

***K*-means clustering**

Algorithm

Herman Kamper

<http://www.kamperh.com/>

K -means clustering algorithm and example

1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :

Calculate the cluster centroid $\boldsymbol{\mu}_k$ as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :

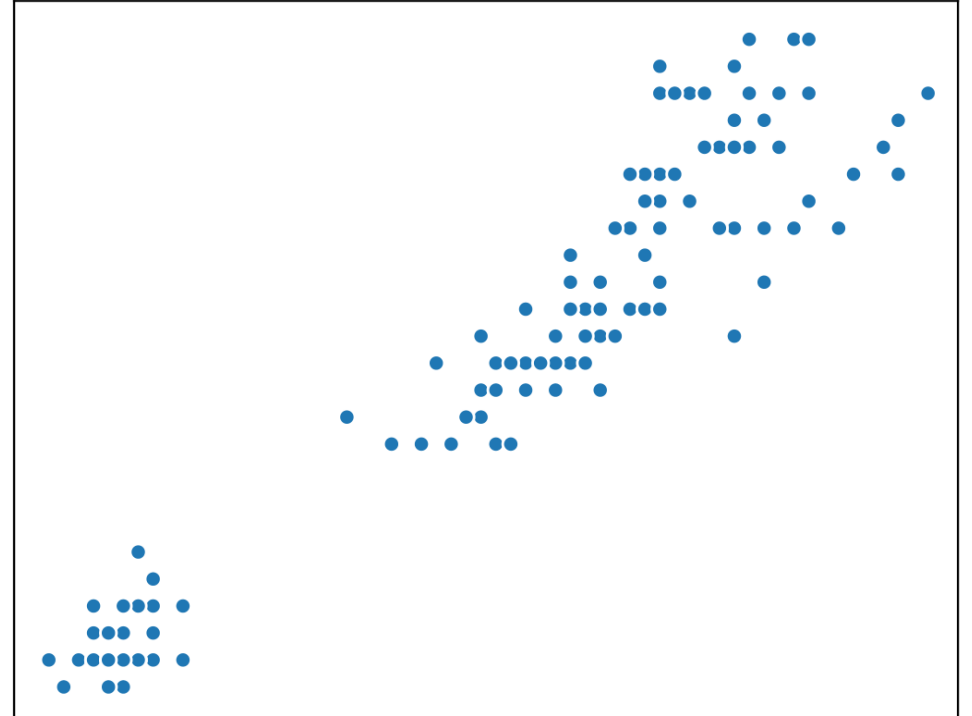
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.

K -means clustering algorithm and example

Example: $K=3$

1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :
Calculate the cluster centroid $\boldsymbol{\mu}_k$ as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.

