# 14. Clustering:
# k-Means Clustering Algorithm

# k-Means Clustering: Overview

## Facts

It is a simple partition based clustering algorithm

It divides a number of data points into a number of clusters, with each cluster represented by a centroid

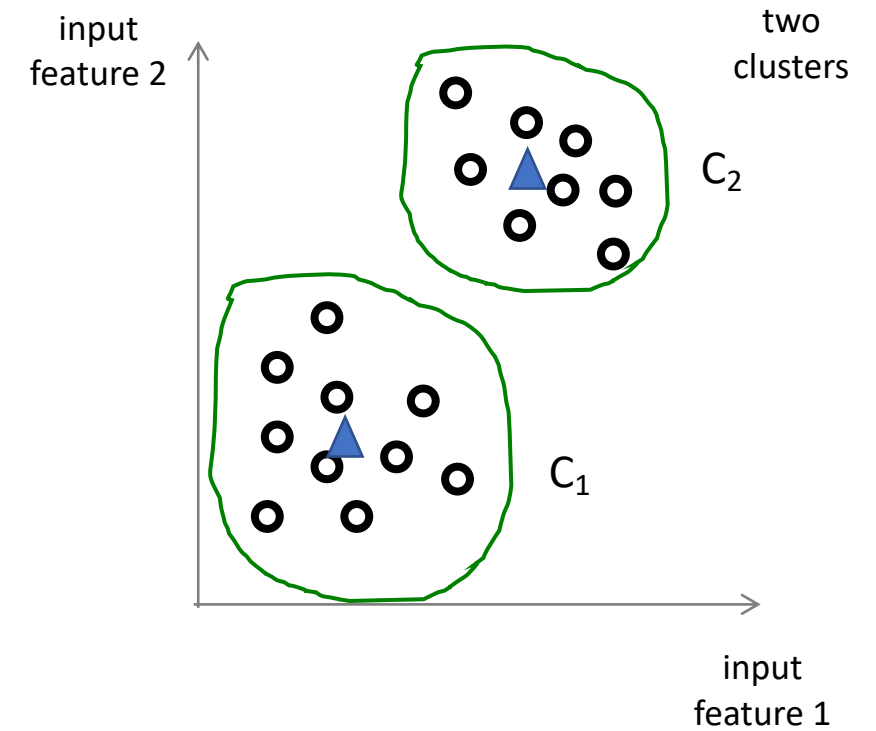It has been proposed by MacQueen in 1967

## Algorithm Overview

Input:

- The collection of data points
- The number of clusters

Output:

- The attribution of the data points to different clusters

## Example



input feature 2

two clusters

$C_2$

$C_1$

input feature 1

# k-Means Clustering: How Does it Work?

## Data

A collection of data points $x_i = (x_{i,k}, x_{i,2}, \ldots, x_{i,m})$, $\mathrm{i} = 1,2, \ldots, n$, that one wishes to organize onto $K$ different clusters, where each data point is characterized by different features.

## Distance Function

Given two data points $x_i$ and $x_j$, we can define their similarity via their squared Euclidean distance:

$$d(x_i, x_j) = \sum_{l=1}^{m} (x_{i,l} - x_{j,l})^2$$

## Cost Function

Given the $K$ clusters $C_1, C_2, \ldots, C_K$, we can then define the cost function

$$J(C_1, C_2, \ldots, C_K) = \frac{1}{K} \cdot \sum_{k=1}^{K} \sum_{x \in C_k} d(x, c_k) \text{ with } c_k = \frac{1}{|C_k|} \cdot \sum_{x \in C_k} x$$

## Optimization Problem

The clustering problem can then be posed as follows:

$$C_1^*, C_2^*, \ldots, C_K^* = \underset{C_1, C_2, \ldots, C_K}{\mathrm{argmin}} \, J(C_1, C_2, \ldots, C_K)$$

Computationally very complex problem

⇩

heuristic algorithm

# k-Means Clustering: Algorithm

## k-Means Clustering Algorithm

This is a very simple iterative algorithm involving various steps:

- The very first step involves choosing some initial cluster representatives (e.g. by choosing K random data samples from the dataset)

- The next steps involve iterating between two of operations

  - Attribution of data samples to clusters

  $$\boldsymbol{x} \in C_k \Longleftarrow d(\boldsymbol{x}, \boldsymbol{c}_k) < d(\boldsymbol{x}, \boldsymbol{c}_l), \forall l \neq k$$

  <span style="color:red">Minimize $J(C_1, C_2, \ldots, C_K)$ keeping constant the centroids</span>

  - Re-computation of cluster centroids

  $$\boldsymbol{c}_k = \frac{1}{|C_k|} \cdot \sum_{\boldsymbol{x} \in C_k} \boldsymbol{x}$$

  <span style="color:red">Chose centroids that minimize $J(C_1, C_2, \ldots, C_K)$</span>

- until convergence (e.g. until nothing changes)

