

# EE2211 Introduction to Machine Learning

## Lecture 11

Wang Xinchao  
[xinchao@nus.edu.sg](mailto:xinchao@nus.edu.sg)

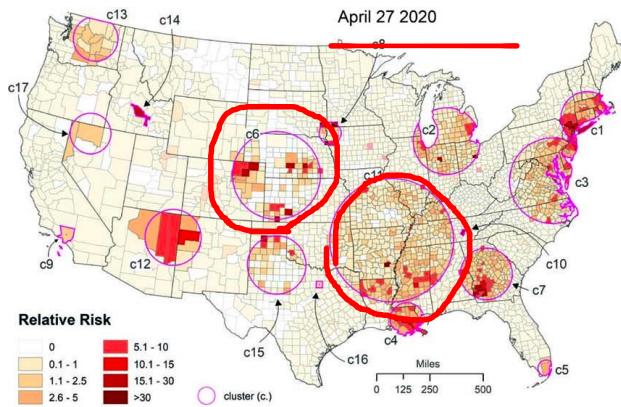
# Course Contents

- Introduction and Preliminaries (Xinchao)
  - Introduction
  - Data Engineering
  - Introduction to Linear Algebra, Probability and Statistics
- Fundamental Machine Learning Algorithms I (Vincent)
  - Systems of linear equations
  - Least squares, Linear regression
  - Ridge regression, Polynomial regression
- Fundamental Machine Learning Algorithms II (Vincent)
  - Over-fitting, bias/variance trade-off
  - Optimization, Gradient descent
  - Decision Trees, Random Forest
- Performance and More Algorithms (Xinchao)
  - Performance Issues
  - K-means Clustering
  - Neural Networks

# Outline

- Introduction of unsupervised learning
- K-means Clustering
  - The most popular clustering technique
- Fuzzy Clustering

# Unsupervised Learning



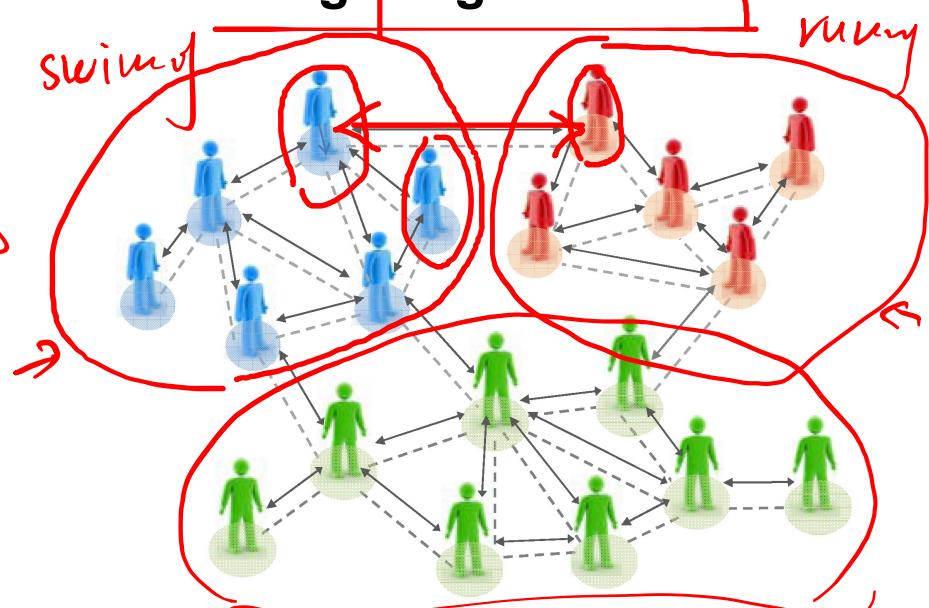
Discovering Covid clusters



Business analysis



Image segmentation



# Unsupervised Learning

## Introduction

**Motivation:** we do not always have labeled data.

In unsupervised learning, the dataset is a collection of unlabeled examples  $\{\mathbf{x}_i\}_{i=1}^M$ .

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$$

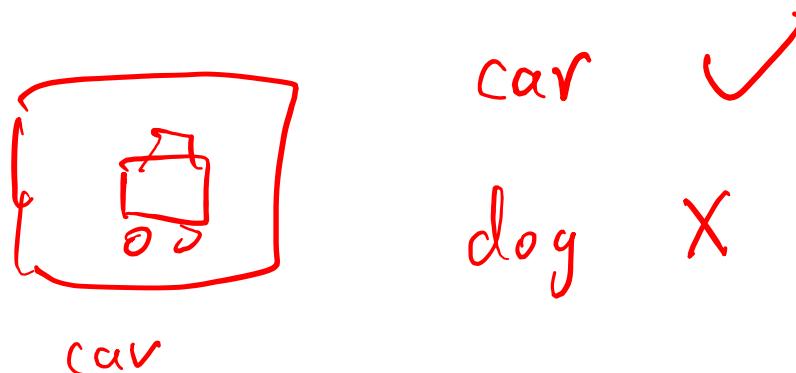


# Unsupervised Learning

## Introduction

Evaluation of unsupervised learning is hard:

- The **absence of labels** representing the desired behavior for your model means **the absence of a solid reference point to judge the quality of your model.**

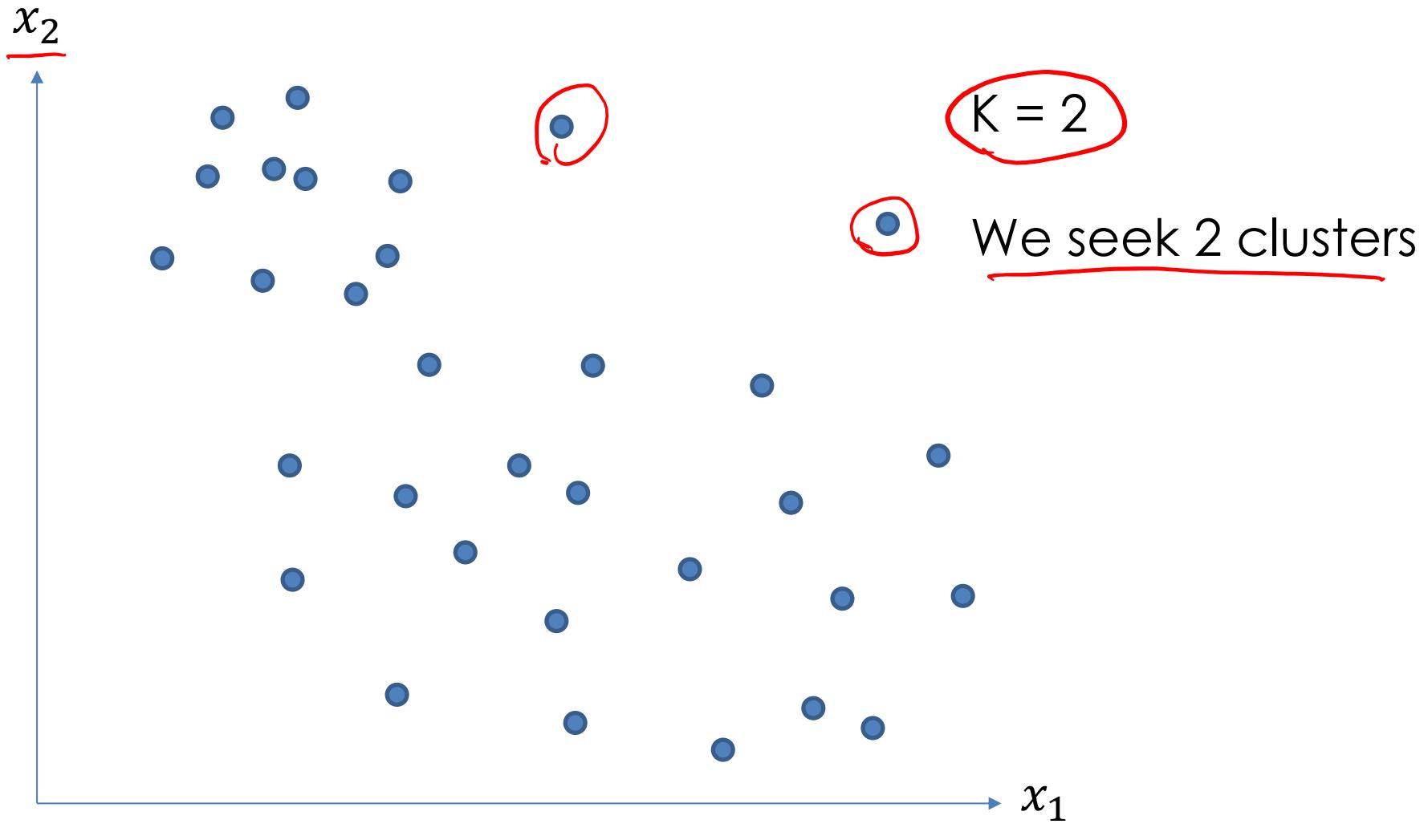


# Unsupervised Learning

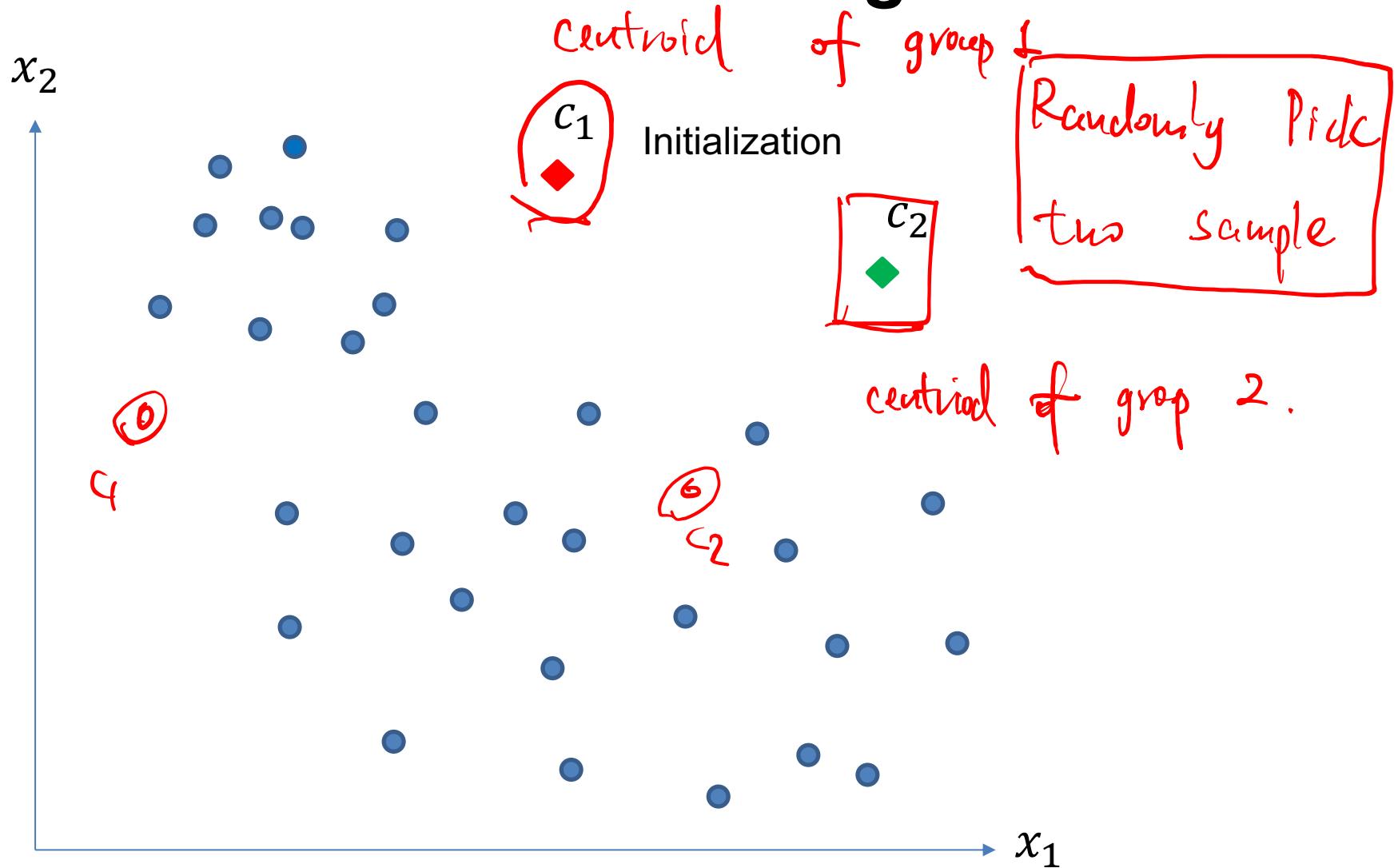
## Main Approaches

- **Clustering**
  - ✓ Groups a set of objects in such a way that objects in the same group (called a **cluster**) are **more similar** (in some sense) to each other than to those in other groups (clusters).
- **Density Estimation**
  - ✓ Models the probability density function (pdf) of the unknown probability distribution from which the dataset has been drawn.
- **Component Analysis**
  - ✓ Breaks down the data from the perspective of signal analysis.
- **Unsupervised Neural Networks**
  - ✓ Autoencoder