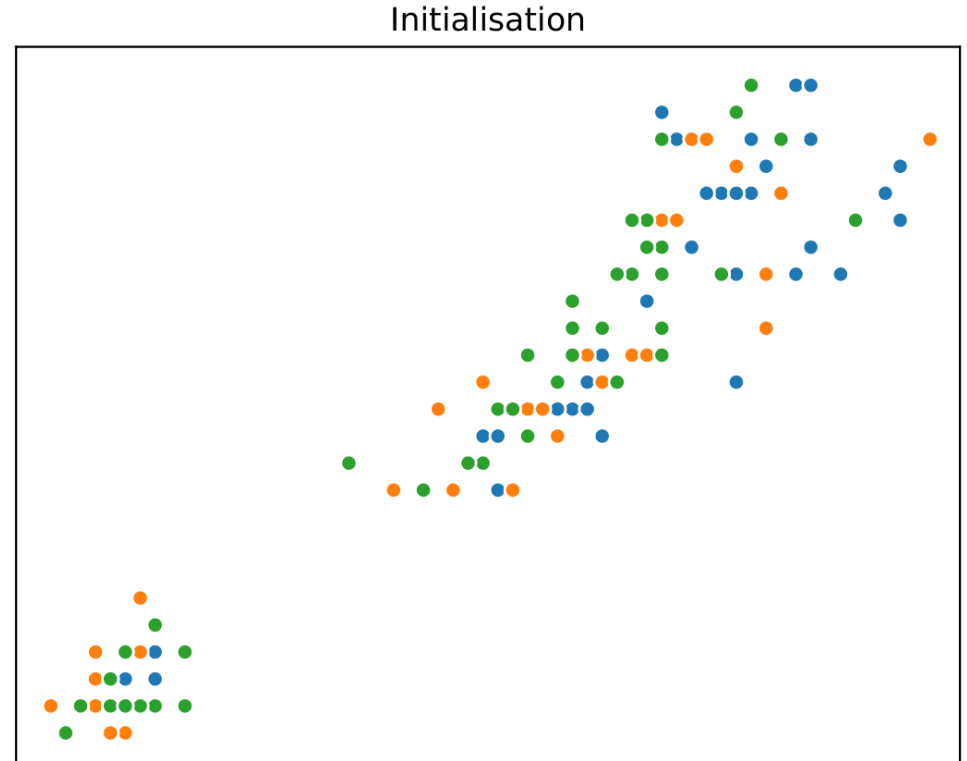
x_2
↑



→
 x_1

K -means clustering algorithm and example

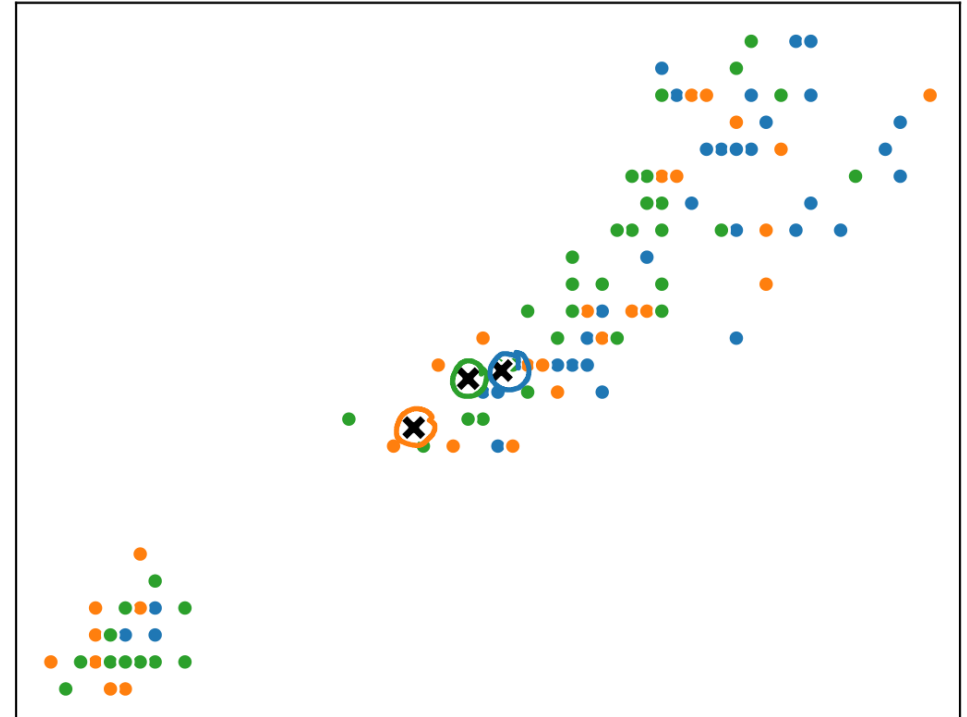
1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :
Calculate the cluster centroid $\boldsymbol{\mu}_k$ as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.