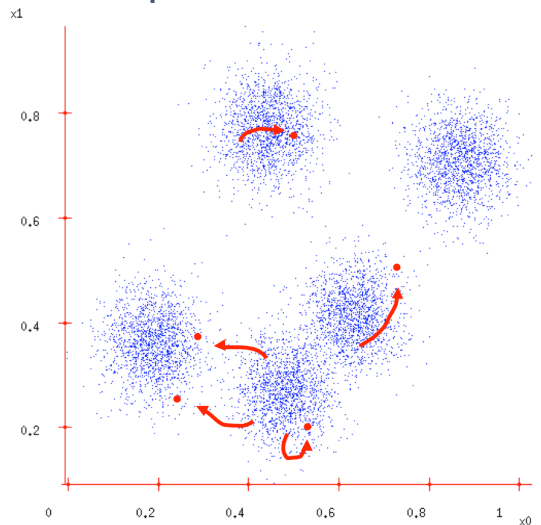# k-Means Clustering: Example

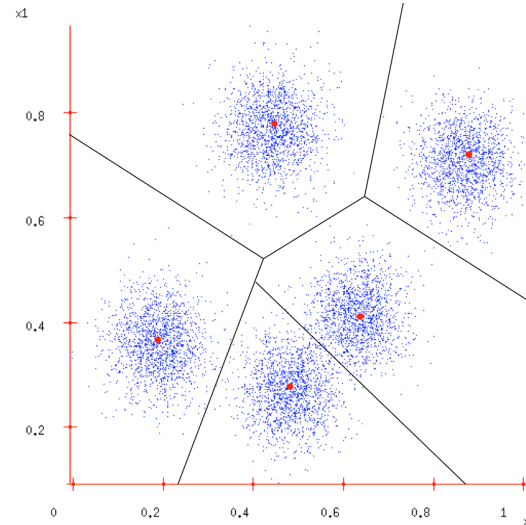

Example

Initialization of Centroids

Attribution of Data to Clusters
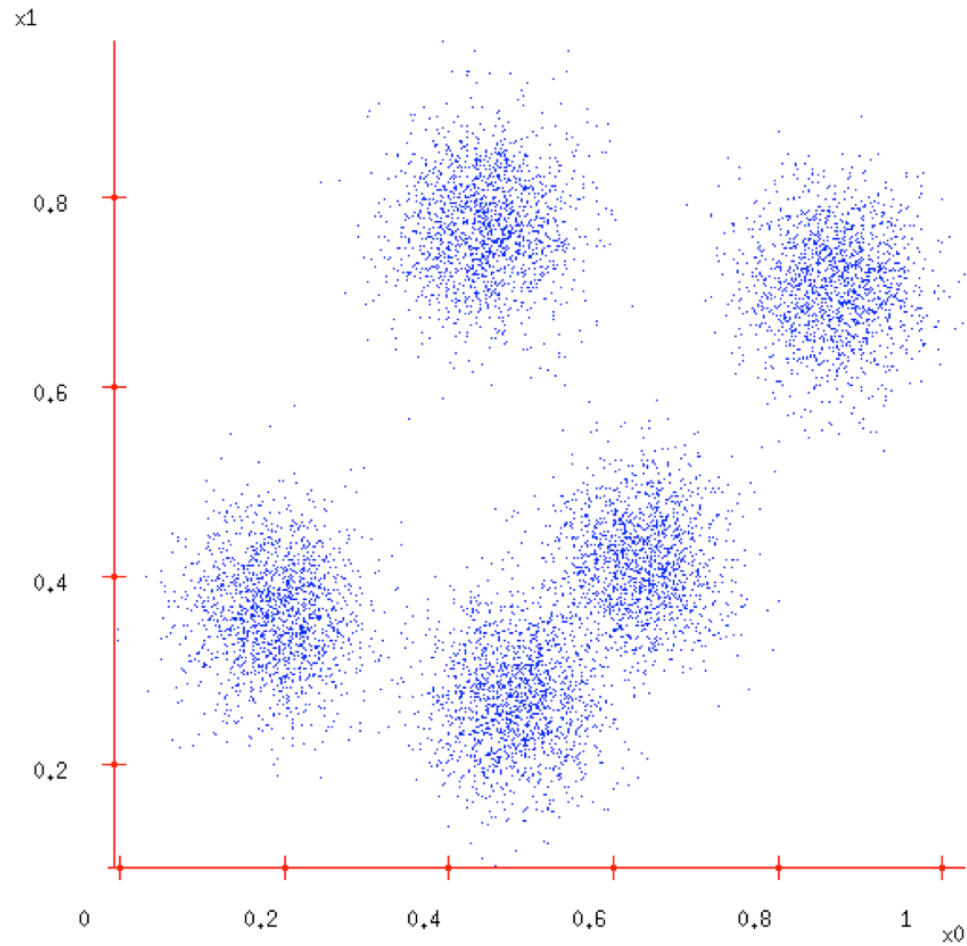
Re-Computation of Centroids

