

# K-means

Caio Corro

## 1 Ré-écriture de la dispersion intra-clusters

Soit  $X \subseteq \mathbb{R}^m$  un ensemble de points et  $\pi = \{C^{(1)}, \dots, C^{(k)}\}$  une partition de  $X$ . Pour simplifier les notations, on notera :

$$\begin{aligned}n &= |X|, \\n_i &= |C^{(i)}|.\end{aligned}$$

Le centroïde du des données  $X$ , que l'on écrit  $\bar{\mathbf{x}}$ , est défini comme :

$$\bar{\mathbf{x}} = \frac{1}{|X|} \sum_{\mathbf{x} \in X} \mathbf{x}.$$

Le centroïde du cluster  $i$ , que l'on écrit  $\mathbf{m}^{(i)}$ , est défini comme :

$$\mathbf{m}^{(i)} = \frac{1}{|C^{(i)}|} \sum_{\mathbf{x} \in C^{(i)}} \mathbf{x}.$$

La dispersion intra-clusters est définie comme :

$$\sum_{i=1}^k \sum_{\mathbf{x} \in C^{(i)}} \|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2.$$

1. Montrer que :

$$\sum_{\mathbf{x}, \mathbf{x}' \in C^{(i)}} (\mathbf{x} - \mathbf{m}^{(i)})^\top (\mathbf{x}' - \mathbf{m}^{(i)}) = 0$$

2. Montrer que :

$$\frac{1}{|C^{(i)}|} \sum_{\mathbf{x}, \mathbf{x}' \in C^{(i)}} \|\mathbf{x} - \mathbf{x}'\|_2^2 = 2 \sum_{\mathbf{x} \in C^{(i)}} \|\mathbf{x} - \mathbf{m}^{(i)}\|_2^2$$

3. Utiliser le résultat précédent pour reformuler la dispersion intra-clusters. Comment pouvez-vous interpréter cette nouvelle formulation ?