K -means clustering algorithm and example

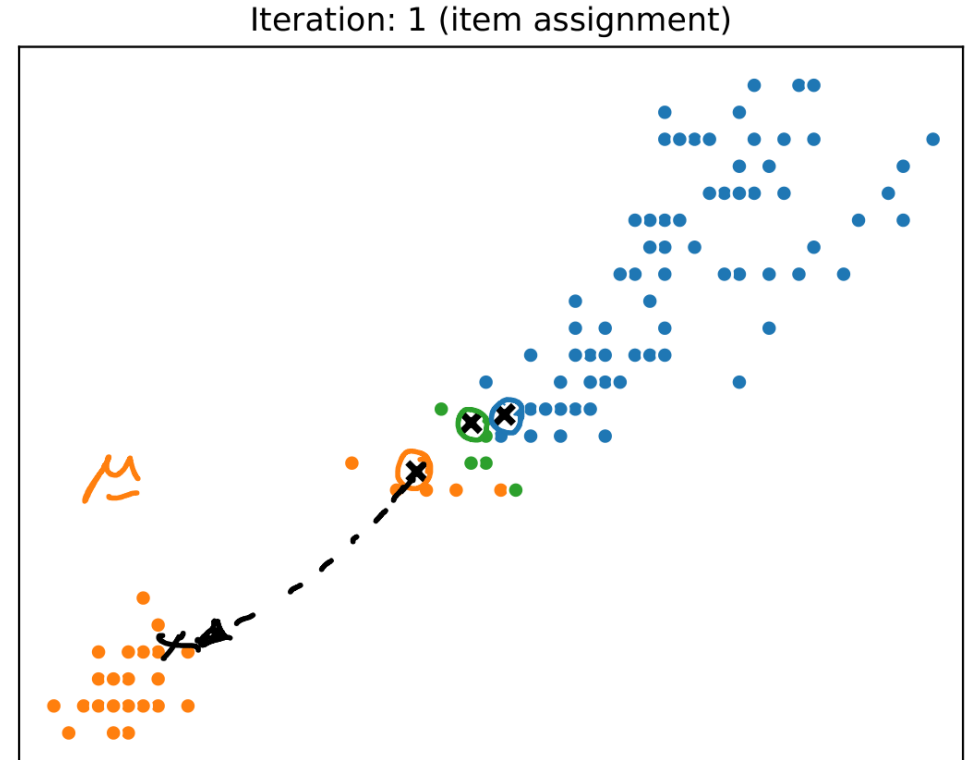
1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :
Calculate the cluster centroid $\boldsymbol{\mu}_k$ as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.

Iteration: 1 (centroid update)



K -means clustering algorithm and example

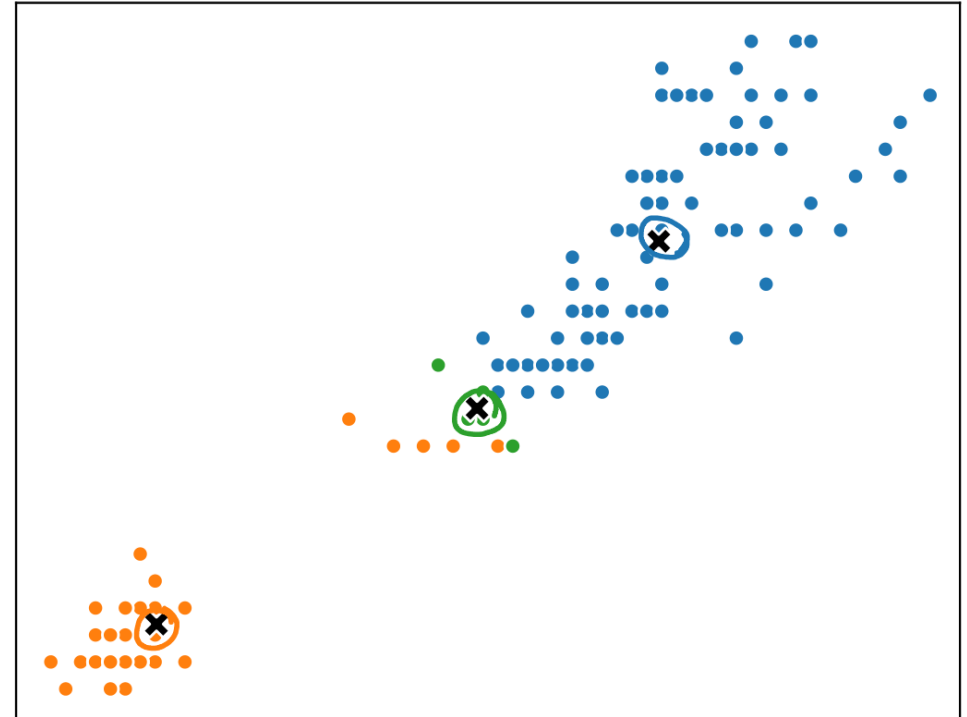
1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :
Calculate the cluster centroid μ_k as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.