Re-Attribution of Data to Clusters

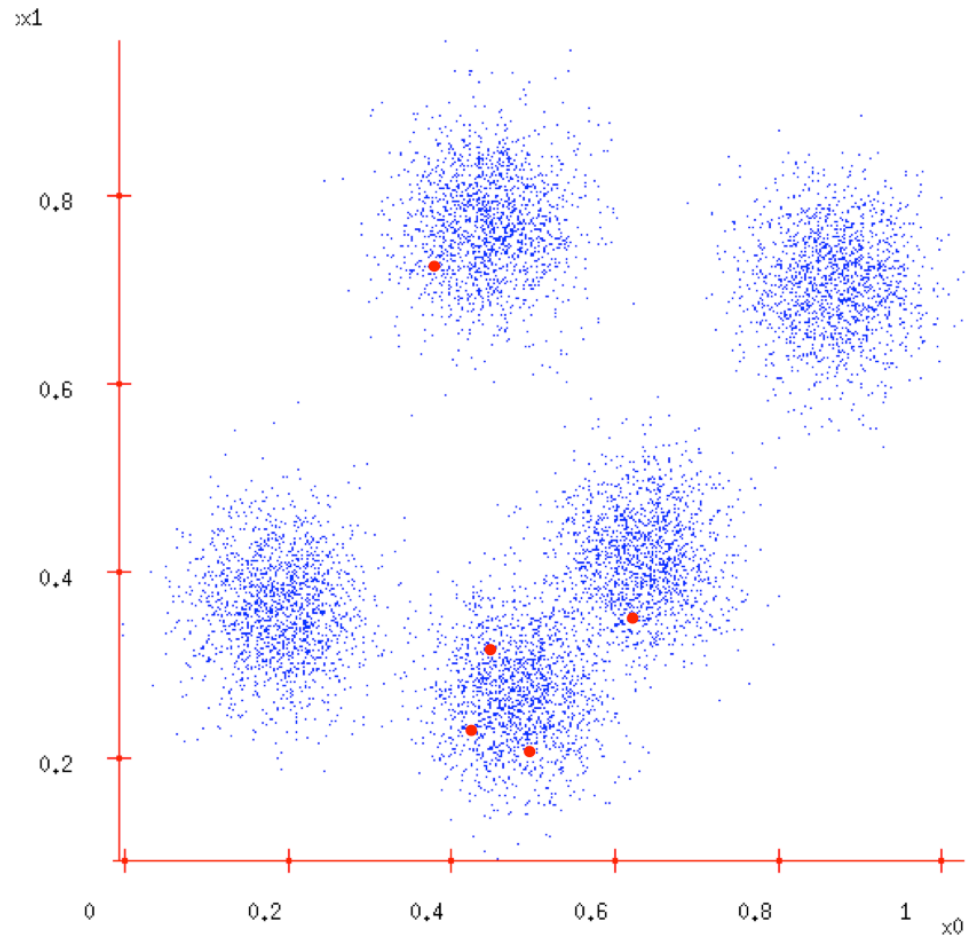This general process is iterated until nothing changes
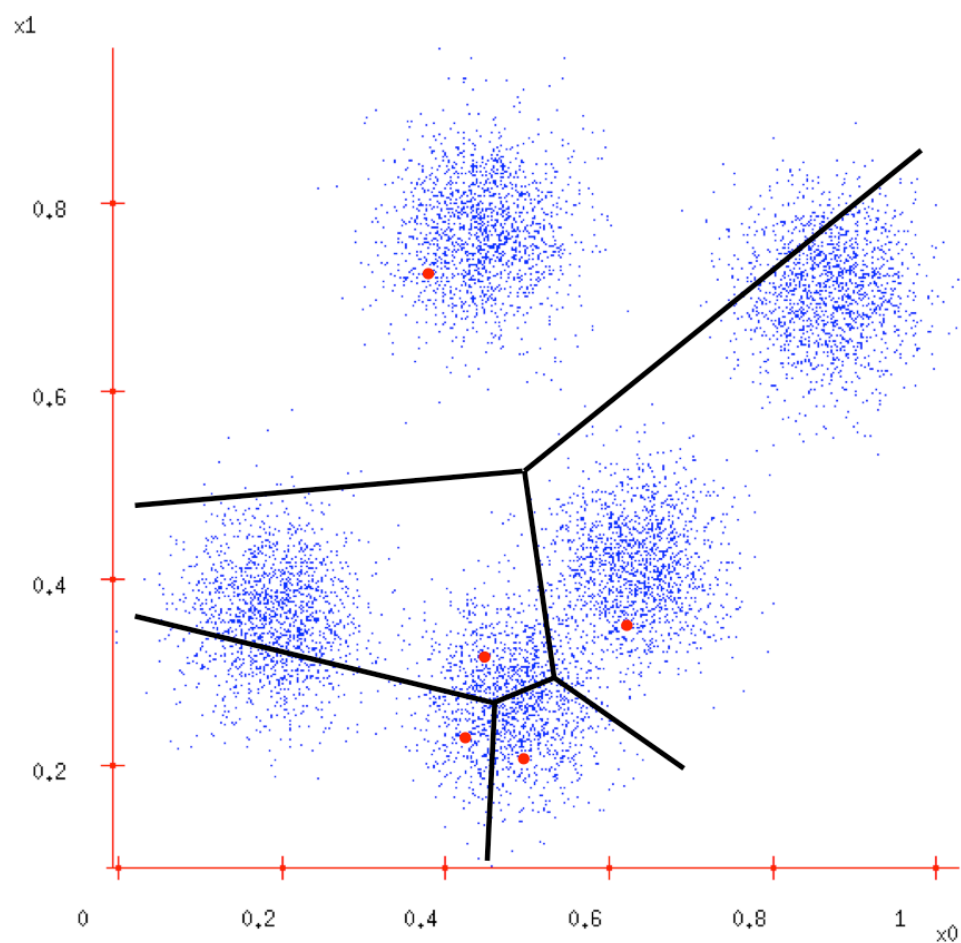
# k-Means Clustering: Example