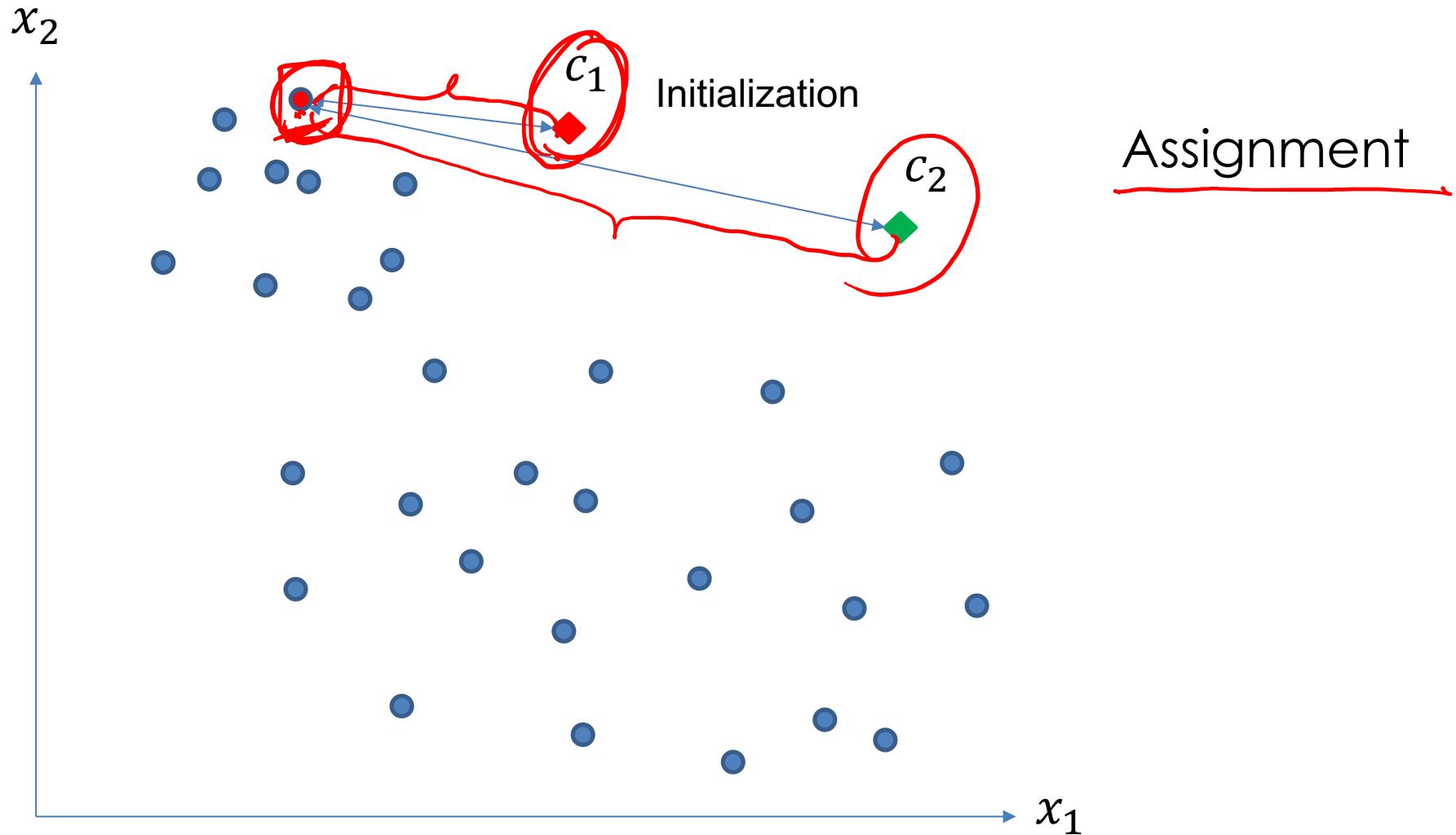
# K-means Clustering (2D)