K -means clustering algorithm and example

1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :
Calculate the cluster centroid $\boldsymbol{\mu}_k$ as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.

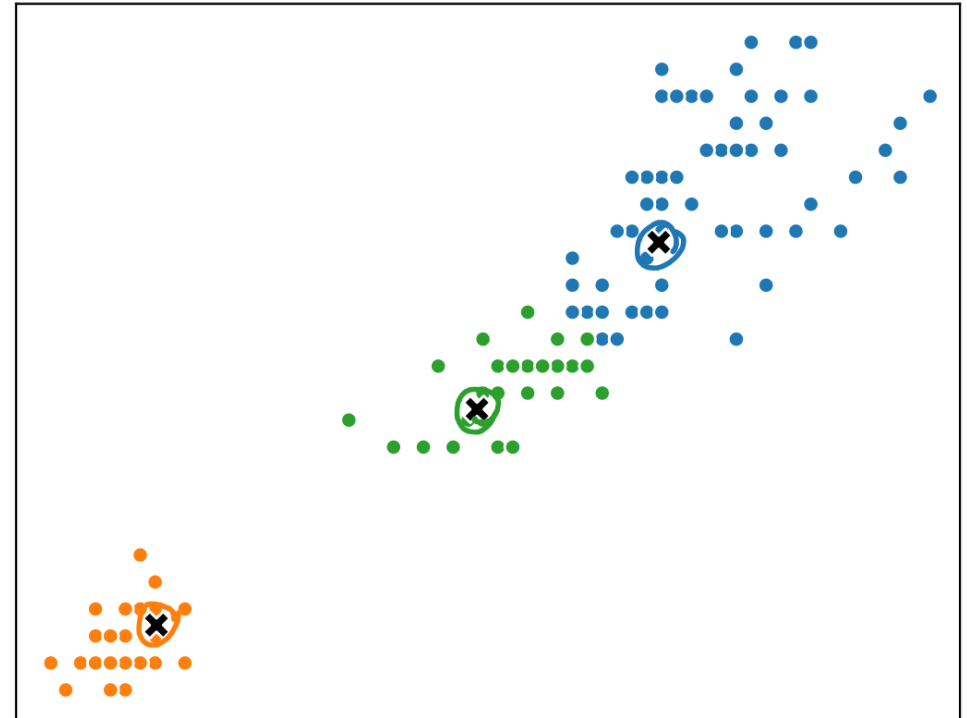
Iteration: 2 (centroid update)



K -means clustering algorithm and example

1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :
Calculate the cluster centroid $\boldsymbol{\mu}_k$ as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.

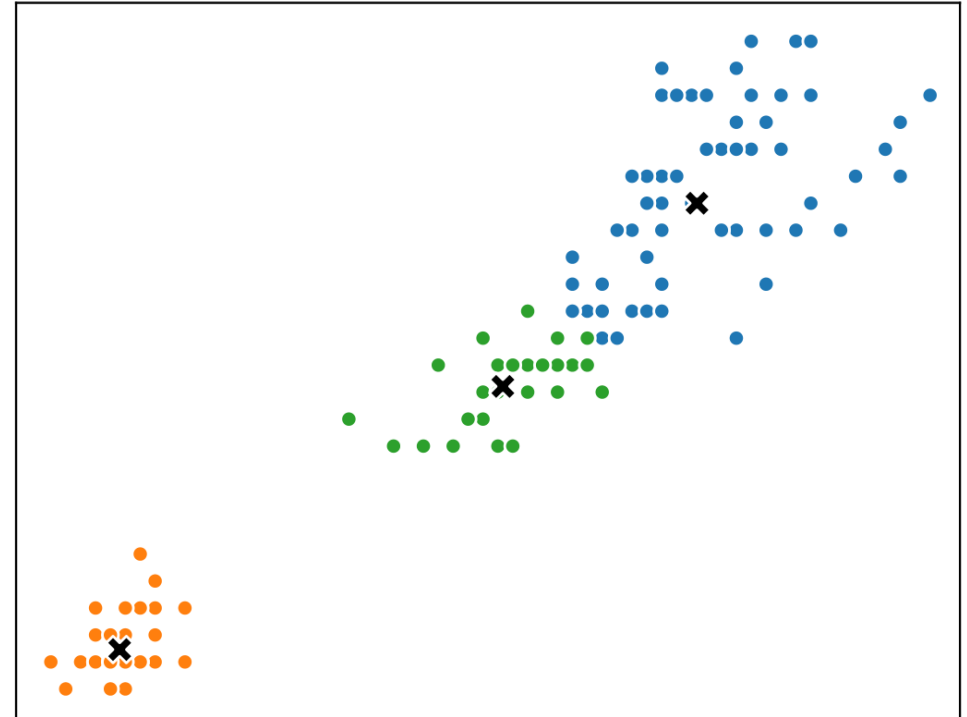
Iteration: 2 (item assignment)



