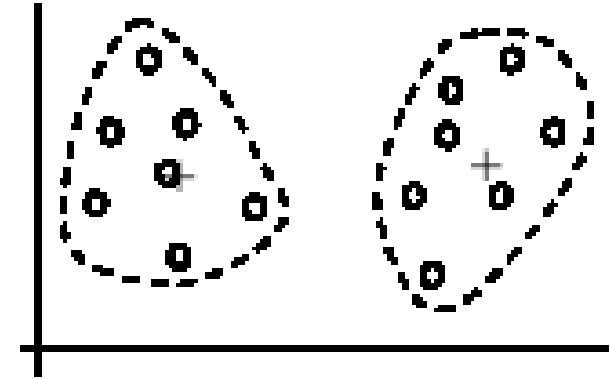
Iteration 2: (D). Cluster assignment



(E). Re-compute centroids



Iteration 3: (F). Cluster assignment



(G). Re-compute centroids

Numerical Example : K- Means Clustering ($K = 2$)

Individual	Variable 1	Variable 2
1	1.0	1.0
2	1.5	2.0
3	3.0	4.0
4	5.0	7.0
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5

Step 1:

Initialization: Randomly we choose following two centroids ($k=2$) for two clusters.

In this case the 2 centroid are: $m1=(1.0,1.0)$ and $m2=(5.0,7.0)$.

Individual	Variable 1	Variable 2
1	1.0	1.0
2	1.5	2.0
3	3.0	4.0
4	5.0	7.0
5	3.5	5.0
6	4.5	5.0
7	3.5	4.5

	Individual	Mean Vector
Group 1	1	(1.0, 1.0)
Group 2	4	(5.0, 7.0)

Step 2:

Thus, we obtain two clusters containing:

$\{1,2,3\}$ and $\{4,5,6,7\}$.

Their new centroids are:

$$m_1 = \left(\frac{1}{3}(1.0 + 1.5 + 3.0), \frac{1}{3}(1.0 + 2.0 + 4.0) \right) = (1.83, 2.33)$$

$$m_2 = \left(\frac{1}{4}(5.0 + 3.5 + 4.5 + 3.5), \frac{1}{4}(7.0 + 5.0 + 5.0 + 4.5) \right) \\ = (4.12, 5.38)$$

Individual	Centroid 1	Centroid 2
1	0	7.21
2 (1.5, 2.0)	1.12	6.10
3	3.61	3.61
4	7.21	0
5	4.72	2.5
6	5.31	2.06
7	4.30	2.92

$$d(m_1, 2) = \sqrt{|1.0 - 1.5|^2 + |1.0 - 2.0|^2} = 1.12$$

$$d(m_2, 2) = \sqrt{|5.0 - 1.5|^2 + |7.0 - 2.0|^2} = 6.10$$

Step 3:

Now using these centroids we compute the Euclidean distance of each object, as shown in table.

Therefore, the new clusters are:

{1,2} and {**3**,4,5,6,7}

Next centroids are:
 $m_1 = (1.25, 1.5)$ and $m_2 = (3.9, 5.1)$

Individual	Centroid 1	Centroid 2
1	1.57	5.38
2	0.47	4.28
3	2.04	1.78
4	5.64	1.84
5	3.15	0.73
6	3.78	0.54
7	2.74	1.08

Step 4 :

The clusters obtained
are:

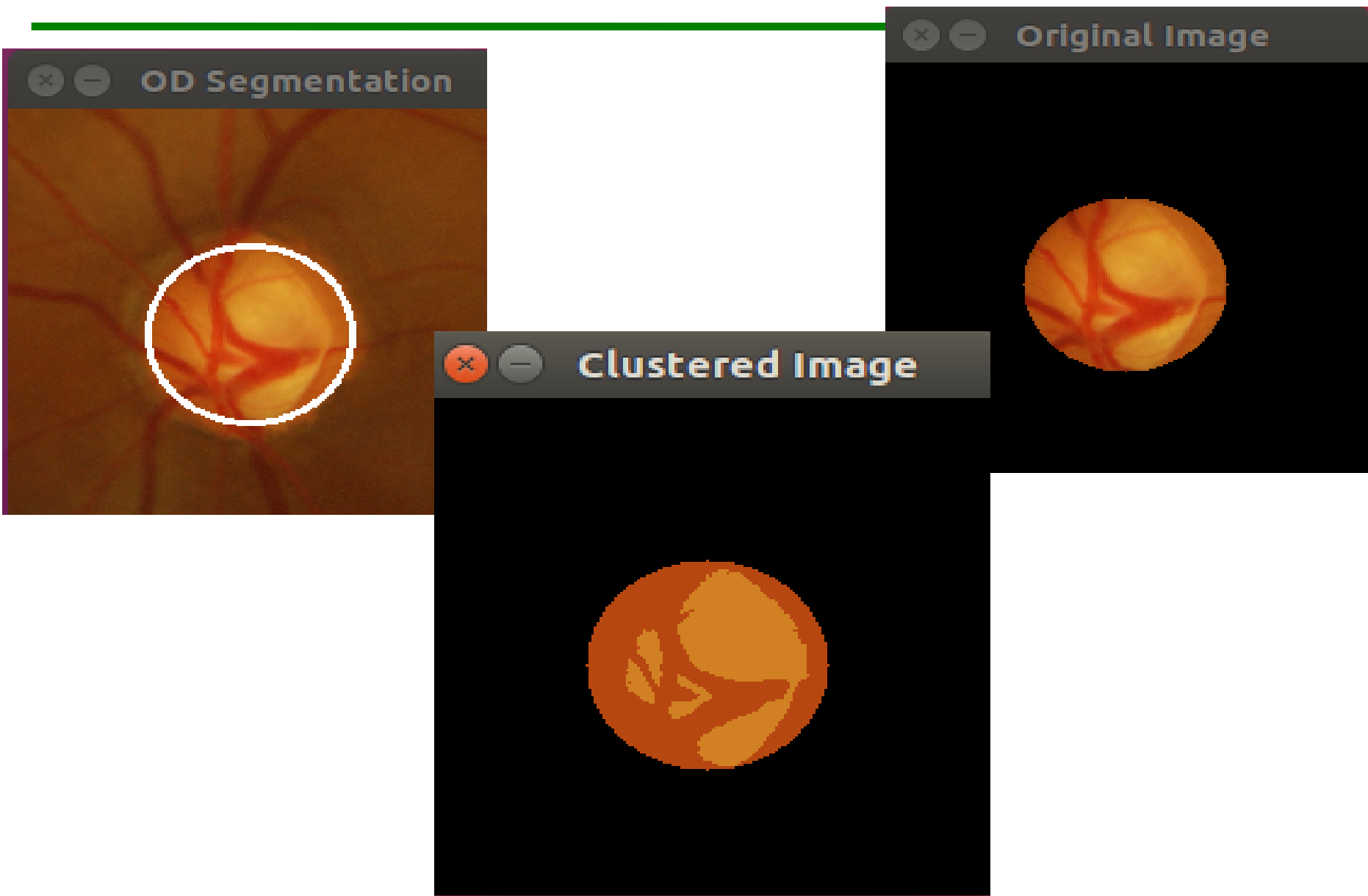
{1,2} and {3,4,5,6,7}

Therefore, there is no
change in the cluster.

Thus, the algorithm
comes to a halt here and
final result consist of 2
clusters {1,2} and
{3,4,5,6,7}.

Individual	Centroid 1	Centroid 2
1	0.58	5.02
2	0.58	3.92
3	3.05	1.42
4	0.88	2.20
5	4.18	0.41
6	4.78	0.81
7	3.75	0.72

Example



Strengths of k -means

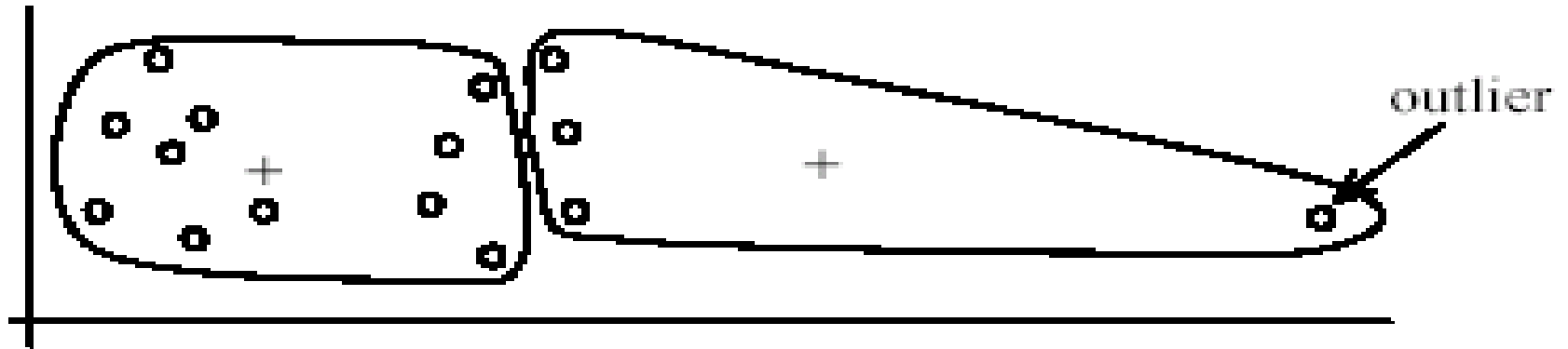
Strengths:

- Simple: easy to understand and to implement
- Efficient: Time complexity: $O(tkn)$,
where n is the number of data points,
 k is the number of clusters, and
 t is the number of iterations.
- Since both k and t are small. k -means is considered a linear algorithm.
- K-means is the most popular clustering algorithm.
- Note that: it terminates at a **local optimum** if SSE is used. The **global optimum** is hard to find due to complexity.

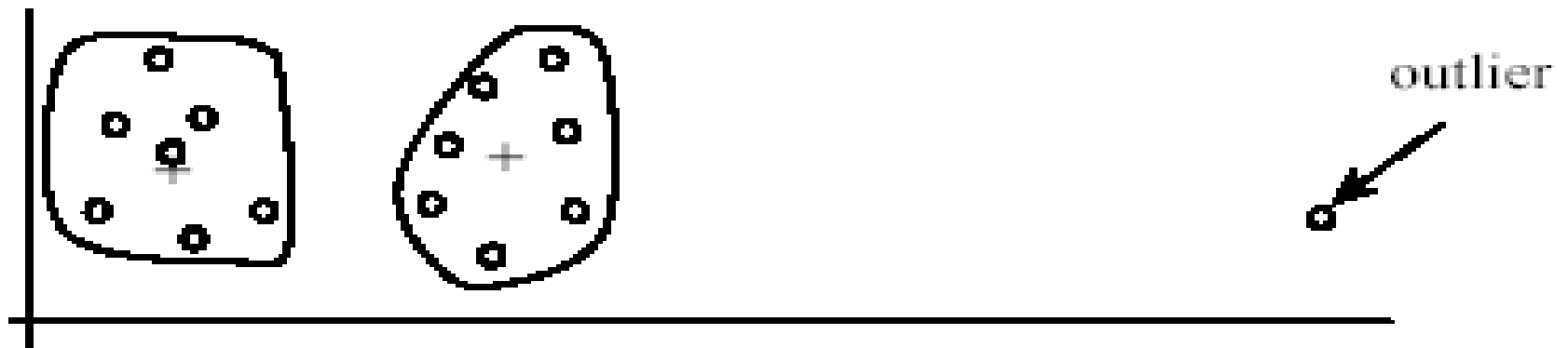
Weaknesses of k -means

- ❏ The algorithm is only applicable if the **mean** is defined.
 - ❏ For categorical data, k -mode - the centroid is represented by most frequent values.
- ❏ The user needs to specify k .
- ❏ The algorithm is sensitive to **outliers**
 - ❏ Outliers are data points that are very far away from other data points.
 - ❏ Outliers could be errors in the data recording or some special data points with very different values.