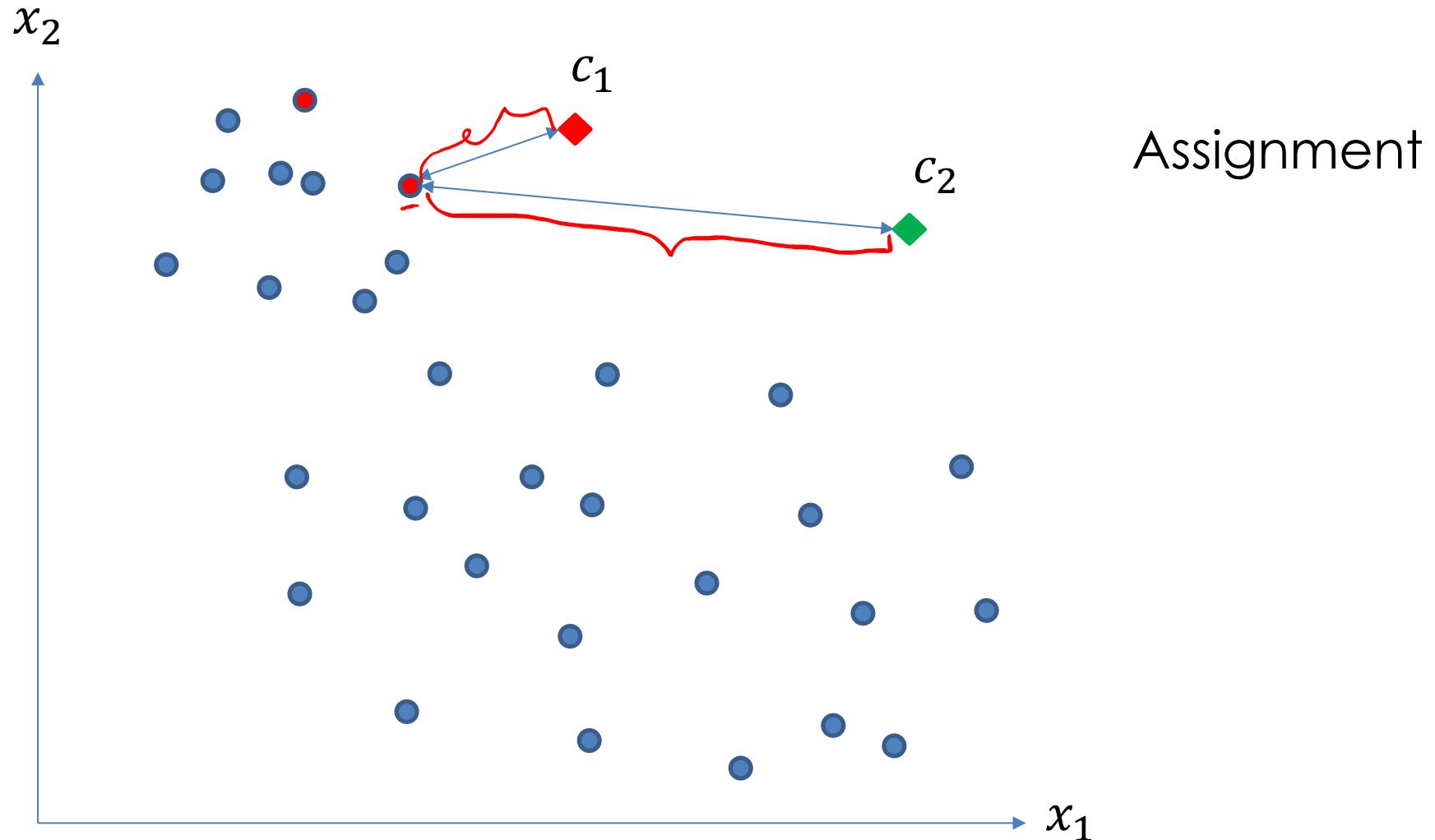
# K-means Clustering



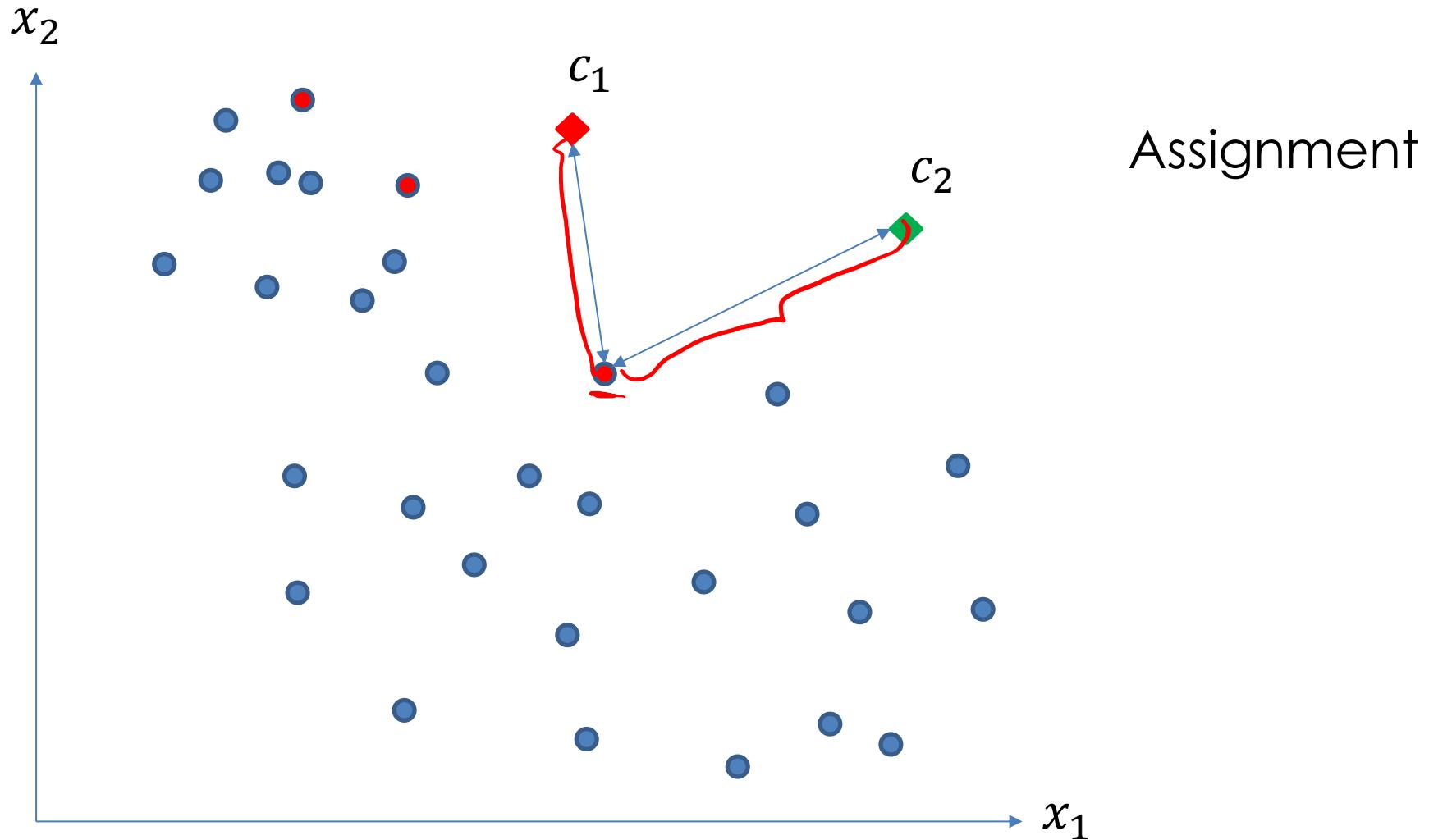
# K-means Clustering



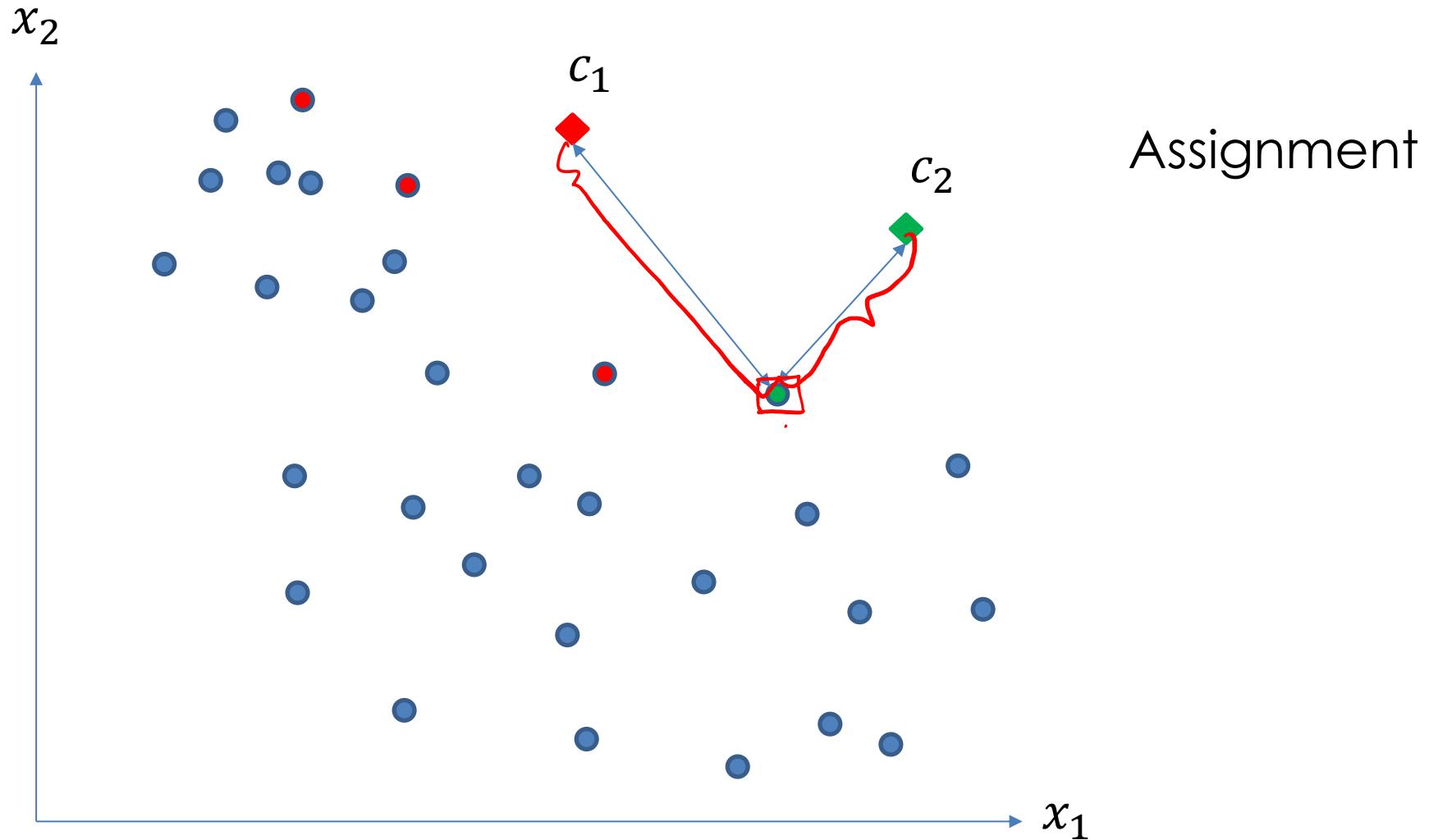
# K-means Clustering



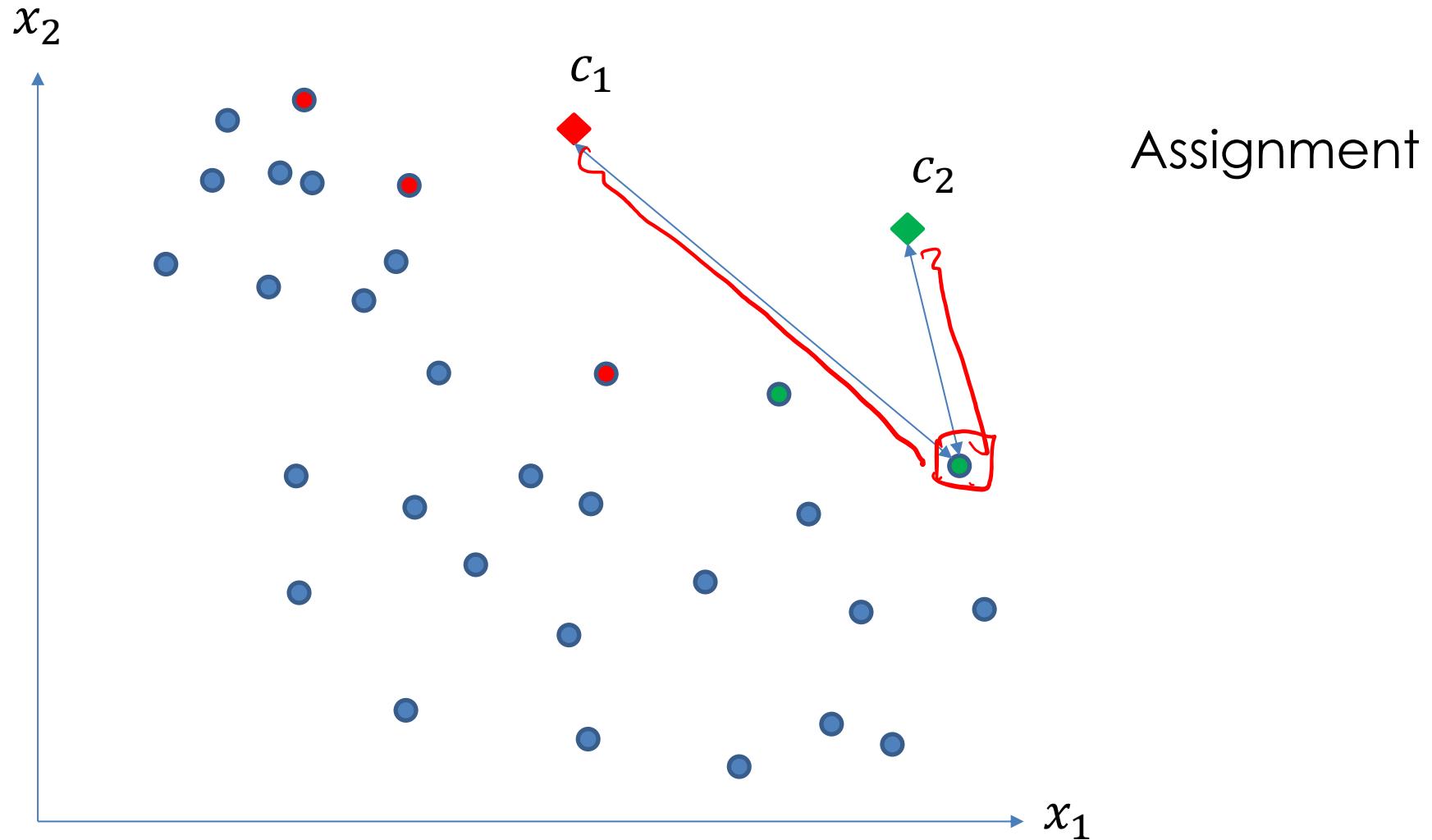
# K-means Clustering



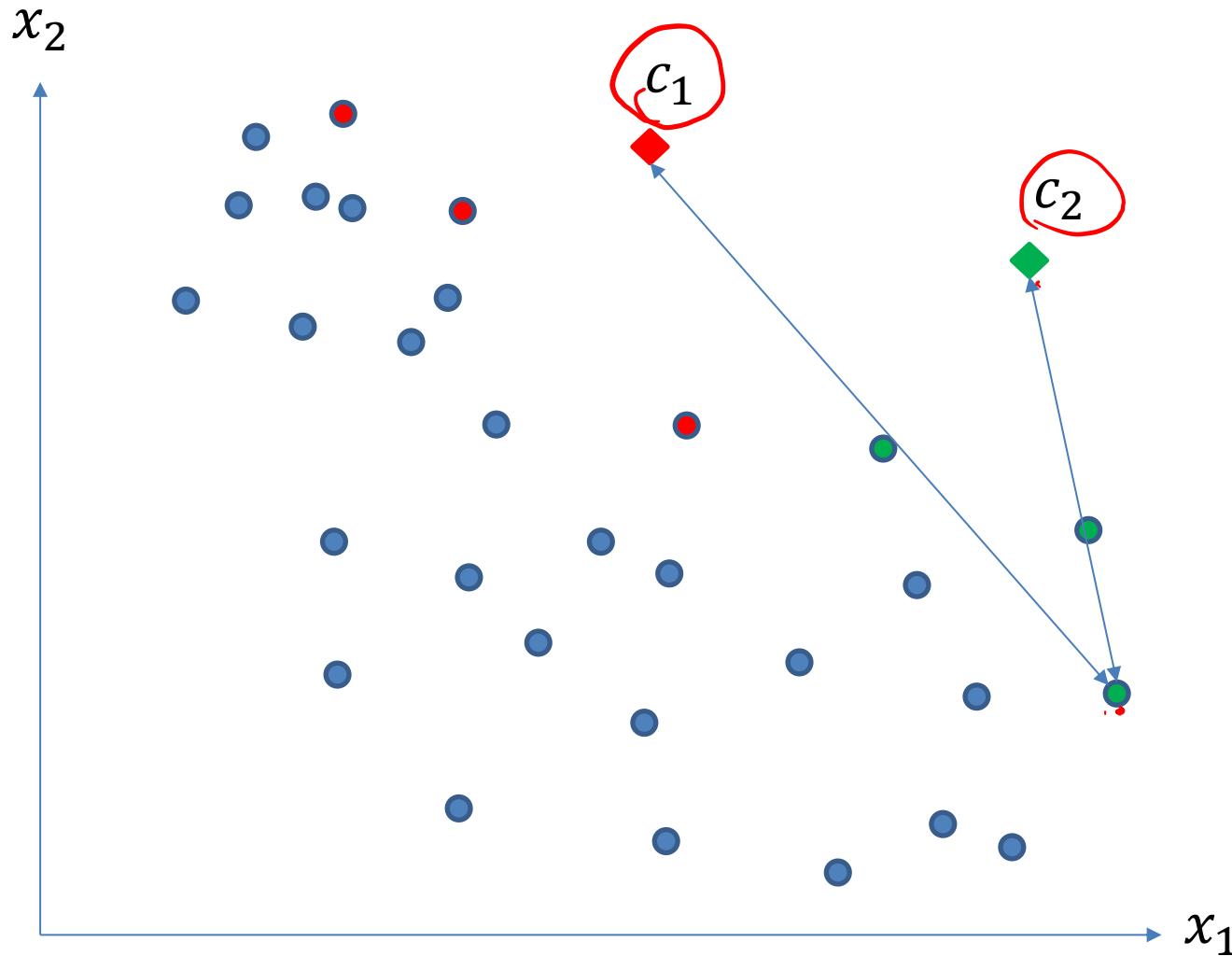
# K-means Clustering



# K-means Clustering



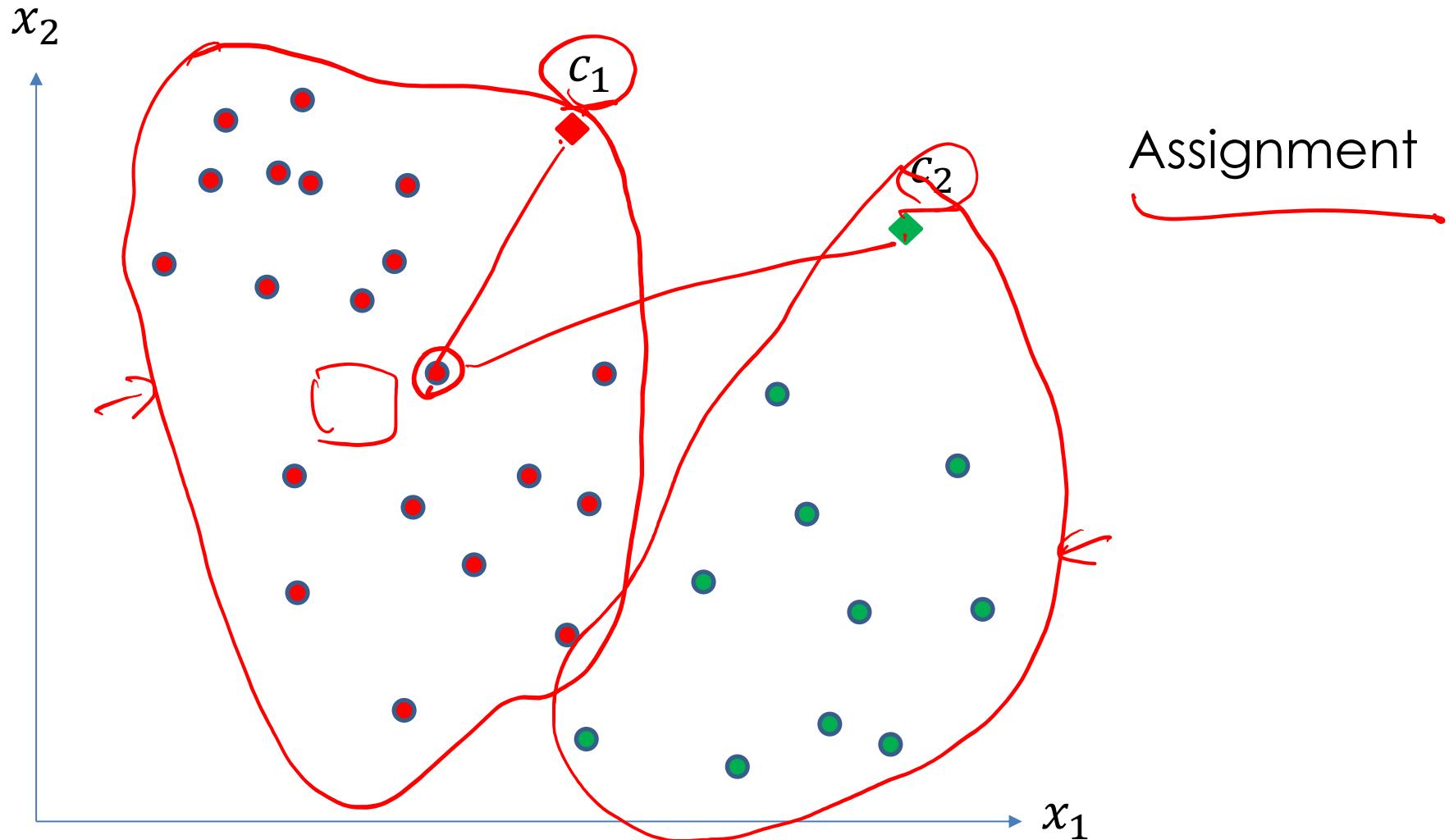
# K-means Clustering



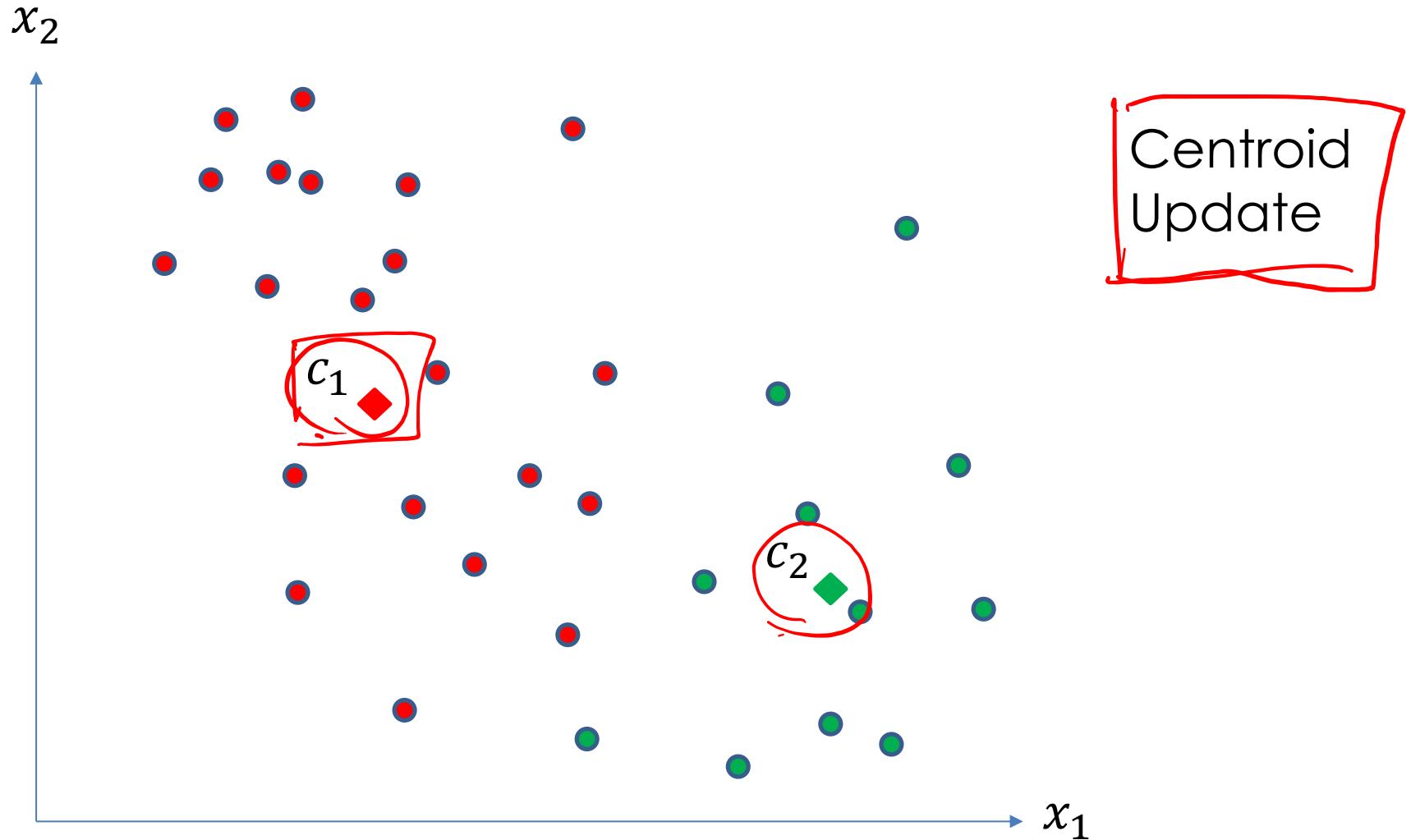
Assignment

Fix  
the  
centroids

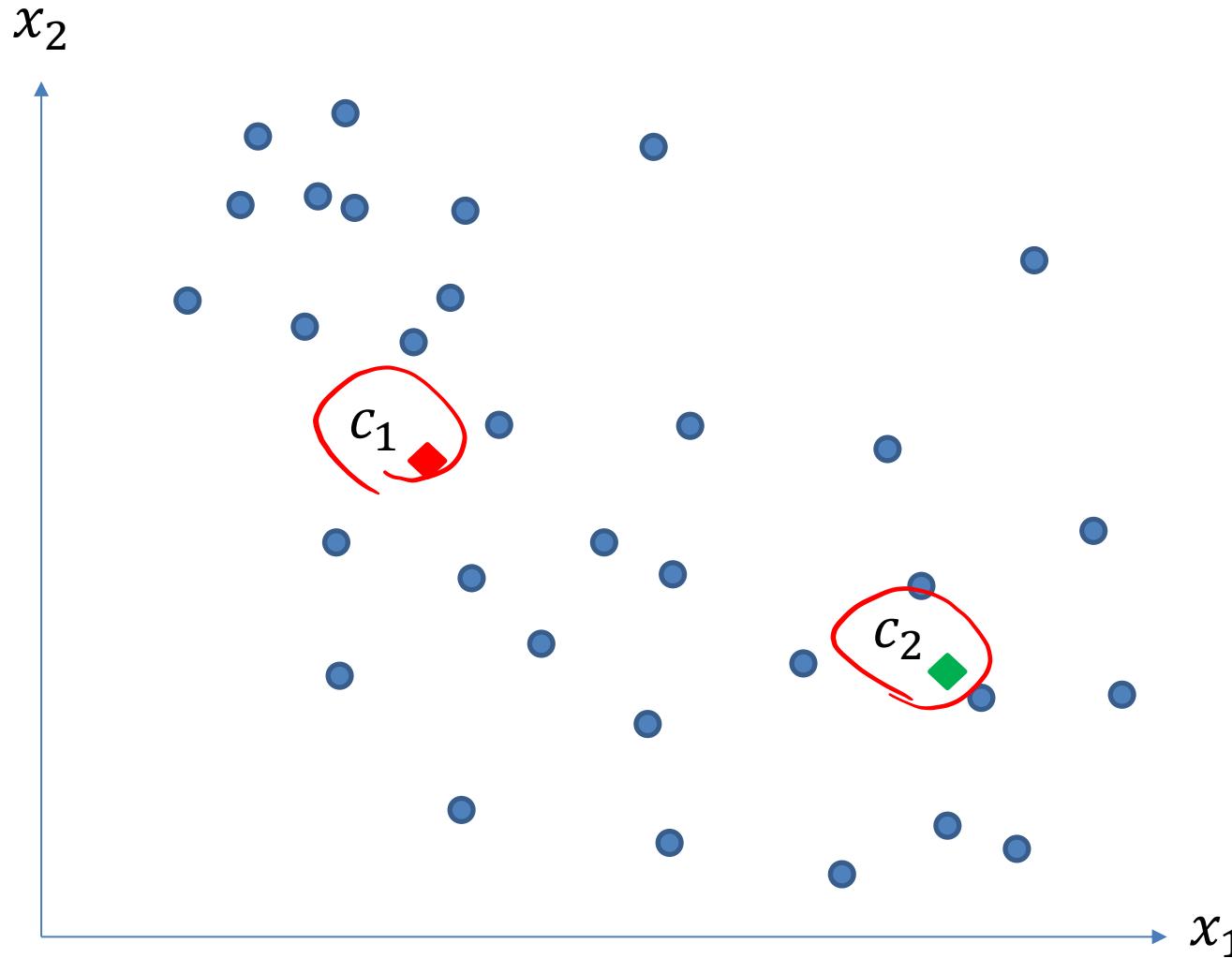
# K-means Clustering



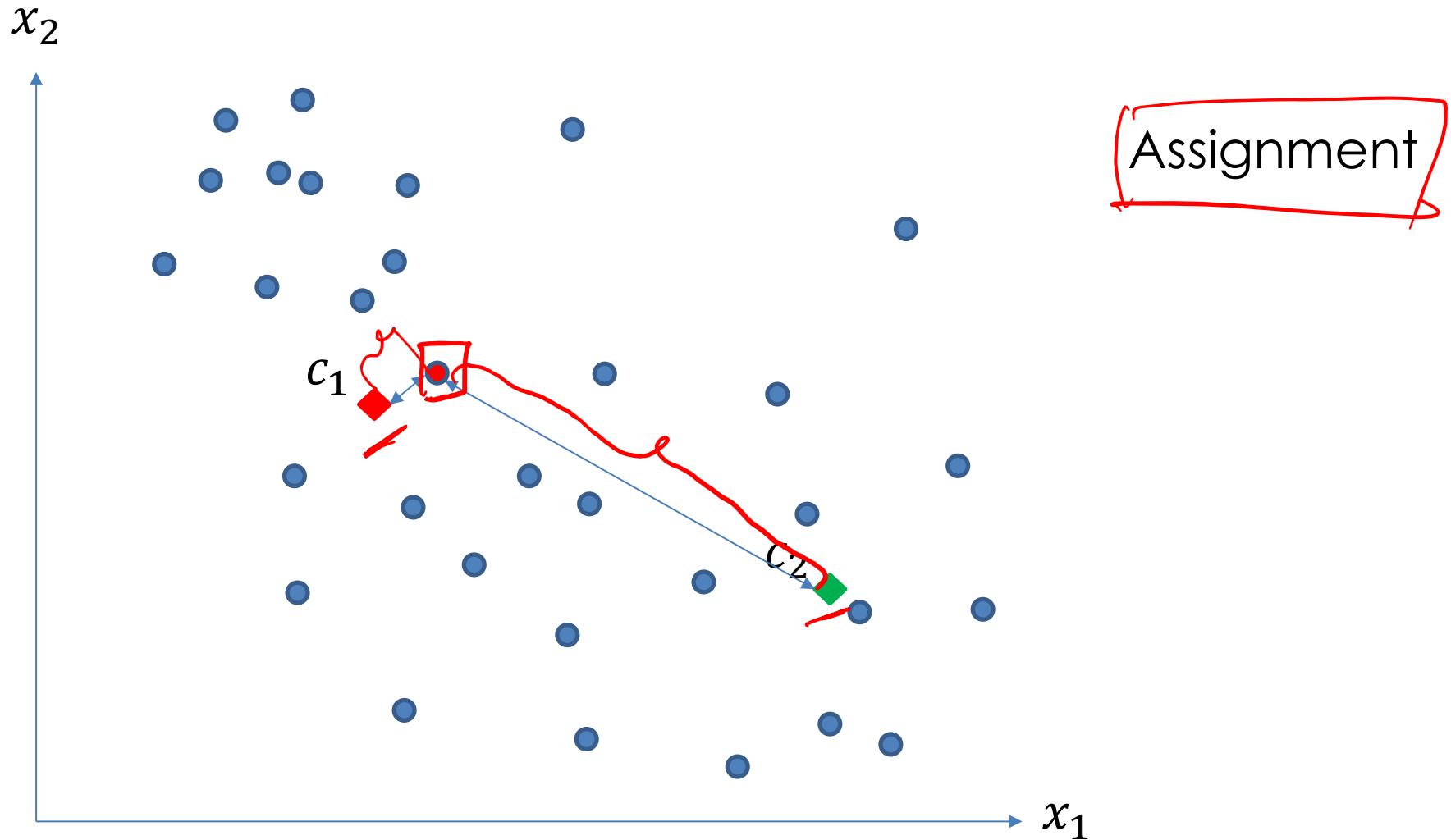
# K-means Clustering



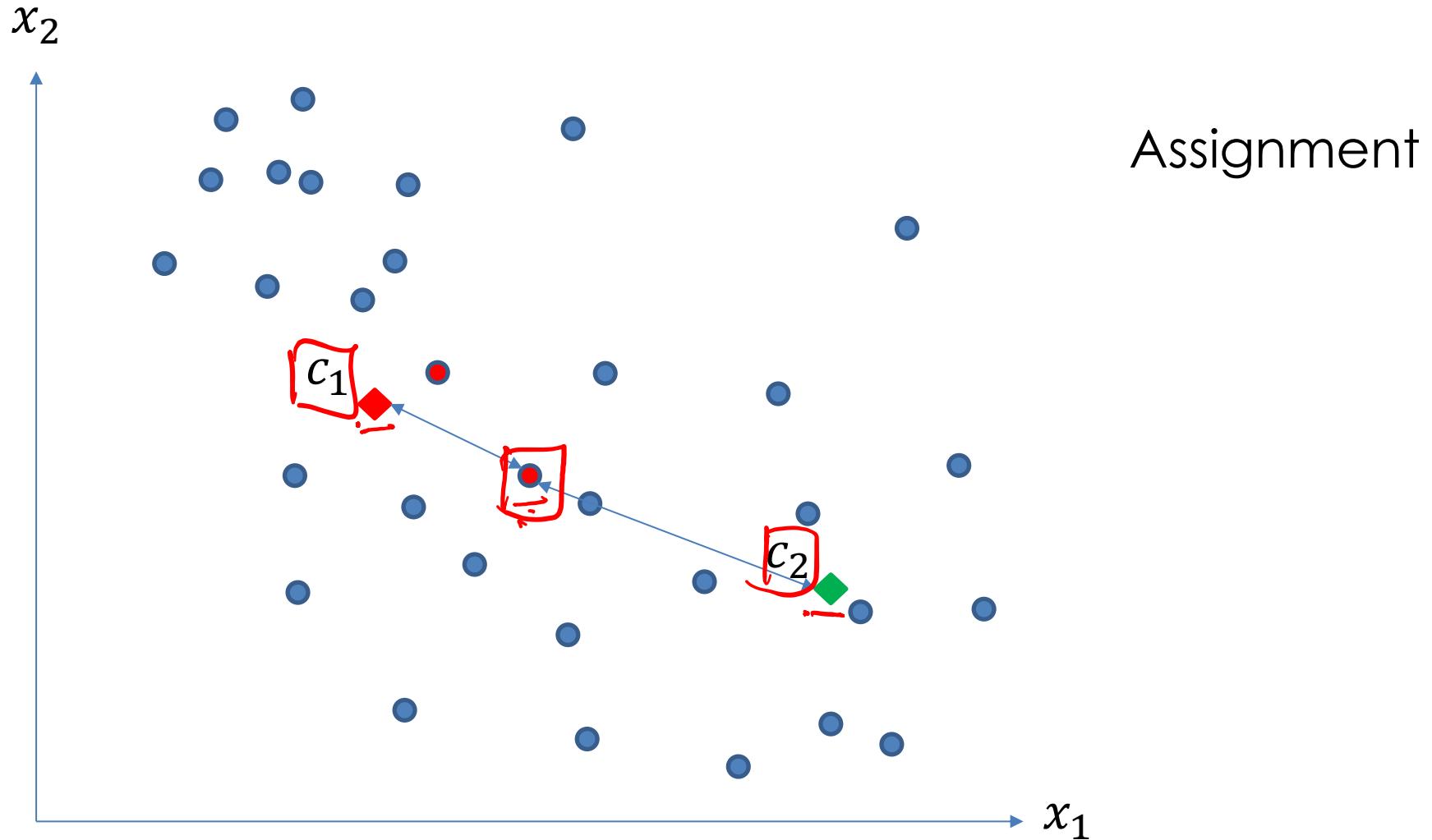
# K-means Clustering



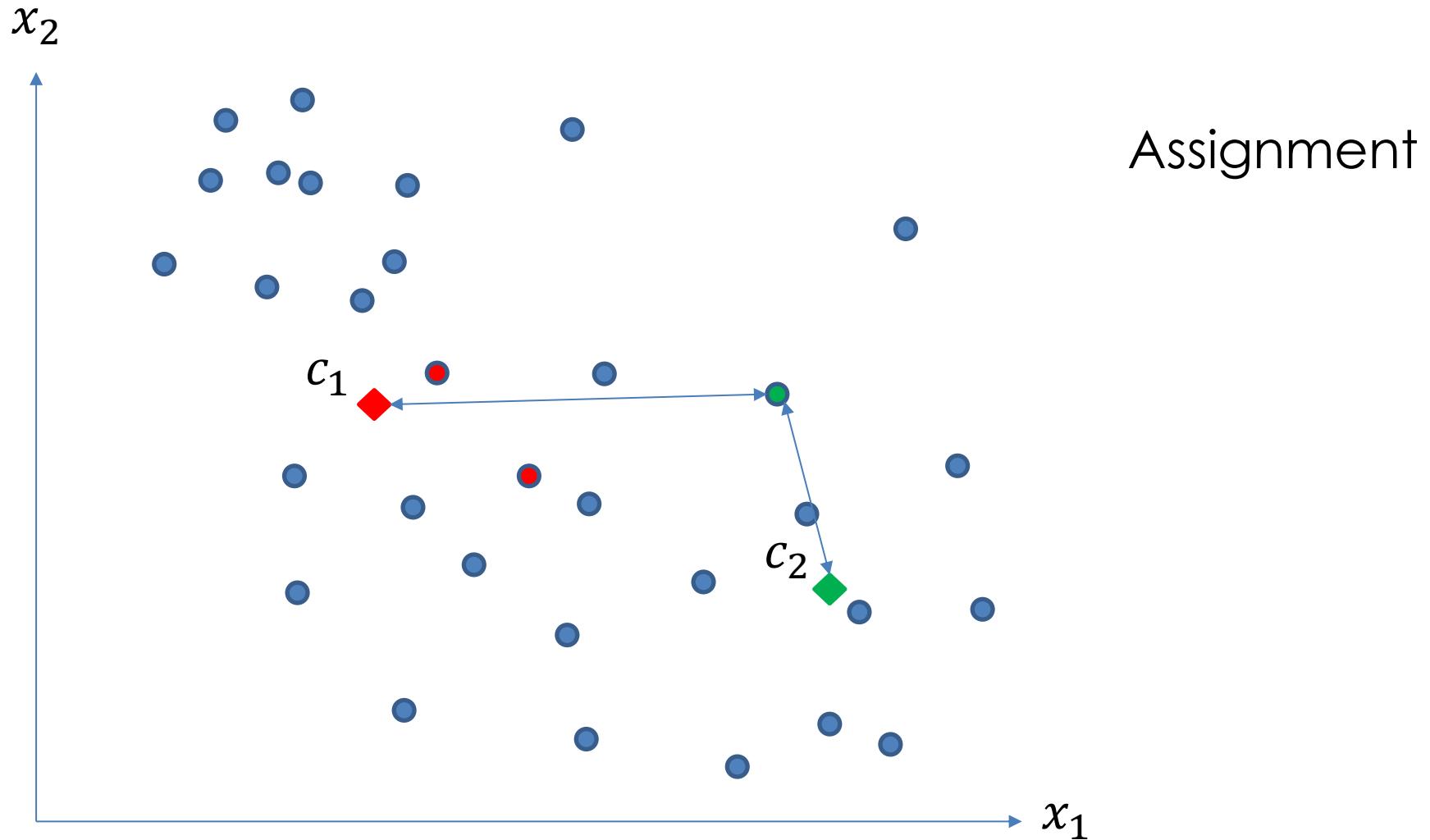
# K-means Clustering



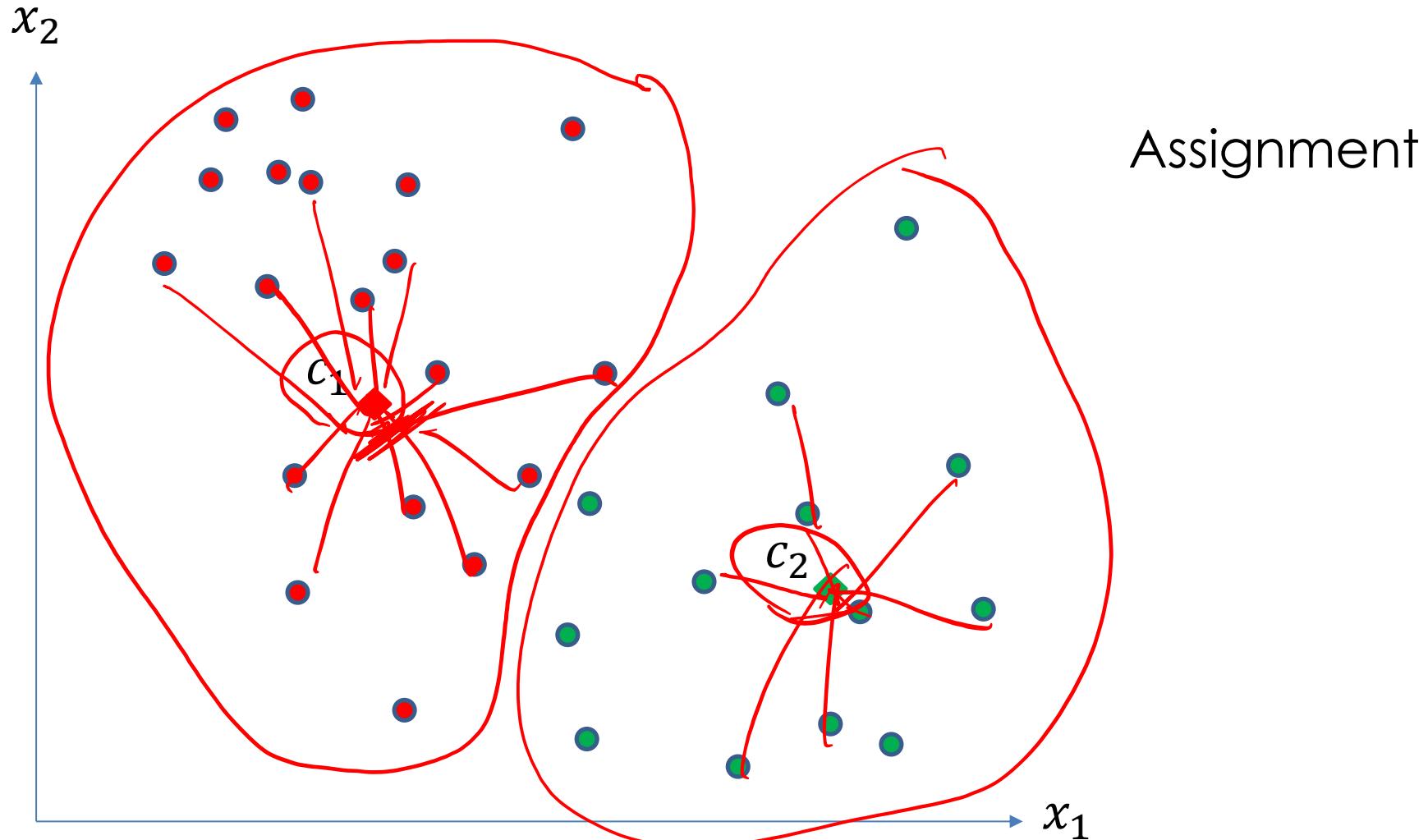
# K-means Clustering



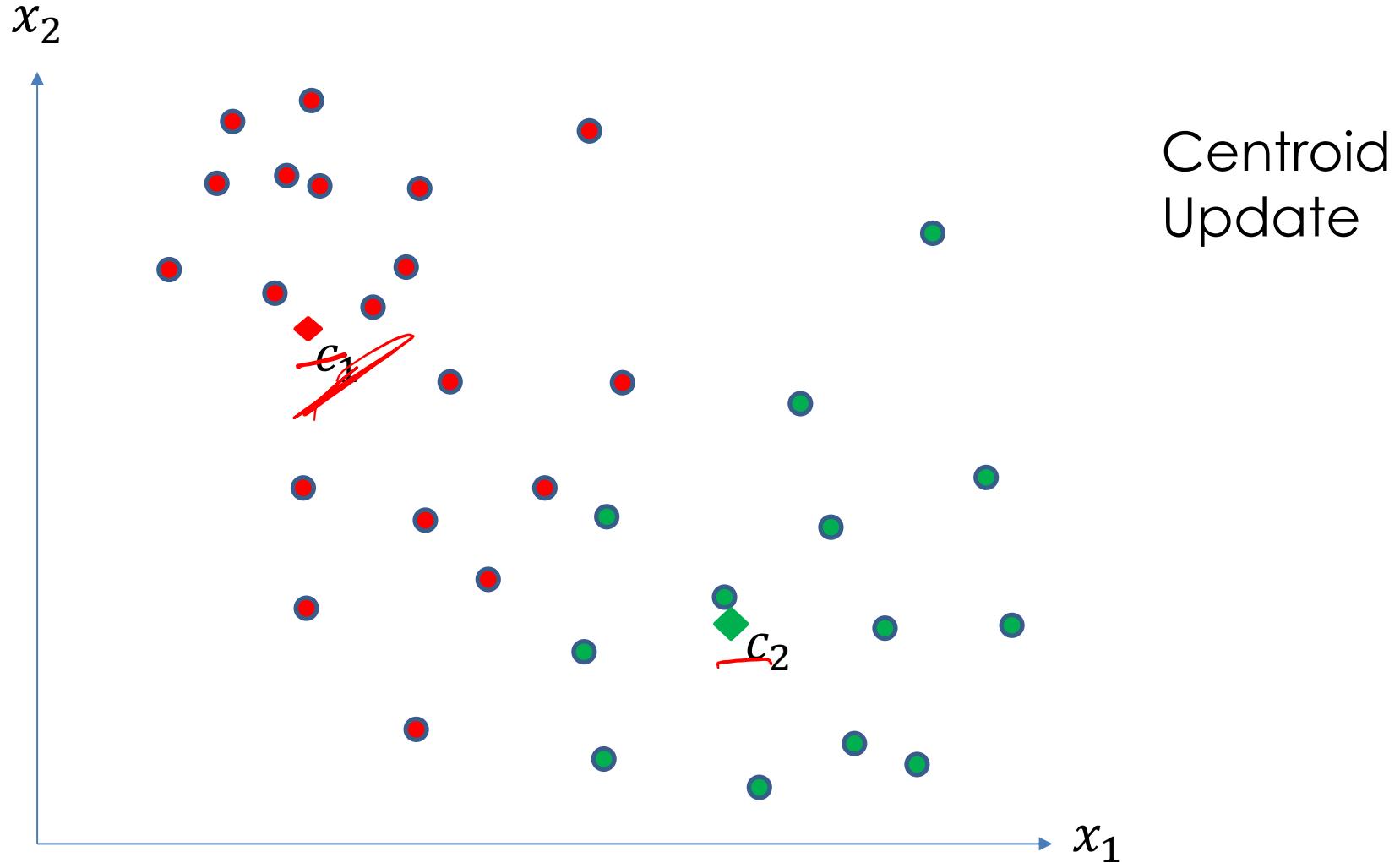
# K-means Clustering



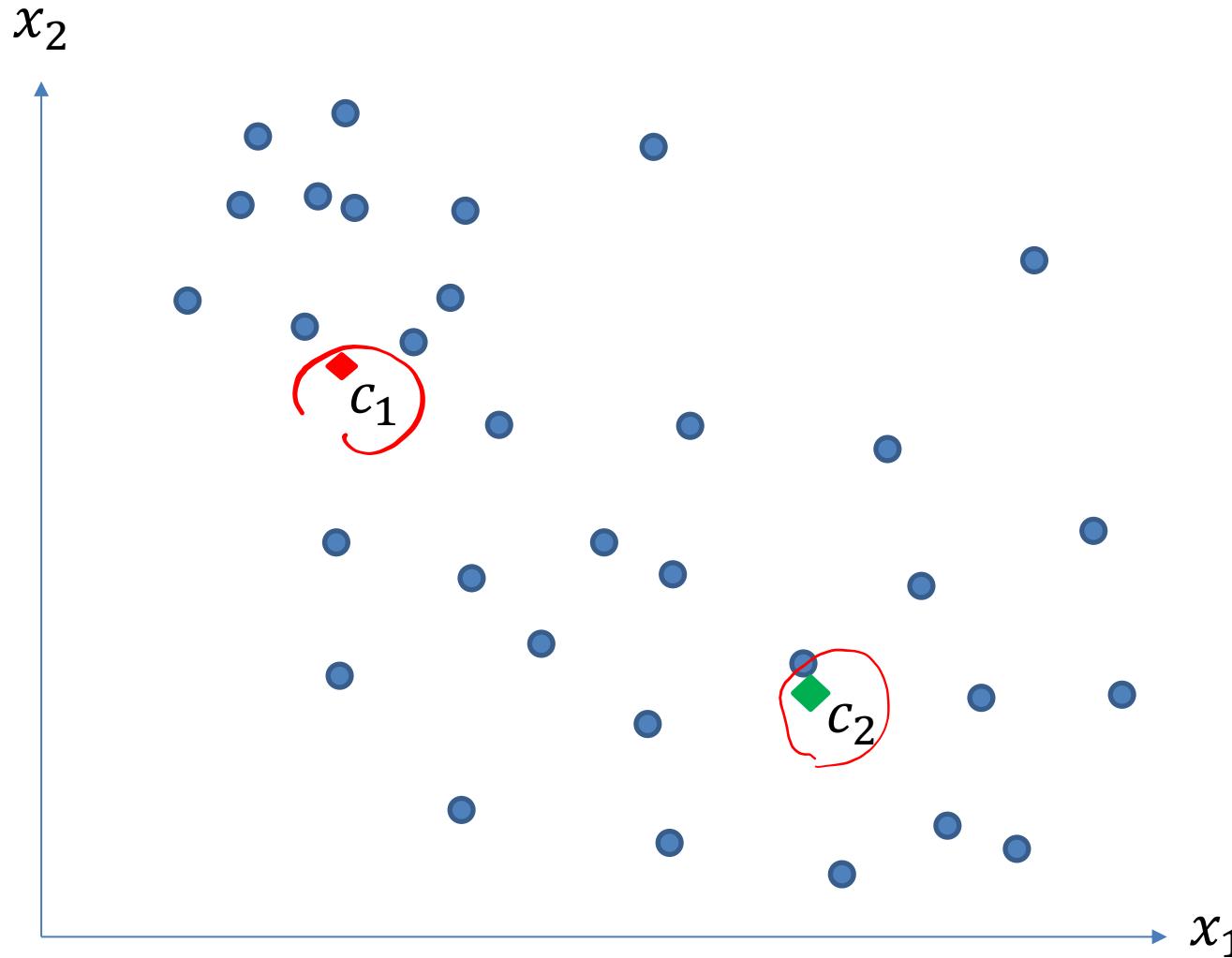
# K-means Clustering



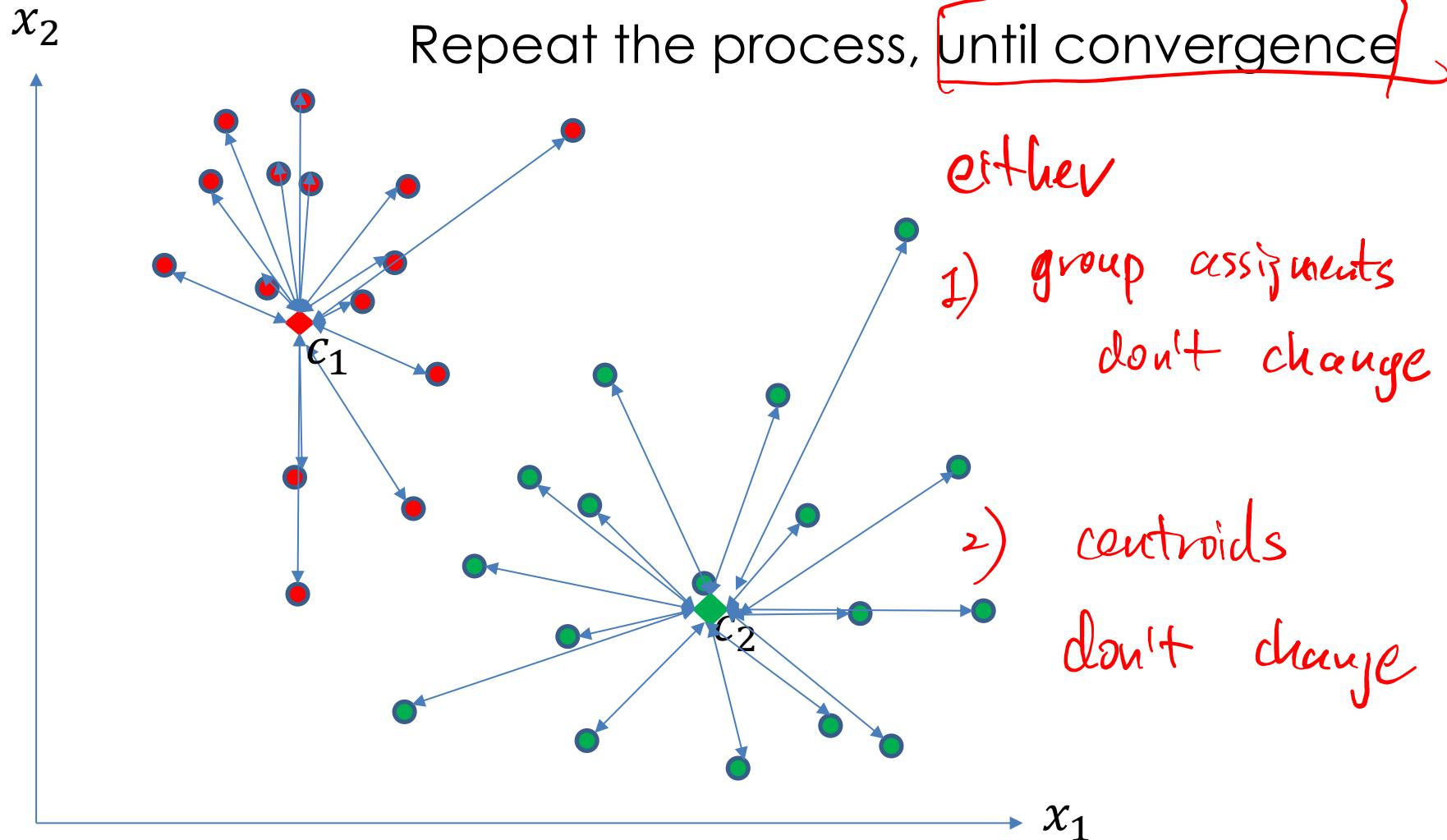
# K-means Clustering



# K-means Clustering



# K-means Clustering



# K-means Clustering

## Basic/Naïve K-means Clustering

Looping between Assignment and Centroid Update

initialization

1. First, we choose  $K$  — the number of clusters. Then we randomly select  $K$  feature vectors, called centroids, to the feature space.

assignment

2. Next, compute the distance from each example  $x$  to each centroid  $c$  using some metric, like the Euclidean distance. Then we assign the closest centroid to each example (like if we labeled each example with a centroid id as the label).

update

