

Stage LIESSE
Introduction à l'IA
B. Apprentissage Non Supervisé

Thomas Bonald

Avril 2021

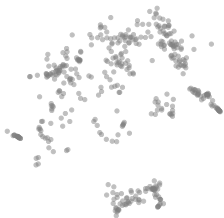


Objectif

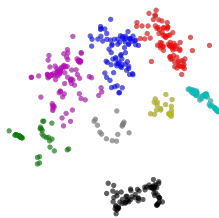
Comprendre la **structure** des données seules
(pas de classe ou de valeur connue).

En pratique, il s'agit de regrouper les données **proches** entre elles
→ attribution d'une **classe**, qu'il faut ensuite **interpréter**.

Données

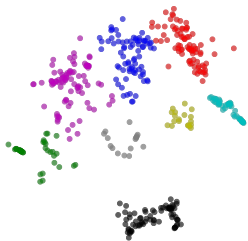


Partition

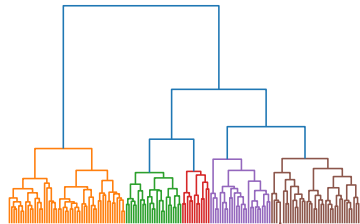


Plan

1. Partitionnement

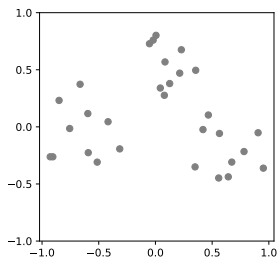


2. Structure hiérarchique



Partitionnement de données

- ▶ **Entrée** : $x_1, \dots, x_n \in \mathbb{R}^d$
 k , nombre de clusters
- ▶ **Sortie** : C_1, \dots, C_k , partition de $\{1, \dots, n\}$



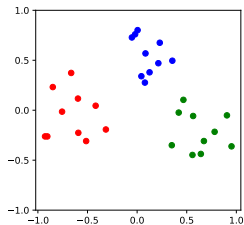
Fonction de coût

Somme des **moments d'inertie** :

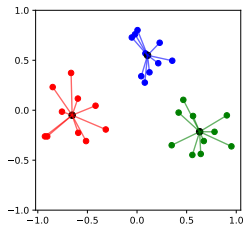
$$\sum_{j=1}^k \sum_{i \in C_j} \|x_i - \mu_j\|^2$$

où $\mu_j = \frac{1}{|C_j|} \sum_{i \in C_j} x_i$ est le **barycentre** du cluster j

Partition



Coût

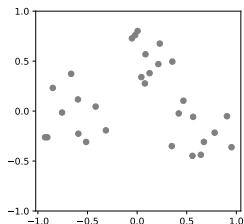


Un problème difficile

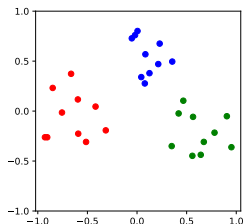
La recherche de la partition optimale est un problème **NP-difficile** :

$$\arg \min_{C_1, \dots, C_k} \sum_{j=1}^k \sum_{i \in C_j} \|x_i - \mu_j\|^2 \quad \text{avec} \quad \mu_j = \frac{1}{|C_j|} \sum_{i \in C_j} x_i$$

Données



Solution



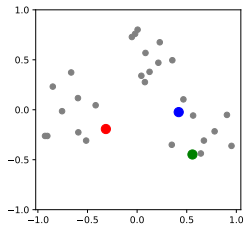
Algorithme des k moyennes

Algorithme **itératif** donnant une solution **approchée**.

Lloyd 1957

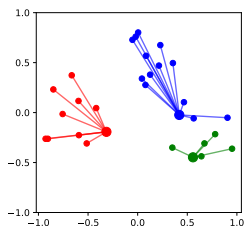
Initialisation

Choix aléatoire des centres



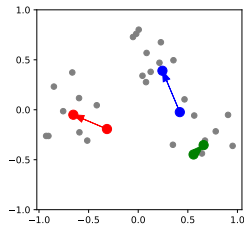
Première partition

Affectation au centre le plus proche



Mise à jour

Nouveaux centres des clusters



Application: Digits

1797 images de chiffres

8×8 pixels

Données

6 1 9 1 1 9 9 8 9 8
5 3 3 9 1 0 7 1 5 6
6 6 4 2 5 9 3 4 0 5
2 8 0 6 1 4 7 5 9 3
5 6 4 3 7 8 5 3 7 0
3 2 8 7 8 4 8 9 5 2
3 7 6 4 9 5 8 7 5 4
6 8 3 8 3 6 7 4 0 6
1 3 4 4 0 2 3 0 7 6
6 7 7 4 5 5 4 1 5 0

