# Clustering (unsupervised technique)

- grouping, similar to each other
- outlier detection (data cleaning / processing)
- filling gaps in your data

## partitional clustering

- n data, k clusters

$$\sum_{k}^{k} \sum_{x_i, x_j \in C_k} d(x_i, x_j)$$

  $k \to k$ clusters

  $\hookrightarrow$ clustering

- Centroids $\to$ center of the cluster

  euclidean distance

$$\sum_{k}^{k} \sum_{x_i, x_j \in C_k} d(x_i, x_j)^2 = \sum_{k}^{k} |C_k| \sum_{x_i \in C_k} d(x_i, \mu_k)^2$$

  mean of $C_k$

- always converge, but can't find the optimal solution every time, depend on the starting point - too close to each other

- use k-means ++ if possible