3. For each centroid, we calculate the average feature vector of the examples labeled with it. These average feature vectors become the new locations of the centroids.

centroid

4. We recompute the distance from each example to each centroid, modify the assignment and repeat the procedure until the assignments don't change after the centroid locations are recomputed.

assignment

5. Finally, we conclude the clustering with a list of assignments of centroids IDs to the examples.



```

# Define the k-means function
def kmeans_step(data, k, centroids):
    # Assign each data point to the closest centroid
    distances = np.sqrt(((data - centroids[:, np.newaxis])**2).sum(axis=2))
    labels = np.argmin(distances, axis=0)

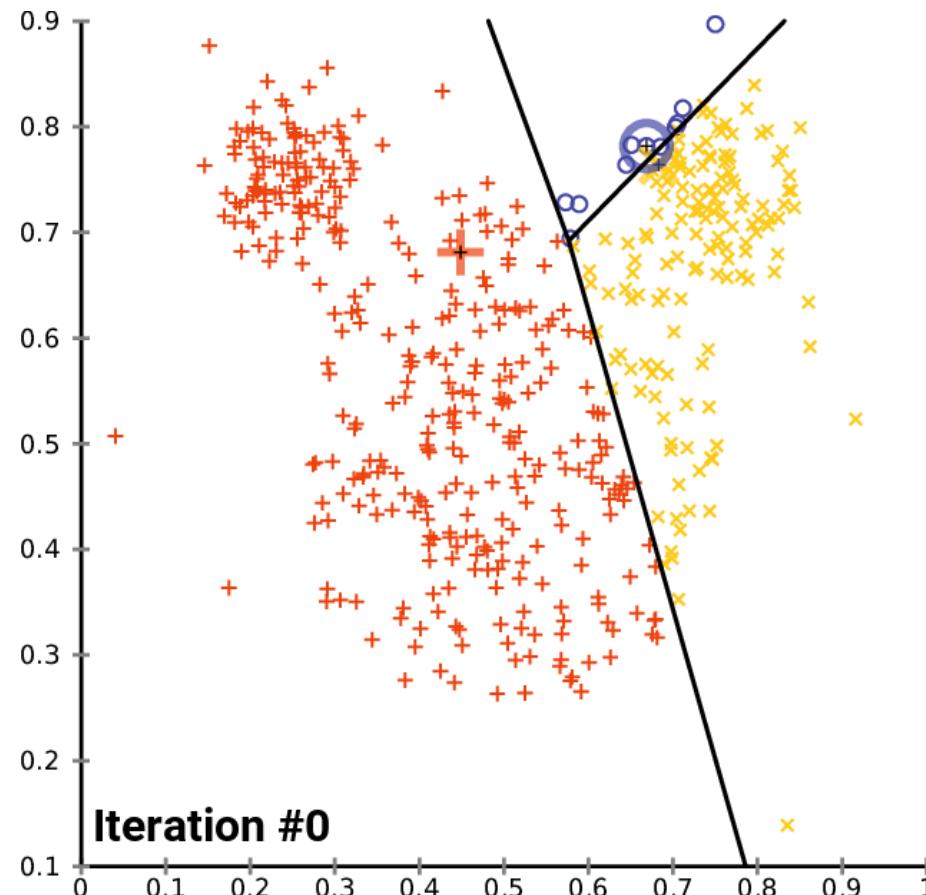
    # Update centroids to be the mean of the data points assigned to them
    new_centroids = np.zeros_like(centroids)
    for j in range(k):
        new_centroids[j] = np.mean(data[labels == j], axis=0)

    # End if centroids no longer change
    if np.linalg.norm(new_centroids - centroids) < tolerance:
        print("End Clustering, Centroids no change.")
        # Return the original centroids and labels, and set end to True
        return centroids, labels, True
    else:
        # Return the centroids and labels, and set end to False
        return new_centroids, labels, False
  
```

- See Python code:
  - lec11.ipynb
- See live demo at:
  - lec11\_kmeans.html

**All available in Canvas**  
**Files\For Students\ Lecture Notes**

# K-means Clustering



[https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering)