Example

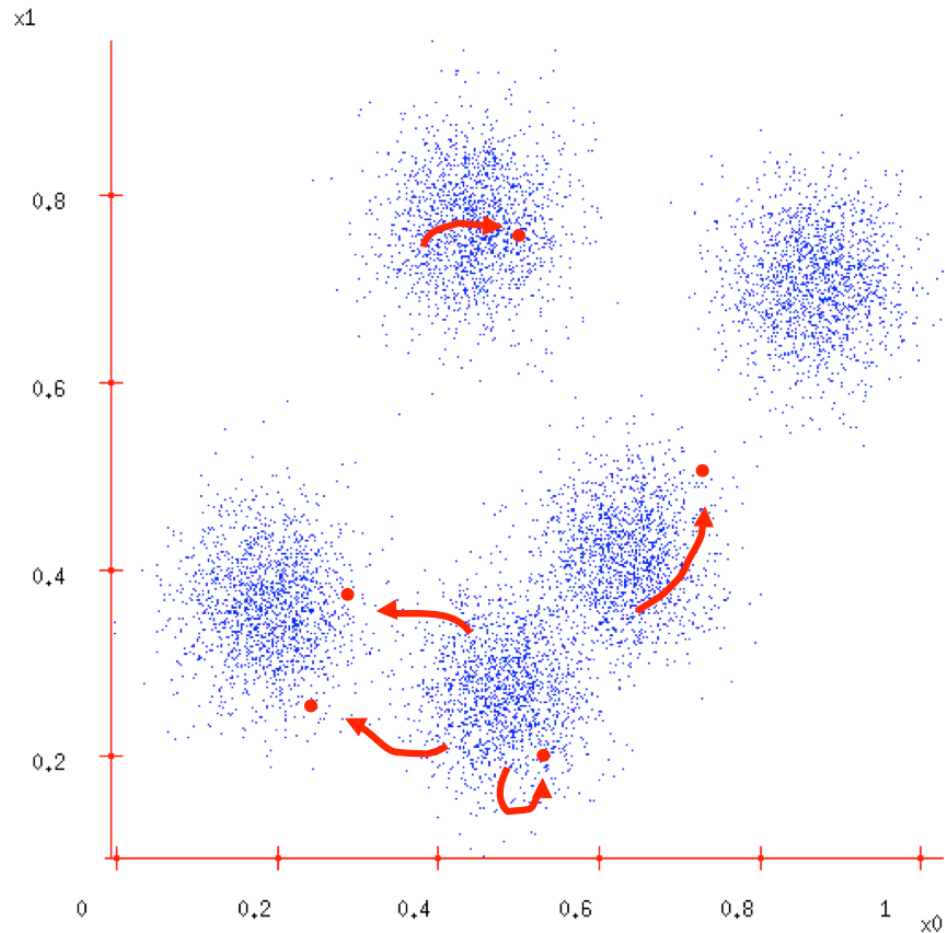# k-Means Clustering: Example



Initialization of Centroids

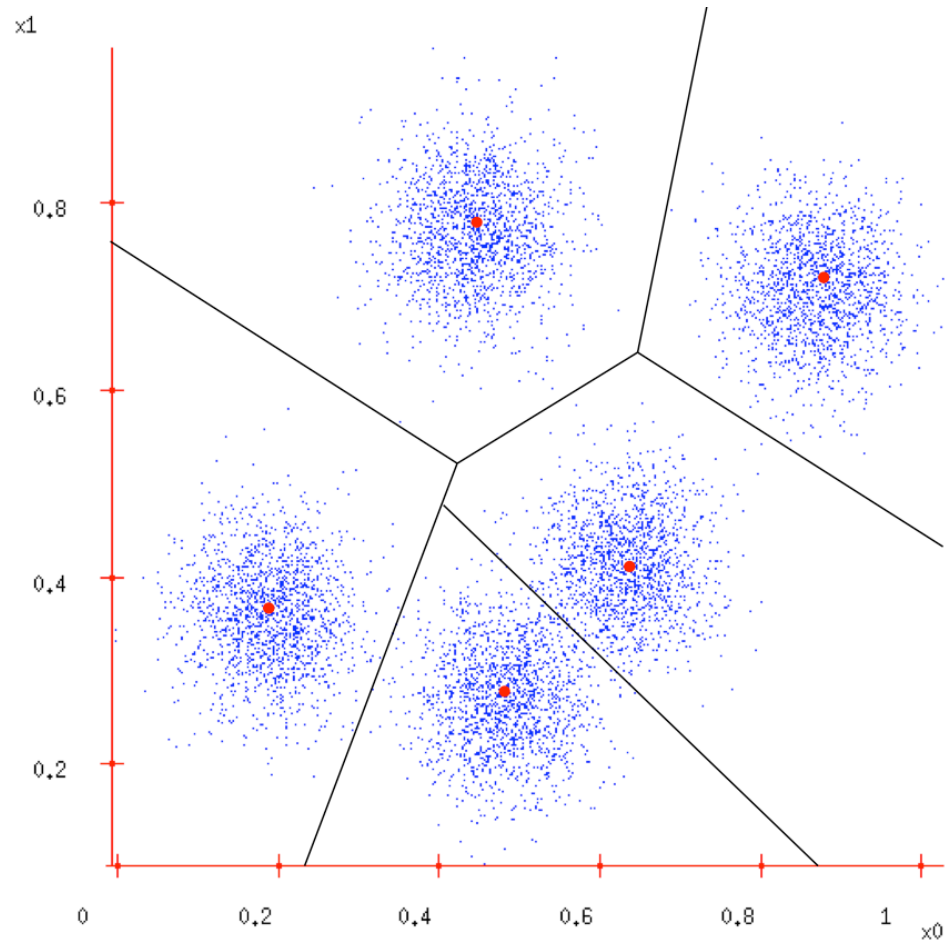# k-Means Clustering: Example

Attribution of Data to Clusters

