



HAL
open science

Classification non-supervisée

Thomas A Rieutord

► **To cite this version:**

Thomas A Rieutord. Classification non-supervisée. École d'ingénieur. France. 2021. meteo-02465143v2

HAL Id: meteo-02465143

<https://meteofrance.hal.science/meteo-02465143v2>

Submitted on 10 Dec 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - ShareAlike 4.0 International License

Classification non-supervisée

Thomas Rieutord

Chercheur au CNRM (Météo-France, CNRS – UMR 3589)
Groupe de Météorologie pour l'Assimilation et la Prévision (GMAP)
Equipe RECYF (prévision d'ensemble, prévisibilité...)

Décembre 2021

Outline

1. Introduction
2. Mesures de dissimilarité
3. Classification en partitions
4. Classification hiérarchique
5. Pour aller plus loin

1. Introduction

De quoi parle-t-on ?

Pourquoi faire des classifications non-supervisées ?

Quelques notations

2. Mesures de dissimilarité

3. Classification en partitions

4. Classification hiérarchique

5. Pour aller plus loin

De quoi parle-t-on ?

Classification automatique : attribuer une classe à des individus, définis par des caractéristiques, à l'aide d'un algorithme.

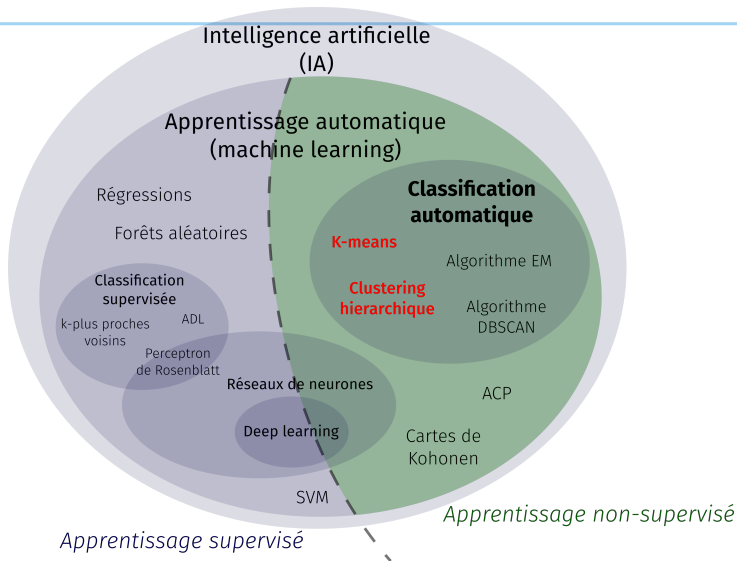
Individu : un élément de notre échantillon. Ils sont représentés par un vecteur.

Caractéristique : quantité commune à tous les individus. Elles sont les composantes des vecteurs qui forment les individus.

Classe : groupement d'individus dont les caractéristiques sont proches.

Algorithme : suite d'opérations explicites permettant de résoudre un problème (ici : regrouper les individus en classes).

De quoi parle-t-on ?



La classification automatique dans le paysage de l'intelligence artificielle et du machine learning.

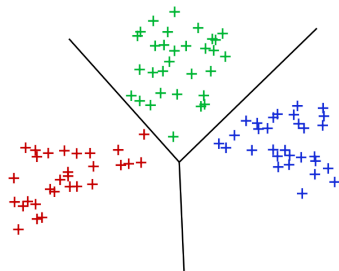
En statistiques, on cherche à expliquer une variable de sortie Y en fonction de variables d'entrée X_1, \dots, X_p . Expliquer cette relation revient à trouver f telle que

$$Y = f(X_1, \dots, X_p)$$

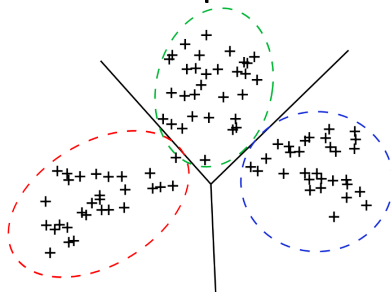
- Lorsque Y est une variable **continue**, on parle de **régression**.
- Lorsque Y est une variable **discrète**, on parle de **classification**.