Clusters

$k = 10$

8 4 9 1 2 9 9 4 4 9
5 5 9 5 5 3 5 5 5 5
4 4 7 7 8 7 7 7 7 7
4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2
0 0 0 0 0 0 0 0 0 0
8 2 8 1 8 8 7 1 1 8
3 3 3 3 3 3 3 3 3 3
9 9 9 8 5 8 8 8 8 9
6 6 6 6 6 6 6 6 6 6

Apprentissage actif

1797 images de chiffres

8 × 8 pixels

Classification par **le plus proche voisin**

après étiquetage **manuel** de 20 échantillons

Choix aléatoire

Précision = 65%

0 1
2 3
4 5
6 7
8 9
0 4
1 3
4 5
6 7
8 1

Centres de clusters

Précision = 91%

2 5
6 3
1 4
8 0
9 8
2 4
1 5
4 7
7 8
2 6

Compression d'image

Image initiale

251 niveaux de gris



Image compressée

4 niveaux de gris



Image initiale

173 060 couleurs



Image compressée

16 couleurs



Segmentation d'image

Image initiale

251 niveaux de gris



Image compressée

4 niveaux de gris



Masque

1 niveau de gris



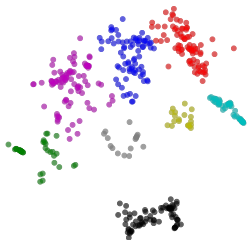
Image segmentée

10 segments

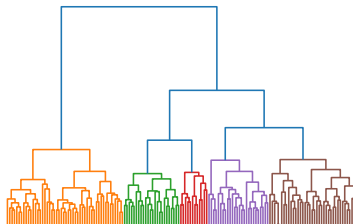


Plan

1. Partitionnement



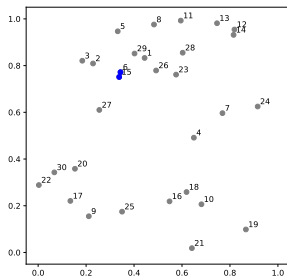
2. Structure hiérarchique



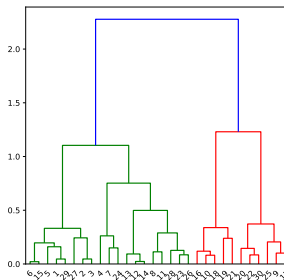
Structure hiérarchique

- ▶ **Entrée** : $x_1, \dots, x_n \in \mathbb{R}^d$
- ▶ **Sortie** : Dendrogramme, arbre binaire de n feuilles

Données



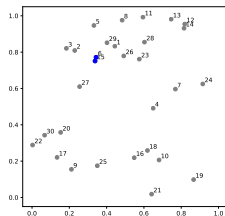
Dendrogramme



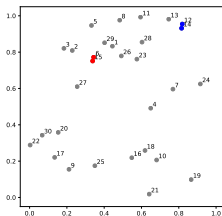
Algorithme agglomératif

Algorithme **glouton** regroupant les clusters les plus proches.

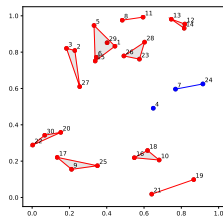
Étape 1



Étape 2



Étape 20

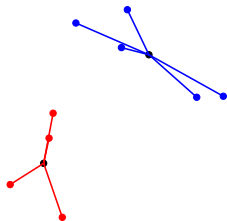


Méthode de Ward

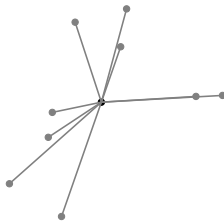
Fusion des clusters minimisant l'augmentation du moment d'inertie

Ward 1963

Clusters séparés



Cluster fusionné



Équivalent à rechercher les clusters C_1, C_2 minimisant :

$$\frac{n_1 n_2}{n_1 + n_2} \|\mu_1 - \mu_2\|^2 \quad \text{avec } n_j = |C_j| \quad \mu_j = \frac{1}{n_j} \sum_{i \in C_j} x_i$$

Application: Digits

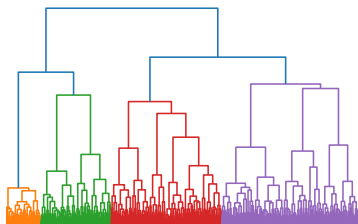
1797 images de chiffres

8×8 pixels

Données

6	4	9	1	1	9	9	8	9	8
5	3	3	9	1	0	7	1	5	6
6	6	4	2	5	9	3	4	0	5
2	8	0	6	2	4	7	5	9	3
5	6	4	3	9	8	5	3	8	0
9	2	8	7	8	4	8	9	5	2
3	7	6	4	9	5	8	7	5	4
6	8	3	8	3	6	7	4	0	6
1	7	4	4	0	2	3	0	7	6
6	4	7	4	5	5	4	1	5	0

Dendrogramme
(méthode de Ward)



Application: Digits

1797 images de chiffres

8×8 pixels

Coupe du dendrogramme (12 clusters)

Centres des clusters

2
7
8
6
0
4
4
9
3
3
1
5
5
4

Dendrogramme