# k-Means Clustering: Example

Re-Computation of Centroids

# k-Means Clustering: Example

Re-Attribution of Data to Clusters
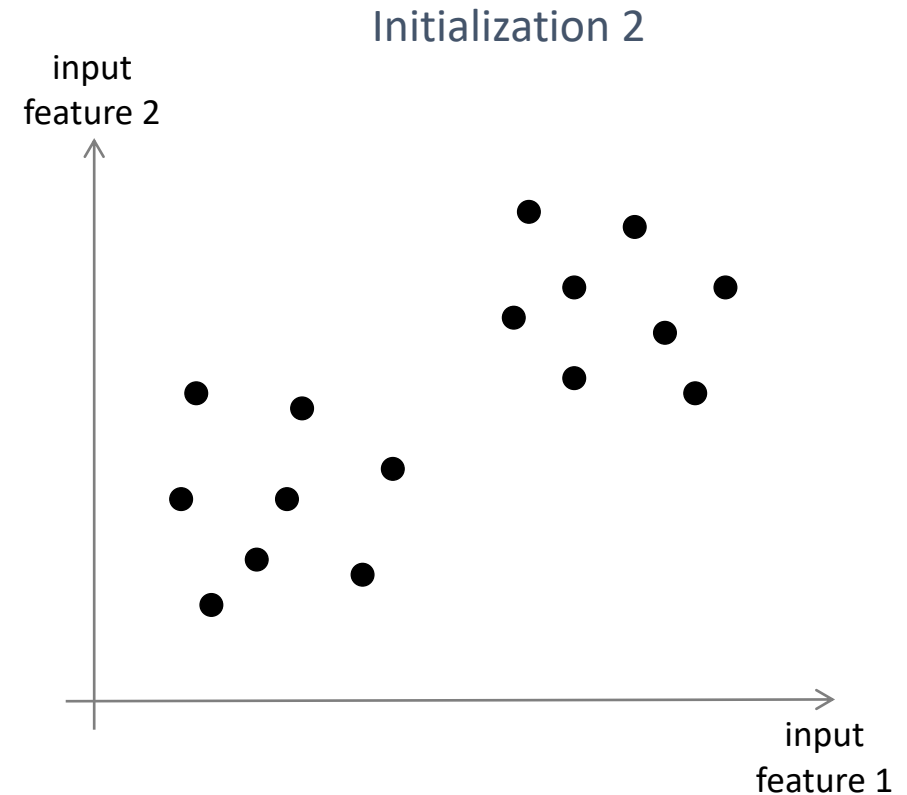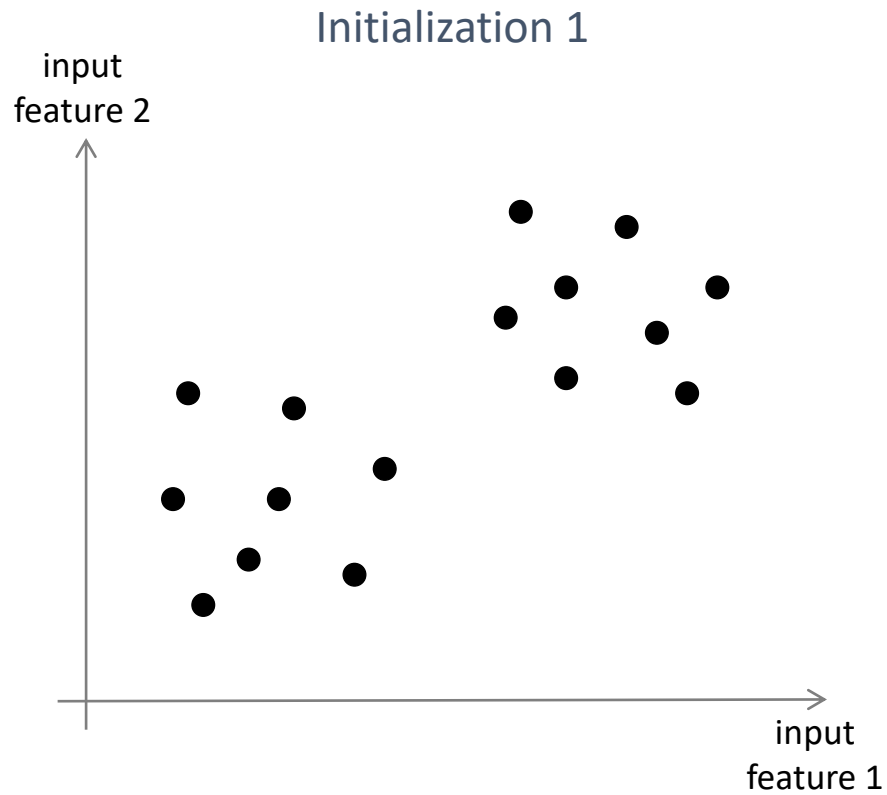
# k-Means Clustering: Considerations