Supervisé vs non-supervisé

Supervisé



Non-supervisé



Trouver les meilleures frontières
à partir de

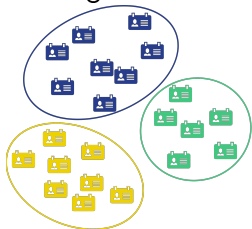
- Position des points
- Couleur des points
(référence)

Trouver les zones de grandes
densité dans

- Position des points
SEULEMENT
- Pas de référence

Pourquoi faire des classifications non-supervisées ?

■ Profilage



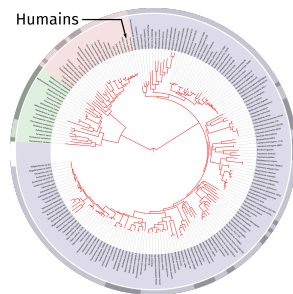
Individus ? personnes physiques.

Caractéristiques ? données socio-économiques (âge, études...) + autres, éventuellement (historique conso, préférences).

■ Taxonomie (classification du vivant)

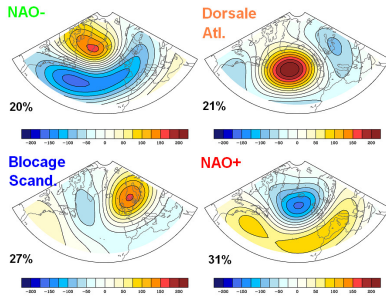
Individus ? espèces vivantes.

Caractéristiques ? gènes.



Pourquoi faire des classifications non-supervisées ?

■ Régimes de temps



© Météo-France

Individus ? carte mensuelle de Z500.

Caractéristiques ? 10 premières composantes principales.

⇒ identification de schémas récurrents qui aident à prévoir le temps à moyenne échéance

De quoi parle-t-on ?

Classification automatique : attribuer une classe à des individus, définis par des caractéristiques, à l'aide d'un algorithme.

Individu (ou observation) : un vecteur $x^i = (x_1^i, \dots, x_p^i)$.

Caractéristique (ou prédicteur) : composante $j \in \llbracket 1, p \rrbracket$ des vecteurs x^i .

Algorithme : retourne une attribution $C(i) = k$.

Classe : ensemble d'objets $C_k = \{i, C(i) = k\}$.

On dispose de N individus ayant tous p caractéristiques :

$$X = \begin{pmatrix} x_1^1 & \dots & x_p^1 \\ \vdots & & \vdots \\ x_1^N & \dots & x_p^N \end{pmatrix} \in \mathbb{R}^{N \times p}$$

Quelques notations

On cherche à grouper N individus en K groupes.

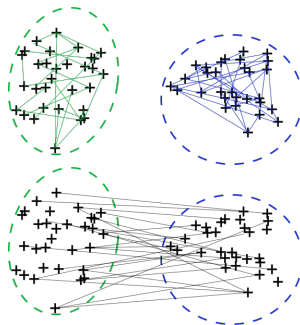
⇒ trouver une fonction d'attribution C (i.e $C(i) = k$) qui minimise une fonction coût $W(C)$.

- Dispersion "within cluster" :

$$W(C) = \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} d(x^i, x^{i'})$$

- Dispersion "between cluster":

$$B(C) = \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i') \neq k} d(x^i, x^{i'})$$



- Dispersion "totale": $T = \sum_{i=1}^N \sum_{i'=1}^N d(x^i, x^{i'}) = W(C) + B(C)$

où $d(x^i, x^{i'})$ mesure la **dissimilarité** entre x^i et $x^{i'}$.

1. Introduction
2. Mesures de dissimilarité
 - Dissimilarité entre caractéristiques
 - Dissimilarité entre individus
3. Classification en partitions
4. Classification hiérarchique
5. Pour aller plus loin

■ Variable quantitative

Exemple: $X_j \in [0, +\infty[$ ou $X_j \in [-1, 1]$...

Toute distance convient. Par exemple, si l est une fonction strictement croissante, on peut prendre

$$d_j(x_j^i, x_j^{i'}) = l(|x_j^i - x_j^{i'}|)$$

La plus répandue est l'écart au carré :

$$d_j(x_j^i, x_j^{i'}) = (x_j^i - x_j^{i'})^2$$

- **Variable ordinale** (catégories ordonnées)

Exemple: $X_j \in \{\text{très mauvais, mauvais, bien, très bien}\} \dots$

Les valeurs prises par X_j sont numérotées de 1 à M . On remplace chaque élément de la colonne j par

$$\frac{k - 1/2}{M}$$

où $k \in \llbracket 1, M \rrbracket$ est le rang de X_j parmi les valeurs prises par X_j . On traite ensuite cette variable transformée comme une variable quantitative¹.