# K-means Clustering

## Optimization Objective Function (within-cluster variance)

Minimize  $J$

$$J = \sum_{i=1}^m \sum_{k=1}^K w_{ik} \|x_i - c_k\|^2 \quad (1)$$

m : # of samples; i: index of samples  
 K : # of clusters; k: index of clusters

assignment      k-th centroid  
 i-th sample

$w_{ik}$   $\rightarrow$  sample  $\leftarrow$  group

The term  $w_{ik}$  is equal to 1 for data point  $x_i$  if the data point belongs to cluster  $S_k$ , else  $w_{ik} = 0$ .

Note: The optimization objective function was called  $C(\mathbf{w})$  in Lecture 8. Here, we use  $J$  (with parameters  $w_{ik}$  and  $c_k$ ) so that it is differentiated from the centroids  $c_k$ .

Ref: <https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a>  
[https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering)

# K-means Clustering

## Optimization Objective Function (within-cluster variance)

Minimize  $J$

$$J = \sum_{i=1}^m \sum_{k=1}^K w_{ik} \|x_i - c_k\|^2 \quad (1)$$

$w_{11} = ?$   $x_1$   
 $w_{12} = ?$   
 $c_1$   
 $x_2$   
 $x_3$   
 $c_2$   
 $x_4$   
 $w_{41} = ?$   
 $w_{42} = ?$   
 $1$

The term  $w_{ik}$  is equal to 1 for data point  $x_i$  if the data point belongs to cluster  $S_k$ , else  $w_{ik} = 0$ .

Note: The optimization objective function was called  $C(\mathbf{w})$  in Lecture 8. Here, we use  $J$  (with parameters  $w_{ik}$  and  $\mathbf{c}_k$ ) so that it is differentiated from the centroids  $\mathbf{c}_k$ .

Ref: <https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a>  
[https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering)