## Does the algorithm terminate?

The k-means clustering algorithm always converges to a solution (cluster assignment), with the cost decreasing from iteration to iteration.

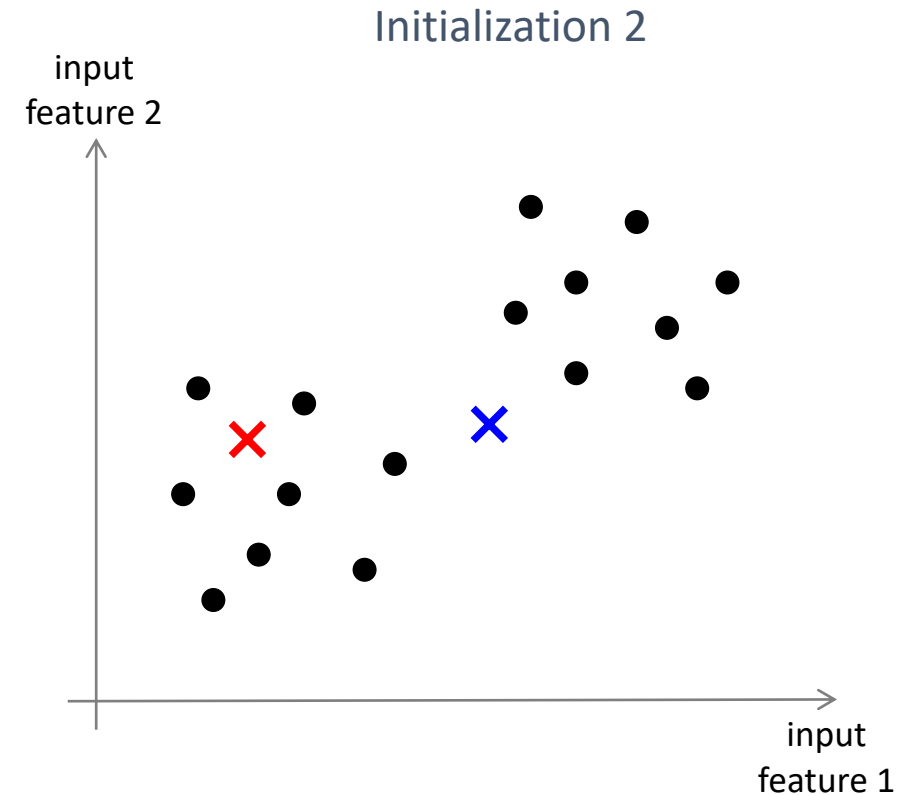However, the resulting cluster assignment may be a local minimum of the cost function *in lieu* of the global one.
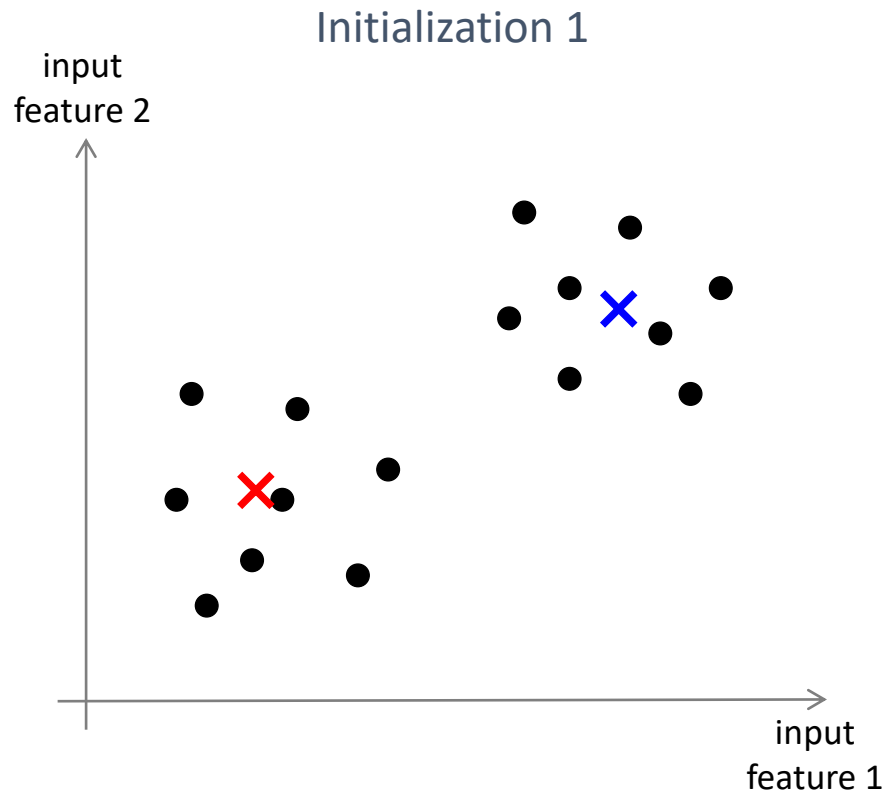
# k-Means Clustering: Considerations