¹NB: Rigoureusement elle ne l'est pas : la variable transformée prend ses valeurs dans $\{\frac{1}{2M}, \dots, 1 - \frac{1}{2M}\}$ qui tend vers $[0, 1]$ quand $M \rightarrow +\infty$

Dissimilarité entre caractéristiques

- **Variable catégorielle** (catégories non-ordonnées)

Exemple: temps sensible à Météo-France $X_j \in \mathcal{E}$ avec $\mathcal{E} = \{\text{pictogrammes pluie, soleil etc.}\} \dots$

Il faut définir une matrice de dissimilarité $(L_{rr'})_{r,r' \in \mathcal{E}}$ telle que $d_j(x_j^i = r, x_j^{i'} = r') = L_{rr'}$

	cl	nu	nu+	pl	nu+	pl	ng	pl	pl+	ng	ng+	av	an	av	an
cl	0	2	6	7	7	8	10	8	10	3	3	6	6		
nu	2	0	3	5	5	7	9	7	9	2	2	4	4		
nu+	6	3	0	2	2	4	7	4	7	4	4	4	4		
pl	7	5	2	0	2	2	4	4	7	4	5	3	4		
nu+	7	5	2	2	0	3	5	2	4	5	4	4	3		
pl	8	7	4	2	3	0	1	2	3	5	6	2	3		
pl+	10	9	7	4	5	1	0	3	2	7	8	4	5		
ng	8	7	4	4	2	2	3	0	1	6	5	3	2		
ng+	10	9	7	7	4	3	2	1	0	8	7	5	4		
av	3	2	4	4	5	5	7	6	8	0	1	3	4		
an	3	2	4	5	4	6	8	5	7	1	0	4	3		
av	6	4	4	3	4	2	4	3	5	3	4	0	1		
an	6	4	4	4	3	3	5	2	4	4	3	1	0		

Le plus souvent, on prend $L_{rr'} = 0$ si $r \neq r'$, 1 sinon.

On dispose de $X = \begin{pmatrix} x_1^1 & \dots & x_p^1 \\ \vdots & & \vdots \\ x_1^N & \dots & x_p^N \end{pmatrix}$

et de $d_j(x_j^i, x_j^{i'})$ pour chaque caractéristique $j \in \llbracket 1, p \rrbracket$.

La dissimilarité entre individus est alors donnée par

$$d_{ii'} = d(x^i, x^{i'}) = \sum_{j=1}^p w_j d_j(x_j^i, x_j^{i'})$$

- Choix des w_j ?
- Problème des unités ?
- Problème des données manquantes ?

Choix des w_j ?

$$\bar{D} = \frac{1}{N^2} \sum_{i=1}^N \sum_{i'=1}^N d(x^i, x^{i'}) = \sum_{j=1}^p w_j \bar{d}_j$$

avec
$$\bar{d}_j = \frac{1}{N^2} \sum_{i=1}^N \sum_{i'=1}^N d_j(x_j^i, x_j^{i'})$$

- La contribution de la j -ème caractéristique à la dissimilarité totale est $w_j \bar{d}_j$.
- Prendre $w_j = 1/\bar{d}_j$ égalise les contributions...



Peut être contre-productif si certaines caractéristiques "séparent mieux" que d'autres.

Choix des w_j ?

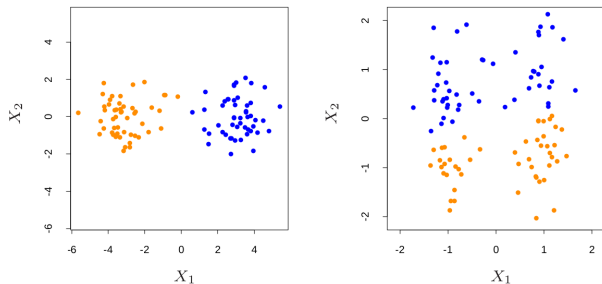


FIGURE 14.5. Simulated data: on the left, K -means clustering (with $K=2$) has been applied to the raw data. The two colors indicate the cluster memberships. On the right, the features were first standardized before clustering. This is equivalent to using feature weights $1/[2 \cdot \text{var}(X_j)]$. The standardization has obscured the two well-separated groups. Note that each plot uses the same units in the horizontal and vertical axes.

