# Clustering

**Task:** Find <span style="color:yellow">clusters/subgroups</span> in a dataset $S = \{x_1, ..., x_n\} \in \mathbb{R}^p$.

▷ Samples within subgroup <span style="color:green">similar/homogeneous</span>

▷ Samples in different subgroups <span style="color:green">"distant"/heterogeneous</span> from each other.

**Ex:** • Customers of a company ⟶ grouping for targeted marketing.

• Biology: Find groups of genes.

**Q:** • What notion of <span style="color:green">similarity</span>?

• Precise definition?

**"Unsupervised Learning":**
No subset of data with "correct" classification/grouping available.

**Note:** Fundamentally different from classification problems!

# K-means clustering: [Steinhaus '56, Lloyd '57]

Given $n$ points $x_1, \ldots, x_n \in \mathbb{R}^p$, find $k$ centroids $c_1, \ldots, c_k \in \mathbb{R}^p$ and a partition $\Gamma_1 \cup \Gamma_2 \cup \ldots \cup \Gamma_k = \{1, \ldots, n\}$ ($\Gamma_i \cap \Gamma_j = \phi$ if $i \neq j$) such that

$$F\left(\{c_i\}_{i=1}^{k}, (\Gamma_i)_{i=1}^{k}\right) := \sum_{j=1}^{k} \sum_{i \in \Gamma_j} d(x_i, c_j)$$

"k-means objective ($\star$)"

is <u>minimized</u>, where $d : \mathbb{R}^p \times \mathbb{R}^p \to \mathbb{R}$ is a "distance" function

Ex: 1. $d(x,y) = \|x - y\|_2^2$ "squared Euclidean" $\leftarrow$ by default!

2. $d(x,y) = \|x - y\|_1$. "k-median".

<u>Observation</u>: K-means with is NP-hard for $k \geq 2$. [Drineas et al. '04].

# Lloyd's algorithm (often called "k-means"):

Input: $\{x_i\}_{i=1}^{n} \in \mathbb{R}^p$, desired nr. of clusters $k$

1. Initialize $c_1, \ldots, c_k \in \mathbb{R}^p$.

   Repeat until convergence:

2. $\forall i = 1, \ldots, n$: Assign $x_i \in \Gamma_j$ if $c_j$

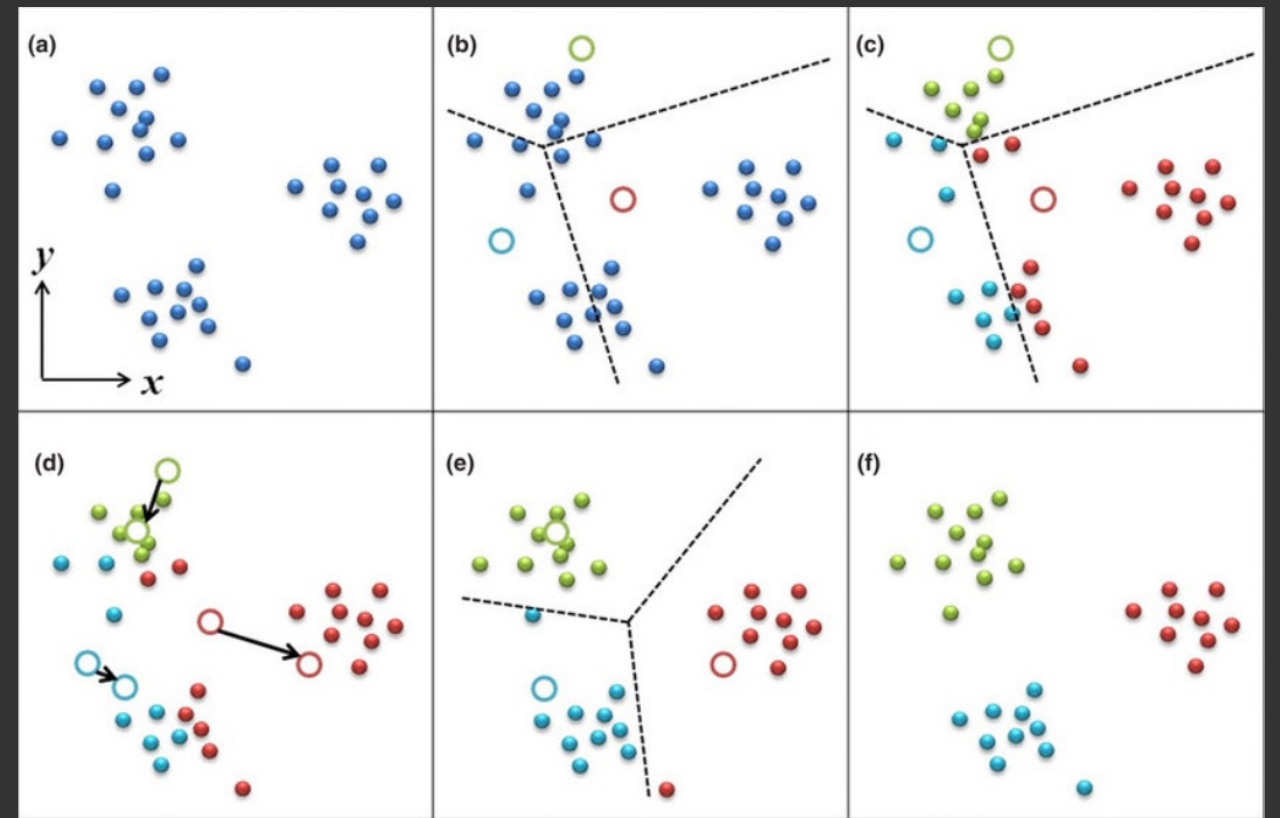   is <u>closest centroid</u> to $x_i$ among $\{c_\ell\}_{\ell=1}^{k}$