K -means clustering algorithm and example

1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :
Calculate the cluster centroid $\boldsymbol{\mu}_k$ as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.

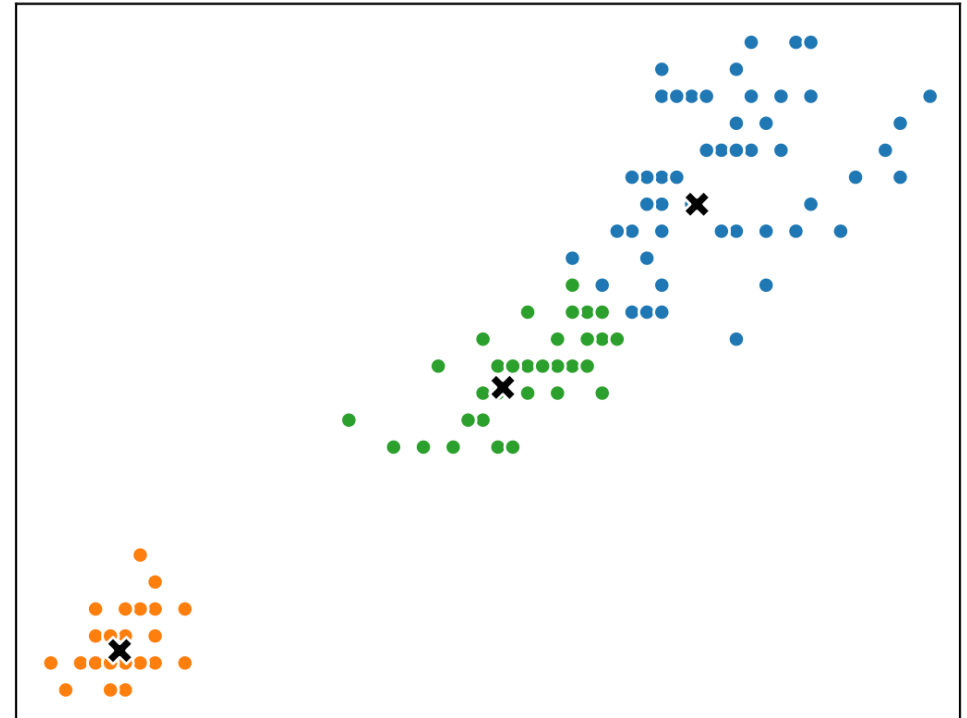
Iteration: 3 (centroid update)



K -means clustering algorithm and example

1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.
2. repeat until cluster assignments stop changing:
 - (a) for cluster $k = 1$ to K :
Calculate the cluster centroid $\boldsymbol{\mu}_k$ as the mean of all the items assigned to cluster k .
 - (b) for item $n = 1$ to N :
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.