Weaknesses of k -means: Problems with outliers



(A): Undesirable clusters



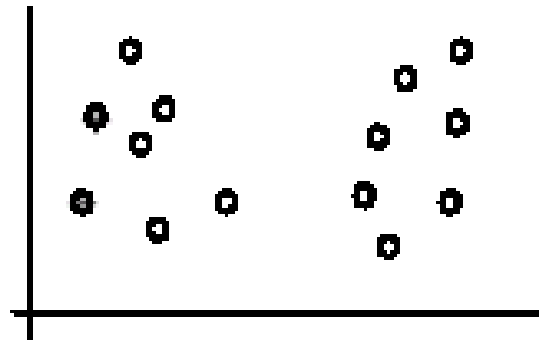
(B): Ideal clusters

Weaknesses of k-means: To deal with outliers

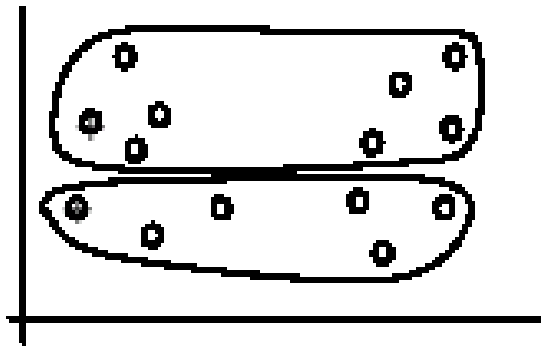
- ❏ One method is to remove some data points in the clustering process that are much further away from the centroids than other data points.
 - ❏ To be safe, we may want to monitor these possible outliers over a few iterations and then decide to remove them.
- ❏ Another method is to perform random sampling. Since in sampling we only choose a small subset of the data points, the chance of selecting an outlier is very small.
 - ❏ Assign the rest of the data points to the clusters by distance or similarity comparison, or classification

Weaknesses of k -means (cont ...)

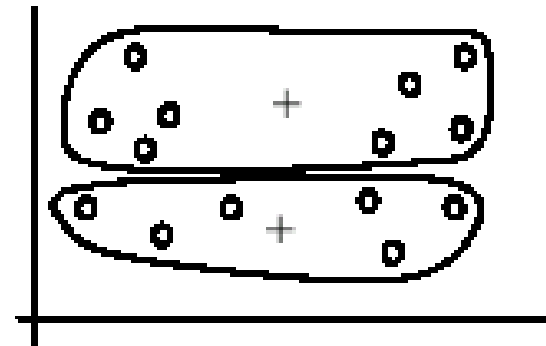
- ❖ The algorithm is sensitive to **initial seeds**.



(A). Random selection of seeds (centroids)



(B). Iteration 1

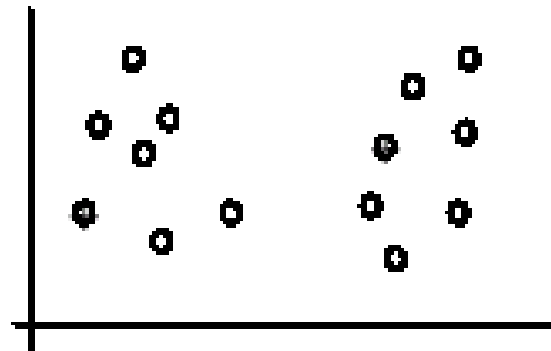


(C). Iteration 2

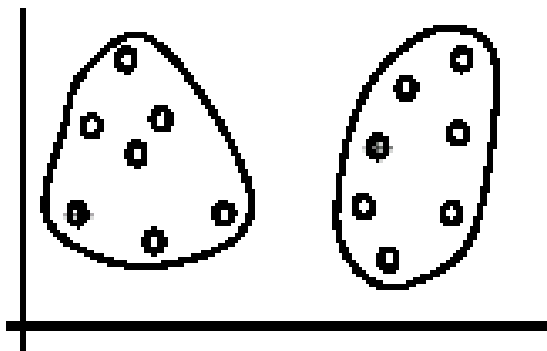
Weaknesses of k -means (cont ...)

- ❖ If we use **different seeds**: good results

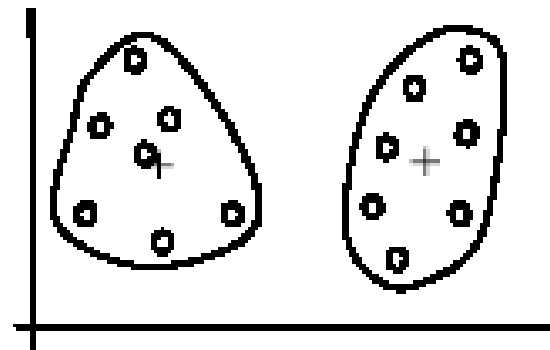
❖ There are some methods to help choose good seeds



(A). Random selection of k seeds (centroids)



(B). Iteration 1



(C). Iteration 2

K-means summary

- ❖ Despite weaknesses, *k*-means is still the most popular algorithm due to its simplicity, efficiency and
 - ❖ other clustering algorithms have their own lists of weaknesses.
- ❖ No clear evidence that any other clustering algorithm performs better in general
 - ❖ although they may be more suitable for some specific types of data or applications.
- ❖ Comparing different clustering algorithms is a difficult task. No one knows the correct clusters!