3. Update $\forall j = 1, \ldots, k$:

   $$c_j = \underset{c \in \mathbb{R}^k}{\operatorname{argmin}} \sum_{i \in \Gamma_j} d(x_i, c) = \frac{1}{|\Gamma_j|} \sum_{i \in \Gamma_j} x_i$$

   $\hookleftarrow$ if $d(\cdot, \cdot) = \|\cdot - \cdot\|_2^2$

$\nearrow$ a) randomly among the $\{x_i\}_{i=1}^{n}$

or

$\searrow$ b) k-means++ : fancier

**1.** **2.**



(a) $y$ $x$  (b)  (c)
(d)  (e)  (f)

**3.** **2.** **3.**

▷ Finds <u>local</u> optimum of (✳)

▷ Works well for "convex" clusters

To consider :

▷ How to choose <u>nr. of clusters K</u>?

▷ <u>Which distance</u> to choose (geometry of underlying space..)?

▷ <u>Initialization</u> : If prior knowledge available, ⟶ might be better than random.

▷ Needs a lot of pairwise distances. <u>If $n \gg 10^5$ or so, slow</u> ⟶ "Minibatch KMeans"

Other clustering methods :

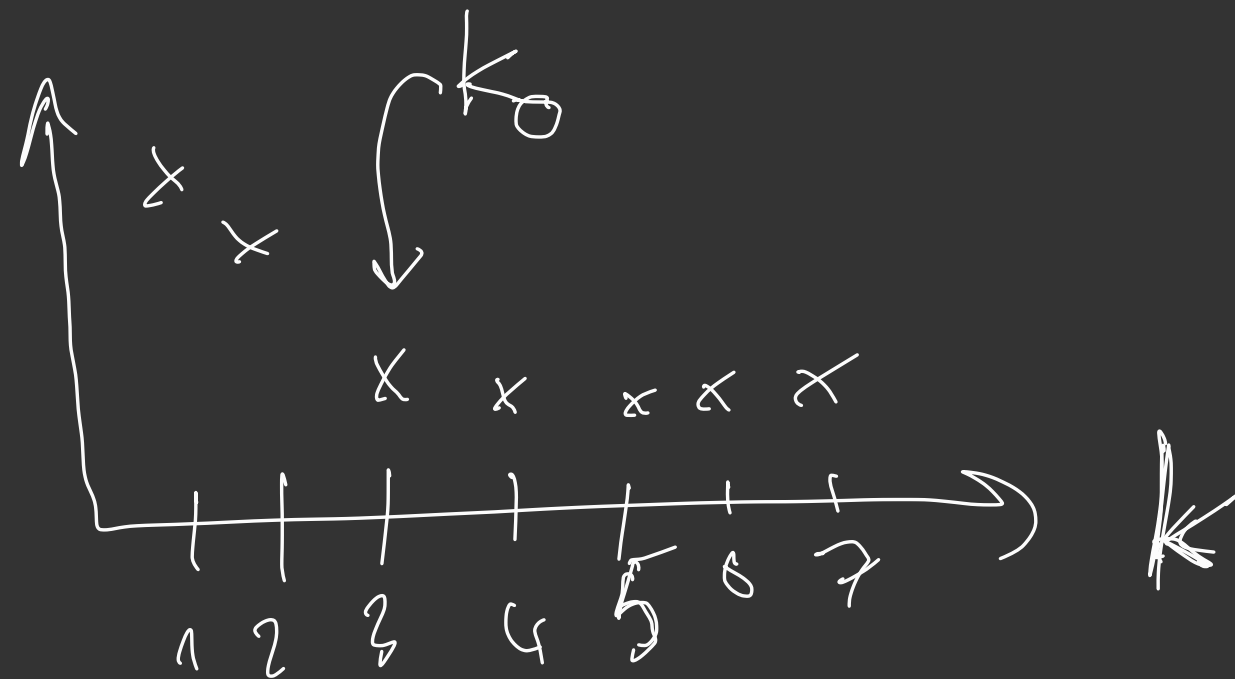– Spectral Clustering : Based on laplacian of similarity graph.

– Hierarchical Clustering : Creates trees

# How to choose parameter "k" in pratice?

One option: Run Lloyd's algorithm for $k = 1, 2, \ldots$ until convergence and then find some $k_0$ s.t. $F(k_0) << F(k_0 - 1)$, but $F(k_0 + \ell) \approx F(k_0)$ for many $\ell = 1, 2, \ldots$,

k-means objective after convergence

$k_0$

$\times \quad \times \times \times \times$

$\times \quad \times$

$\uparrow$

1 2 3 4 5 6 7 $\longrightarrow k$

number of clusters

"Elbow plot"