# K-means Clustering

## Naïve K-means Algorithm

### 1. Assignment Step (fix $c$ and update $w$ ):

$\mathbf{x}_i \in S_k$  ( $w_{ik} = 1$ ) if  $\|\mathbf{x}_i - \mathbf{c}_k\|^2 < \|\mathbf{x}_i - \mathbf{c}_j\|^2$  (else  $w_{ik} = 0$ ),  
 $i = 1, \dots, m; j, k = 1, \dots, K.$

Computing distances to all centroids

### 2. Update Step (fix $w$ and update $c$ ):

$$\frac{\partial J}{\partial \mathbf{c}_k} = -2 \sum_{i=1}^m w_{ik} (\mathbf{x}_i - \mathbf{c}_k) = 0 \Rightarrow \mathbf{c}_k = \frac{\sum_{i=1}^m w_{ik} \mathbf{x}_i}{\sum_{i=1}^m w_{ik}}$$

Solving an optimization, i.e., setting derivative to 0

Note:  $\|\mathbf{x} - \mathbf{c}\| = \sqrt{\sum_{d=1}^D (x_d - c_d)^2}$  is called the Euclidean distance.  
 where  $\mathbf{x} = (x_1, x_2, \dots, x_D)$ ,  $\mathbf{c} = (c_1, c_2, \dots, c_D)$

# K-means Clustering

## 1. Assignment Step (fix $c$ and update $w$ ):

$$\underline{\mathbf{x}_i \in S_k} \quad (w_{ik} = 1) \text{ if } \|\mathbf{x}_i - \mathbf{c}_k\|^2 < \|\mathbf{x}_i - \mathbf{c}_j\|^2 \text{ (else } w_{ik} = 0), \quad i = 1, \dots, m; j, k = 1, \dots, K.$$

## 2. Update Step (fix $w$ and update $c$ ):

$$\boxed{\frac{\partial J}{\partial \mathbf{c}_k}} = -2 \sum_{i=1}^m w_{ik} (\mathbf{x}_i - \mathbf{c}_k) = 0 \Rightarrow \boxed{\mathbf{c}_k} = \frac{\sum_{i=1}^m w_{ik} \mathbf{x}_i}{\sum_{i=1}^m w_{ik}}$$

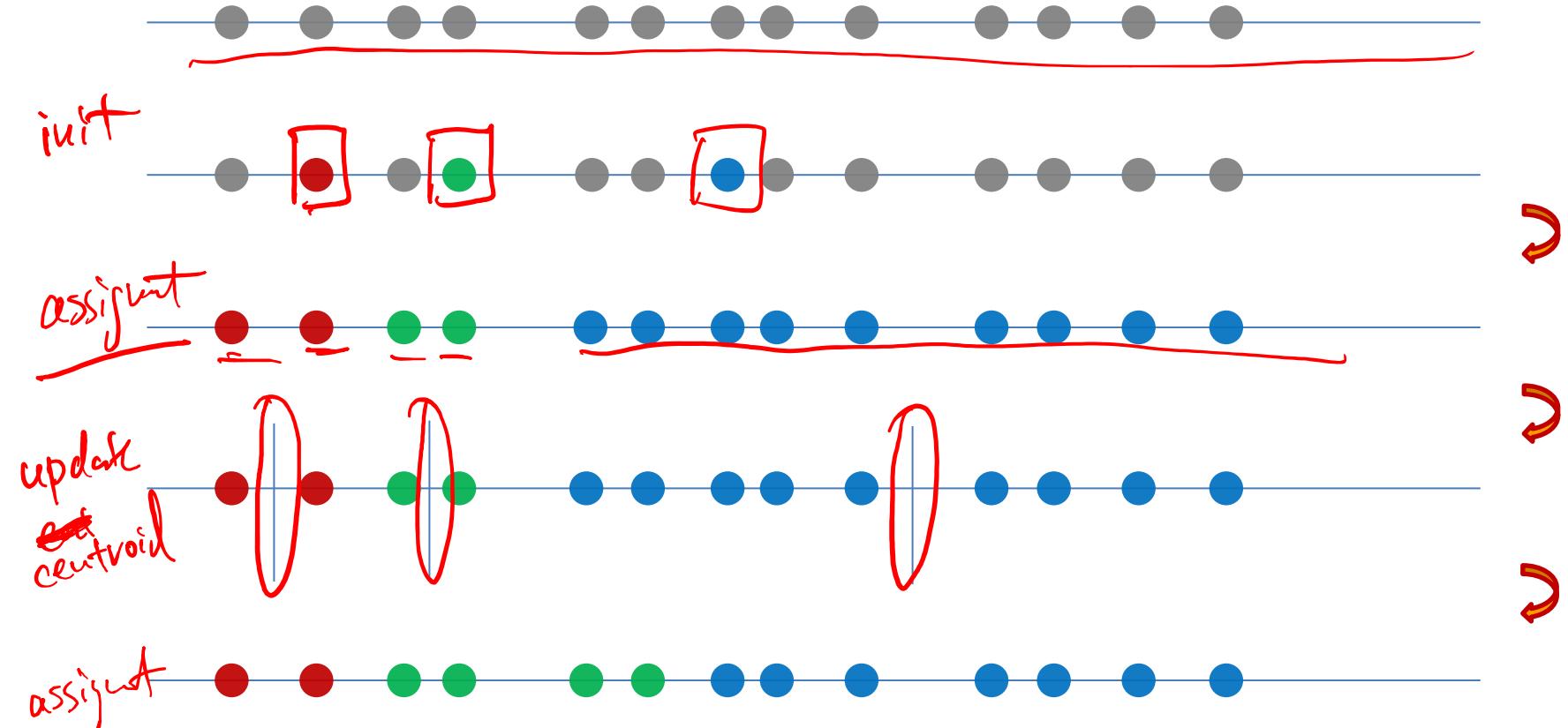
By repeating this two steps, the total loss  $J = \sum_{i=1}^m \sum_{k=1}^K w_{ik} \|\mathbf{x}_i - \mathbf{c}_k\|^2$ , is guaranteed to NOT increase (i.e., remain the same or decrease) until convergence.