Iteration: 3 (item assignment)



***K*-means clustering**

Details and loss

Herman Kamper

<http://www.kamperh.com/>

K -means clustering algorithm details

1. Randomly assign each item $\underline{x}^{(n)}$ to one of the K clusters.

2. repeat until cluster assignments stop changing:

(a) for cluster $k = 1$ to K :

Calculate the cluster centroid $\underline{\mu}_k$ as the mean of all the items assigned to cluster k .

$$\underline{\mu}_k = \frac{1}{|C_k|} \sum_{i \in C_k} \underline{x}^{(i)}$$

(b) for item $n = 1$ to N :

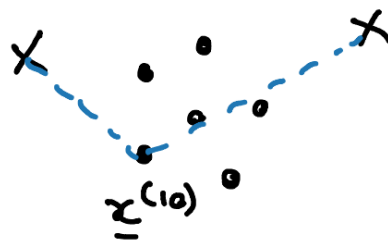
Assign item $\underline{x}^{(n)}$ to the cluster with the closest centroid.

$$\arg \min_k \|\underline{x}^{(n)} - \underline{\mu}_k\|^2$$

C_k : Set of indices of items assigned to cluster k

e.g. $C_4 = \{205, 12, 303\}$; $|C_4| = 3$
 $\underline{x}^{(205)}$

$|C_k|$: # items in cluster k



K -means clustering algorithm details

$$\text{Loss: } J(C_1, \dots, C_K, \underline{\mu}_1, \dots, \underline{\mu}_K) = \sum_{k=1}^K \underbrace{\sum_{i \in C_k} \|\mathbf{x}^{(i)} - \underline{\mu}_k\|^2}$$

1. Randomly assign each item $\mathbf{x}^{(n)}$ to one of the K clusters.

2. repeat until cluster assignments stop changing:

(a) for cluster $k = 1$ to K :

Calculate the cluster centroid $\underline{\mu}_k$ as the mean of all the items assigned to cluster k .

Centroid update:

Update $\underline{\mu}_1, \dots, \underline{\mu}_K$ while
keeping C_1, \dots, C_K

(b) for item $n = 1$ to N :

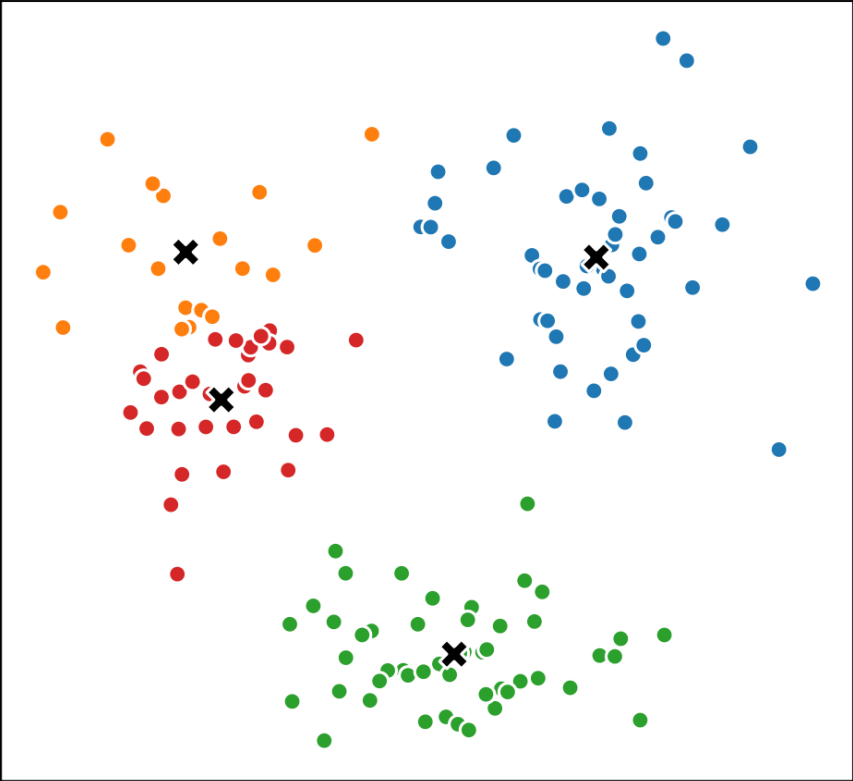
Assign item $\mathbf{x}^{(n)}$ to the cluster with the closest centroid.

Cluster assignments:

Update C_1, \dots, C_K while
keeping $\underline{\mu}_1, \dots, \underline{\mu}_K$ fixed

Random initialisation leads to different local optima

Sum of squared distances to centroids: 68.26



Sum of squared distances to centroids: 66.97

