# Unit 4 - Lecture 14: Clustering 2

## Objectives:

- Understand the limitation of the K-means algorithm
- Understand how K-Medoids algorithm is different from the k-Means algorithm
- Understand the computational complexity of the k-means and the K-Medoids algorithms
- Understand the importance of choosing the right number of clusters
- Understand elements that can be supervised in unsupervised learning

## Limitations of the k-Means Algorithm:

The cost we introduced last time is:

$$\text{Cost}(C_1, \ldots C_k, z^{(1)} \ldots z^{(k)}) = \sum_{j=1}^{k} \overbrace{\sum_{i \in C_j} \| x^{(i)} - z^{(j)} \|^2}^{\text{in each cluster}}$$

go through all the clusters from $1 \to k$

measure the distance between the representative and each point

This K-mean Algorithm works to find the best partitioning, which optimizes this cost.

1. Randomly initialize $z^{(1)} \ldots z^{(k)}$

2. Iterate until no change in cost
   2a. for $i = 1 \ldots m$    $C_j = \{ i \mid s.t. \ z^{(j)} \text{ is closest to } x^{(i)} \}$
   2b. for $j = 1 \ldots k$    $z^{(j)} = \dfrac{\sum x^{(i)}}{C_j}$

The limitations of this algorithm:

- the z's are actually not guaranteed to be members of the original set of points x.
- In order for us to say that the representative is actually a centroid of the points, we had to utilize the squared Euclidean distance. So if we want to use another methode, this algorithm will not cut.

## Introduction to the k-Medoids Algorithm:

1. Randomly initialize $\{ z^{(1)} \ldots z^{(k)} \} \subseteq \{ x^{(1)}, \ldots x^{(m)} \}$

2. Iterate until there is no change in cost
   2.a. for $i = 1 \ldots m$    $C_j = \{ i \mid s.t. \ z^{(j)} \text{ is closest to } x^{(i)} \}$
   2.b. for $j = 1 \ldots k$    $z^{(j)} \in \{ x^1 \ldots x^{(m)} \mid s.t. \sum_{c \in C_j} \text{dist}(x^{(i)}, z^{(j)}) \text{ is minimal.} \}$

# Computation complexity of K-Means and K-Medoids:
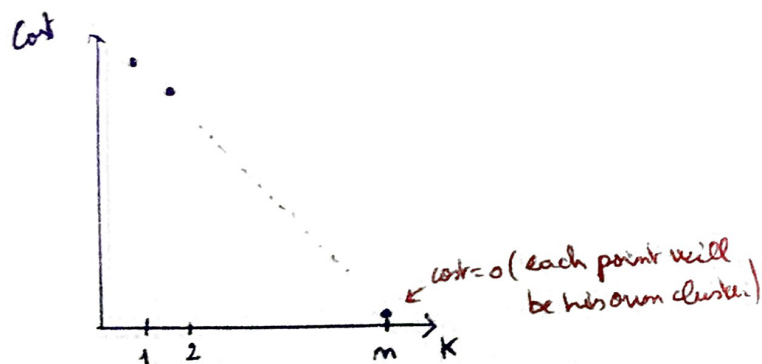
Differences between these two algorithms?

(2.b) in K-Medoids clearly seems to us more expensive than the computations that we are doing in K-Means. If we are going to use $\mathcal{O}$ notation which tell us about the order of growth, let's at these algorithms and compare them in terms of their time complexity. (for one iteration)

$$K\text{-Mean} : \mathcal{O}(\underset{\uparrow}{n} . \underset{\uparrow}{k} . \underset{\uparrow}{d})$$

all points ─ n, size of cluster ─ d, nb of clusters ─ k

$$K\text{-Medoids} : \mathcal{O}(n^2 kd)$$

whenever your are selecting clustering algorithms which fits for your application, you may want to think about different consideration when you're finding the best clustering algorithm for your needs.

## Determining the Number of clusters:



cost = 0 (each point will be his own cluster)

Let's say we have a supervised task

the k will always be related to the performance on the final supervised task.

People always feel that unsupervised means that we don't provide our system with any knowledge However, even though it's unsupervised in the sense that we don't provide any annotated points, we as people who develop these algorithms actually provide quite a bit of indirect supervision.

We decide which similarity measure to use, how many clusters to give.