Why? At Step 2: we compute the new mean, by solving an optimization, i.e., compute the derivative and set to zero, and solve  $\mathbf{c}_k$ .

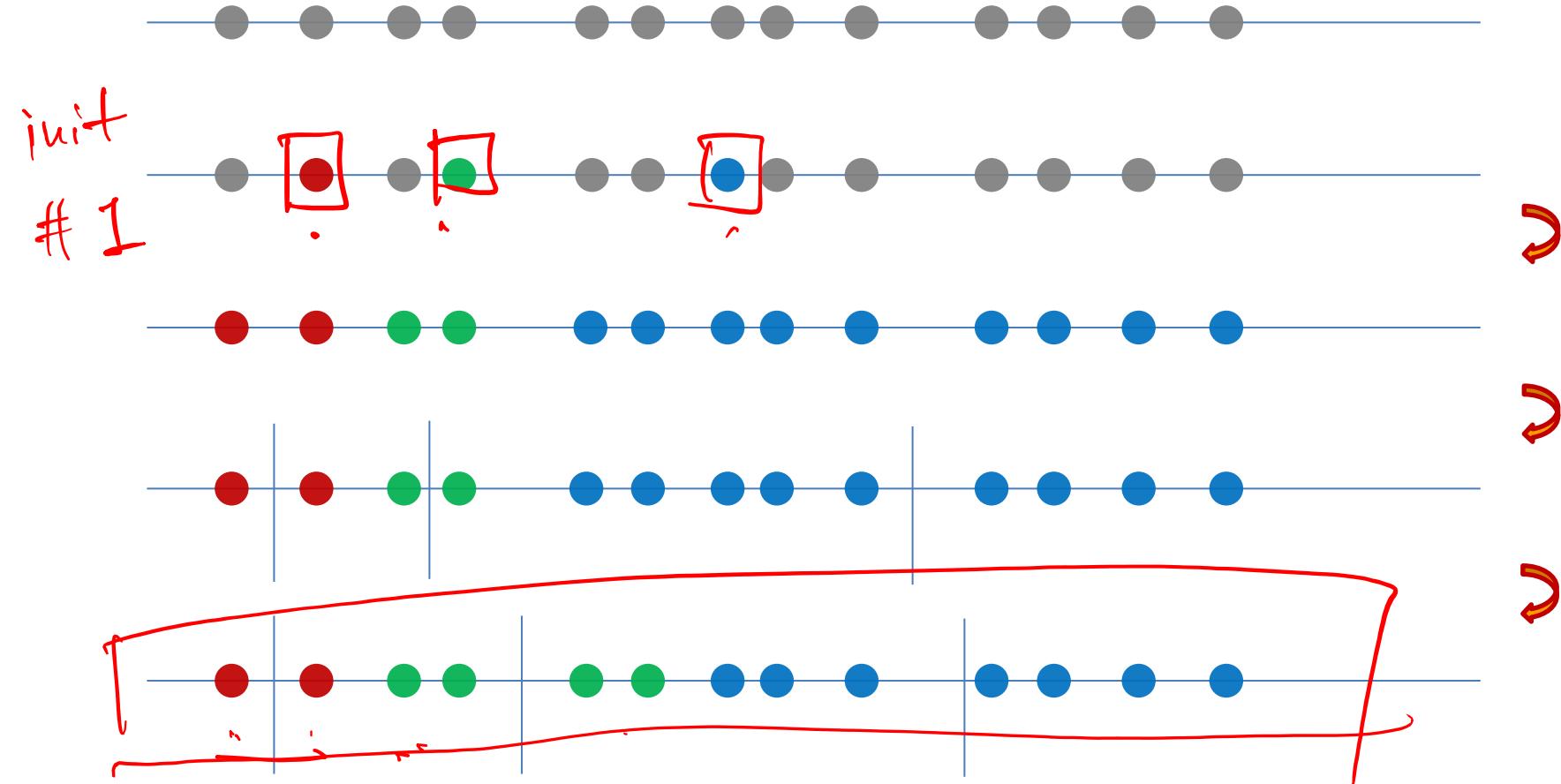
This means that, the new  $\mathbf{c}_k$  is guaranteed to give a smaller  $J$  value.

At Step 1: we only change the assignment, if the distance to the new centroid is smaller! In other words, we either remain in the old group, or change to a new group that is closer (i.e., gives a smaller  $J$ )

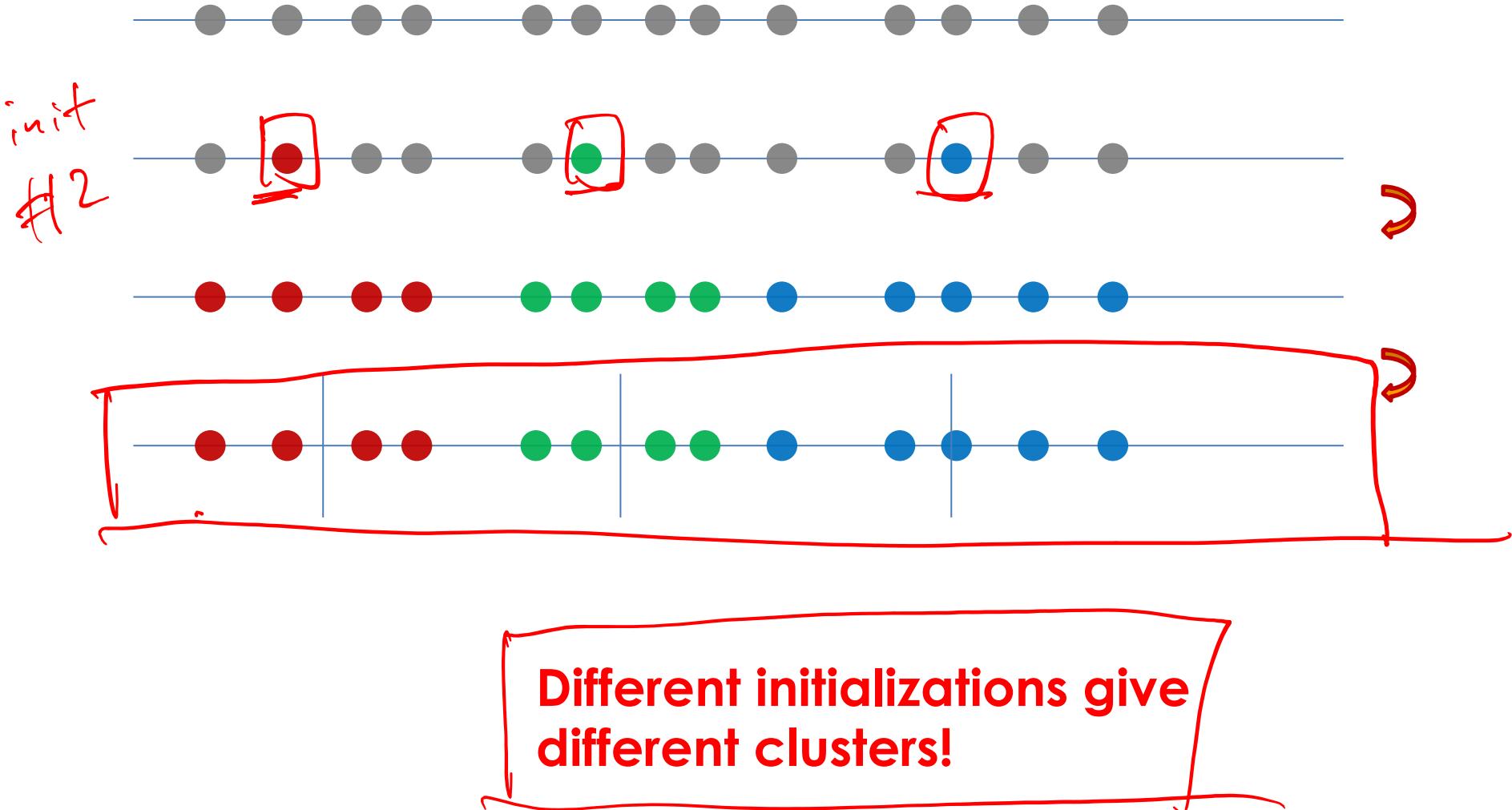
# K-means Clustering (1 D)



# K-means Clustering (1 D)



# K-means Clustering (1 D)



# K-means Clustering

- Unfortunately, k-means is not guaranteed to find a global minimum, it finds only local minimum.
- Example:  $K=2$



- Finding the optimal  $J$  is NP-hard\*
- In practice, k-means clustering usually performs well
- It can be very efficient, and its solution can be used as a starting point for other clustering algorithms

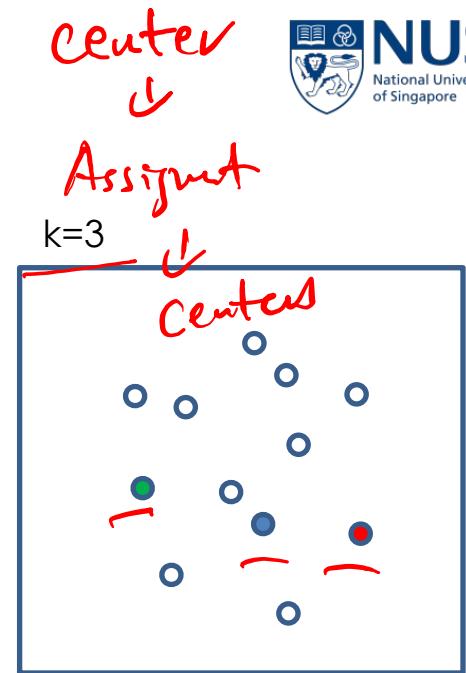
\*<https://en.wikipedia.org/wiki/NP-hardness>

# K-means Clustering

- Initialization

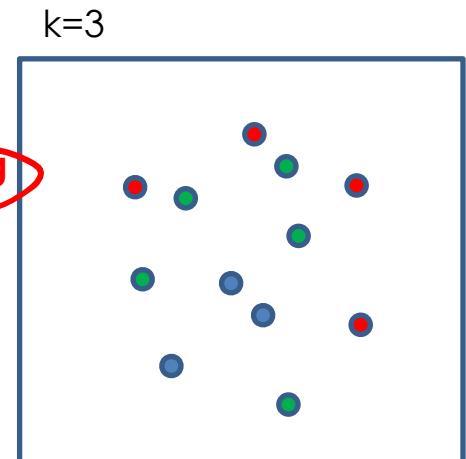
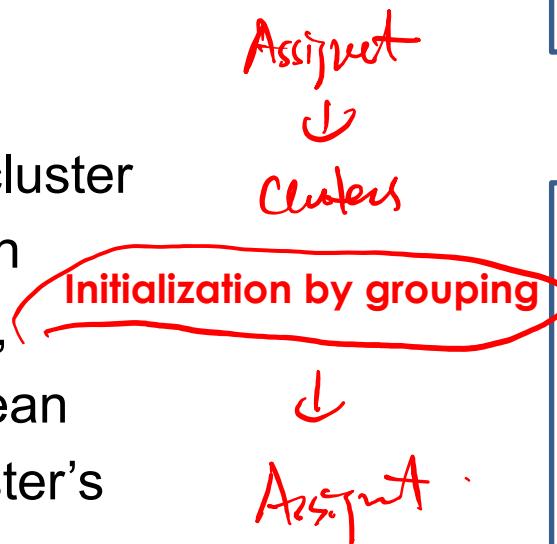
**Forgy method:** *Default K-means*

- Randomly chooses  $k$  observations from the dataset and uses these as the initial means.



**Random partition:**

- First randomly assigns a cluster to each observation and then proceeds to the update step, thus computing the initial mean to be the centroid of the cluster's randomly assigned points



Ref: [https://en.wikipedia.org/wiki/K-means\\_clustering#Standard\\_algorithm\\_\(naive\\_k-means\)](https://en.wikipedia.org/wiki/K-means_clustering#Standard_algorithm_(naive_k-means))

# Hard vs Soft Clustering

## Hard clustering:

Each data point can belong only one cluster, e.g. K-means

- For example, an apple can be red OR green (hard clustering)

## Soft clustering (also known as Fuzzy clustering):

Each data point can belong to more than one cluster.

- For example, an apple can be red AND green (fuzzy clustering)
- Here, the apple can be red to a certain degree as well as green to a certain degree.
- Instead of the apple belonging to green [green = 1] and not red [red = 0], the apple can belong to green [green = 0.5] and red [red = 0.5]. These values are normalized between 0 and 1; however, they do not represent probabilities, so the two values do not need to add up to 1.

Ref: [https://en.wikipedia.org/wiki/Fuzzy\\_clustering](https://en.wikipedia.org/wiki/Fuzzy_clustering)

# Hard vs Soft Clustering

## Objective Function for Fuzzy C-means

**Minimize  $J$**

$$J = \sum_{i=1}^m \sum_{k=1}^C (w_{ik})^r \|x_i - c_k\|^2$$

where  $w_{ik} = \frac{1}{\sum_{j=1}^c \left( \frac{\|x_i - c_k\|}{\|x_i - c_j\|} \right)^{\frac{2}{r-1}}}$

**No need to memorize**

Each element,  $w_{ik} \in [0,1]$ , tells the degree to which element,  $x_i$ , belongs to cluster  $c_k$ .

The fuzzifier  $r > 1$  determines the level of cluster fuzziness; usually  $1.25 \leq r \leq 2$ .