## What is the impact of initialization?

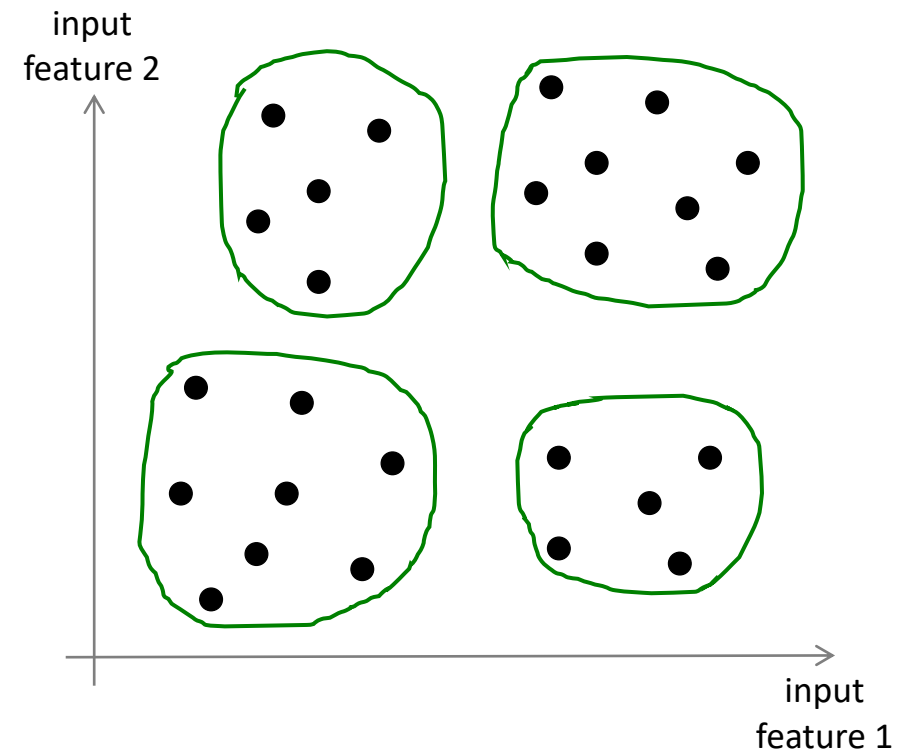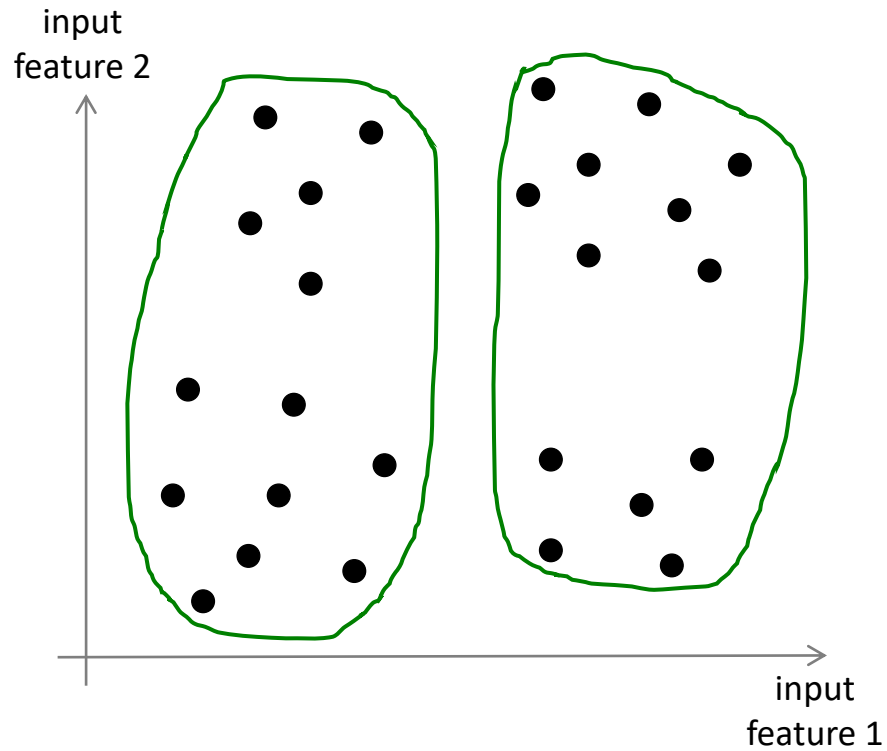# k-Means Clustering: Considerations

What is the impact of initialization?



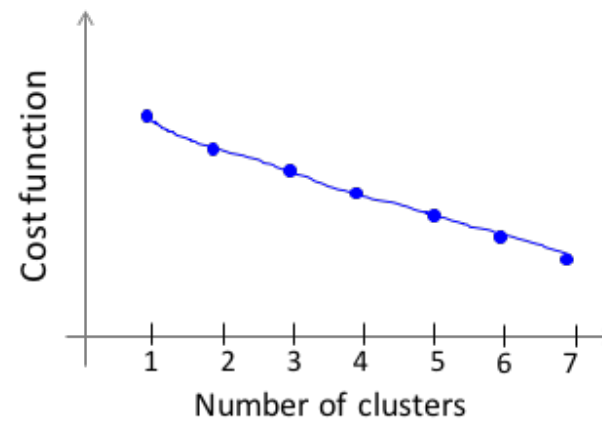Solution: Run the algorithm with multiple initializations, but select the result with minimum cost
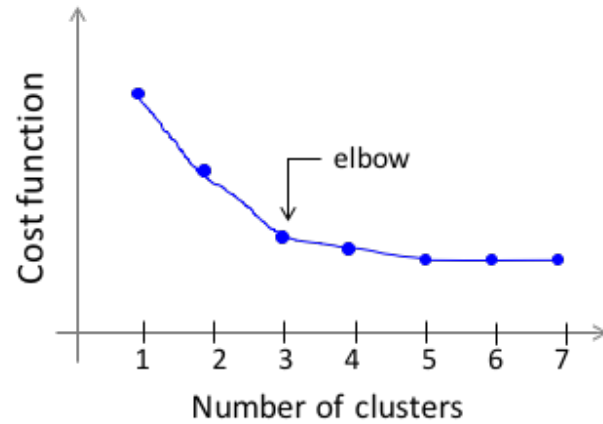
# k-Means Clustering: Considerations

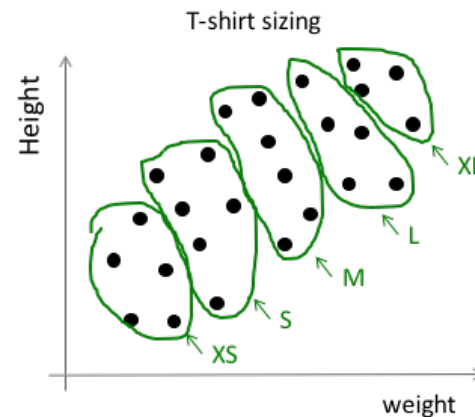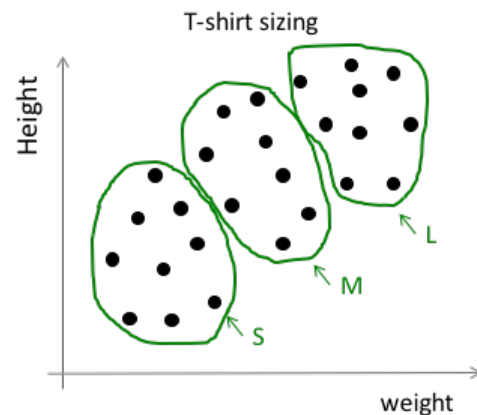How to select the number of clusters?

# k-Means Clustering: Considerations

How to select the number of clusters?



Elbow Method

Application Informed