# Hard vs Soft Clustering

## Objective Function for Fuzzy C-means

Minimize  $J$

$$J = \sum_{i=1}^m \sum_{k=1}^c (w_{ik})^r \|x_i - c_k\|^2$$

where  $w_{ik} = \frac{1}{\sum_{j=1}^c \left( \frac{\|x_i - c_k\|}{\|x_i - c_j\|} \right)^{r-1}}$

$$\begin{aligned} w_{11} &= 0.6 \\ w_{12} &= 0.2 \end{aligned}$$

$x_1$  ●  
 $c_1$  ◆  
 $x_2$  ○

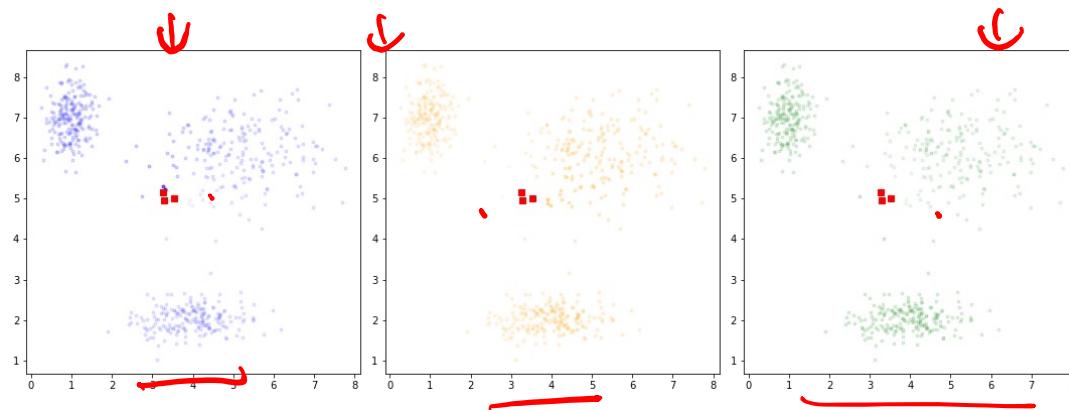
●  
 ■  $c_2$   
 ○  $x_4$

$$\begin{aligned} w_{41} &= 0.18 \\ w_{42} &= 0.75 \end{aligned}$$

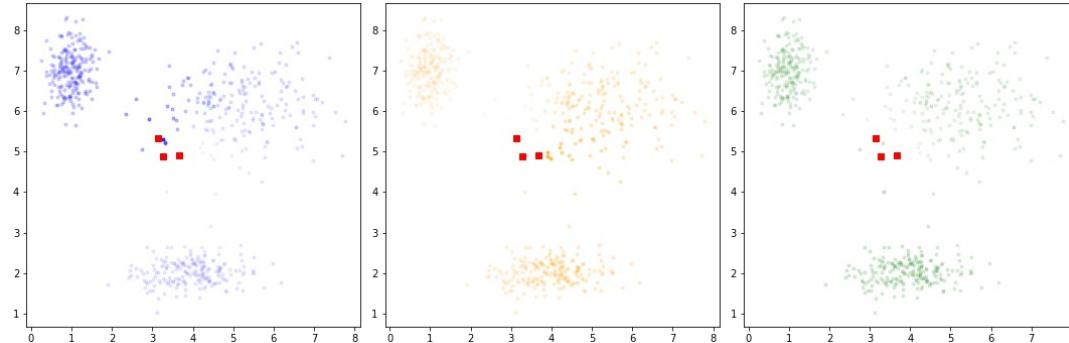
Each element,  $w_{ik} \in [0,1]$ , tells the degree to which element,  $x_i$ , belongs to cluster  $c_k$ .

The fuzzifier  $r > 1$  determines the level of cluster fuzziness; usually  $1.25 \leq r \leq 2$ .

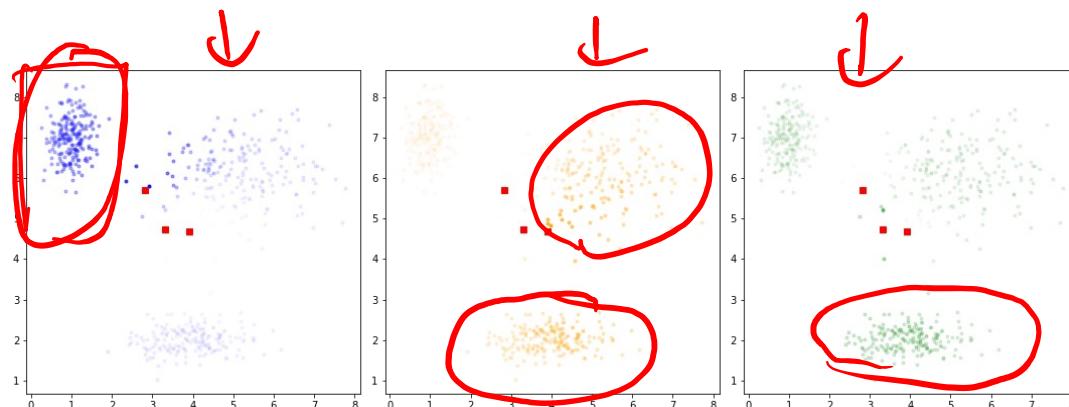
Iteration 1



Iteration 2

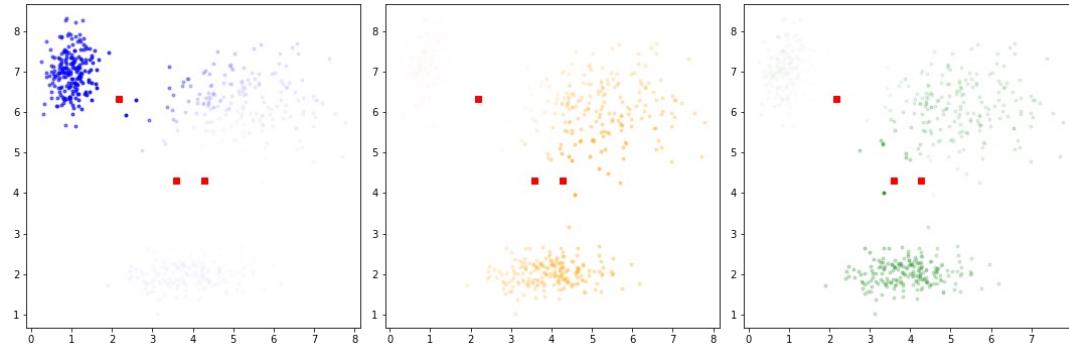


Iteration 3

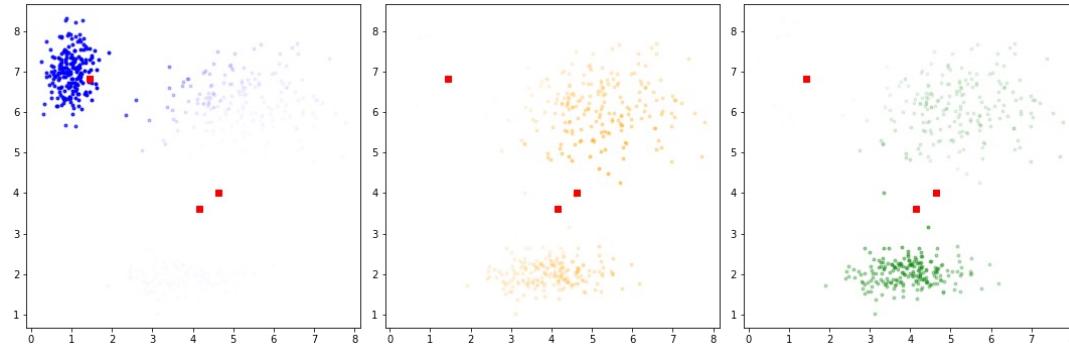


## Visualization of Fuzzy C-means Iterations

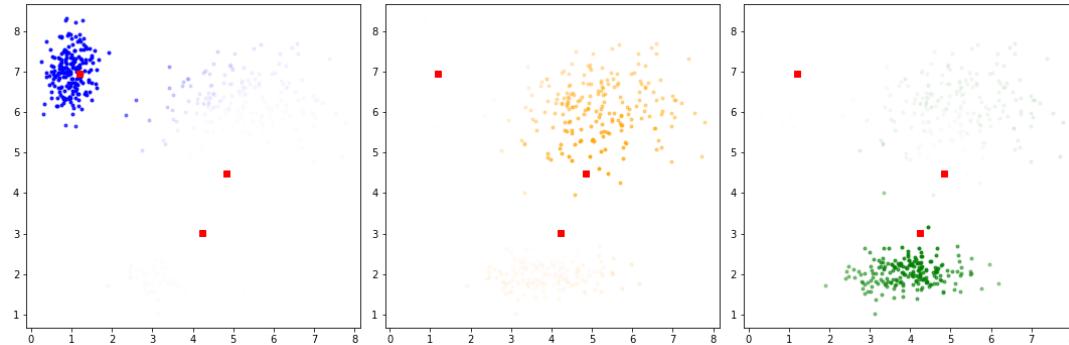
Iteration 4



Iteration 5

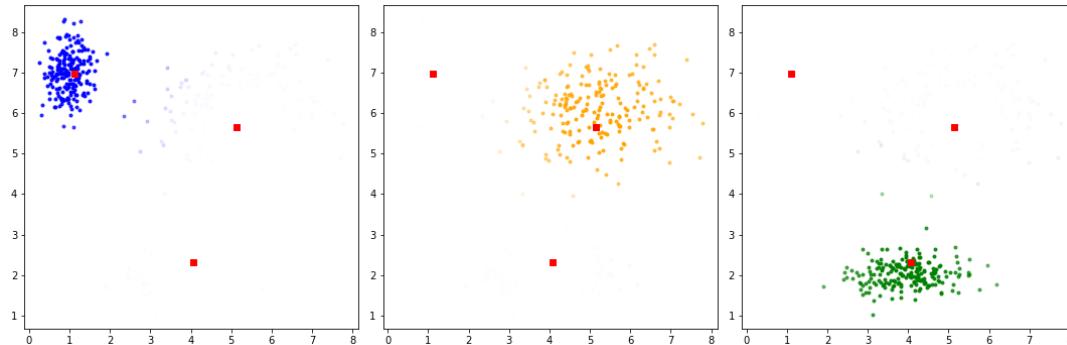


Iteration 6

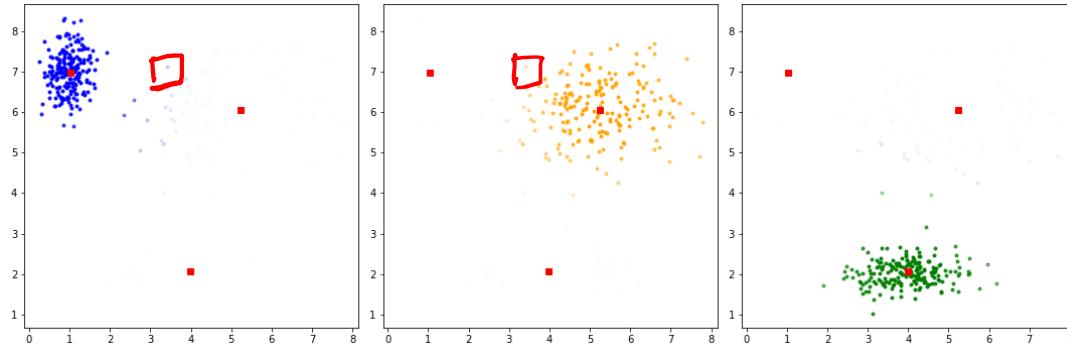


## Visualization of Fuzzy C-means Iterations

Iteration 7



Iteration 8

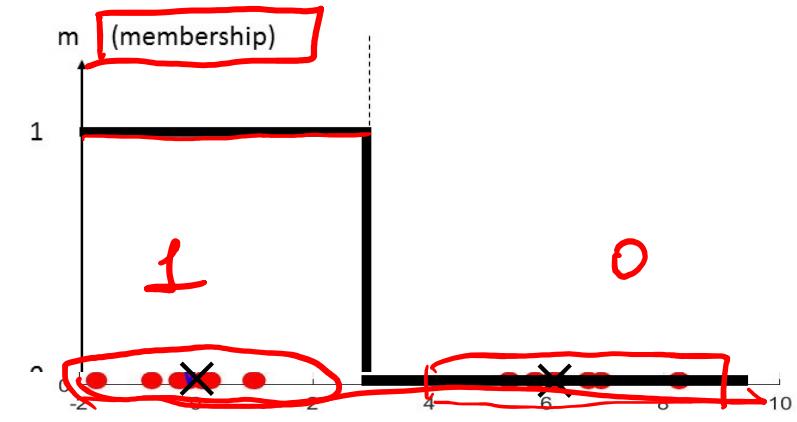


## Visualization of Fuzzy C-means Iterations

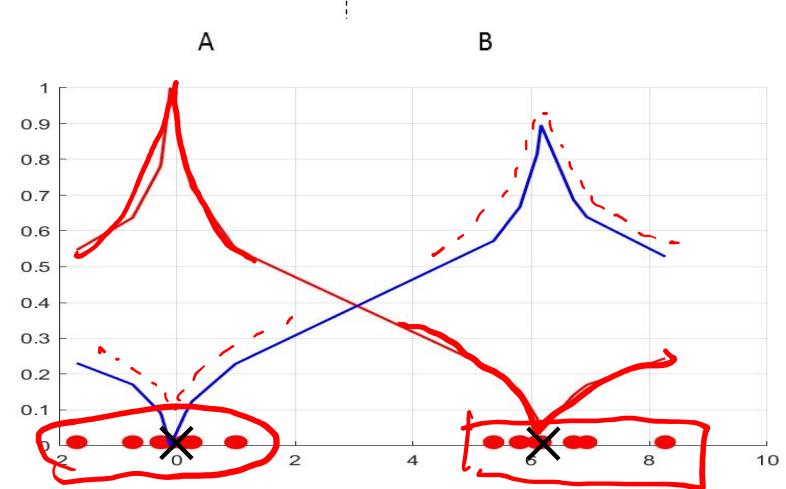
# Hard vs Soft Clustering

## Naïve K-means versus Fuzzy C-means

Naïve K-means:  $w_{ik} \in \{0,1\}$



Fuzzy C-means:  $w_{ik} \in [0,1]$



Ref: [https://en.wikipedia.org/wiki/Fuzzy\\_clustering](https://en.wikipedia.org/wiki/Fuzzy_clustering)

# Summary

- Introduction of unsupervised learning
- K-means Clustering      Hard clustering
  - The most popular clustering technique
- Fuzzy Clustering      Soft clustering