©figure de Hastie, Tibshirani, Friedman (2001).

Problème des unités ?

Si X_1 est une P_{mer} en hectopascal (~ 1000) et X_2 est une surcôte en mètres (~ 0.1), et qu'on calcule la dissimilarité directement : $d(x^i, x^{i'}) = \underbrace{(x_1^i - x_1^{i'})^2}_{\sim 100} + \underbrace{(x_2^i - x_2^{i'})^2}_{\sim 0.01}$

X_1 est prépondérant dans la dissimilarité à cause de son unité...

- Les caractéristiques doivent être adimensionnées dans le calcul de la dissimilarité.

⇒ Normalisation...

- "uniforme" : $\tilde{X}_j = \frac{X_j - \min(X_j)}{\max(X_j) - \min(X_j)}$

- "normale" : $\tilde{X}_j = \frac{X_j - \overline{X_j}}{\sigma(X_j)}$

- autre...

- Ce choix a une grande influence sur les résultats.

⇒ Nécessite d'être fait par un utilisateur expert.

Problème des données manquantes ?

En pratique, il arrive que certaines caractéristiques ne soient pas disponibles pour certains individus (panne capteur, erreur, valeur aberrante...).

Soient deux individus x^i et $x^{i'}$.

- Si x^i et $x^{i'}$ ont des caractéristiques $J_{commun} \subset \llbracket 1, p \rrbracket$ en commun

→ on omet les caractéristiques manquantes

$$d_{ii'} = \sum_{j \in J_{commun}} w'_j d_j(x^i, x^{i'}) \text{ avec } \sum_{j \in J_{commun}} w'_j = \sum_{j=1}^p w_j$$

- Si x^i et $x^{i'}$ n'ont aucune caractéristique en commun

→ $d_{ii'}$ ne peut pas être calculée..

- retirer un des deux individus (e.g. celui avec le plus de données manquantes)
- compléter les données manquantes (*data imputation*)

Problème des données manquantes ?

Prenons par exemple $X = \begin{pmatrix} \times & x_2^1 & x_3^1 & x_4^1 & x_5^1 \\ x_1^2 & x_2^2 & x_3^2 & x_4^2 & x_5^2 \\ \times & \times & x_3^3 & x_4^3 & x_5^3 \\ x_1^4 & x_2^4 & \times & \times & \times \end{pmatrix}$

où " \times " est une donnée manquante.

- $d_{12} = \sum_{j=2}^5 w'_j d_j(x^1, x^2)$ avec $\sum_{j=2}^5 w'_j = \sum_{j=1}^5 w_j$
- d_{34} ? \rightarrow retirer x^4 de X (individu avec le plus de données manquantes)

Résumé sur la dissimilarité

- Pour chaque caractéristique, prendre une dissimilarité adaptée à sa nature.
- Adimensionner les caractéristiques (normalisation ou choix des w_j).
- Définir une stratégie pour les données manquantes.



Définir la dissimilarité est une question qui nécessite une **expertise** sur le problème à résoudre.

Outline

1. Introduction
2. Mesures de dissimilarité
3. Classification en partitions
 - L'algorithme des K-means
 - Convergence
 - Initialisation
 - Choix de K
 - Variantes
4. Classification hiérarchique
5. Pour aller plus loin

Rappels

On cherche à grouper N individus en K groupes.

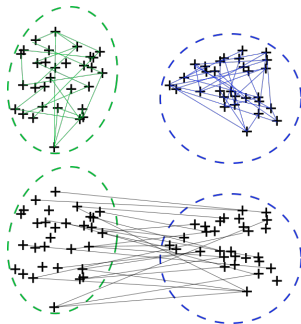
⇒ trouver une fonction d'attribution C (i.e $C(i) = k$) qui minimise une fonction coût $W(C)$.

- Dispersion "within cluster" :

$$W(C) = \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} d(x^i, x^{i'})$$

- Dispersion "between cluster":

$$B(C) = \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i') \neq k} d(x^i, x^{i'})$$



- Dispersion "totale": $T = \sum_{i=1}^N \sum_{i'=1}^N d(x^i, x^{i'}) = W(C) + B(C)$

où $d(x^i, x^{i'})$ mesure la **dissimilarité** entre x^i et $x^{i'}$.

L'algorithme des K-means

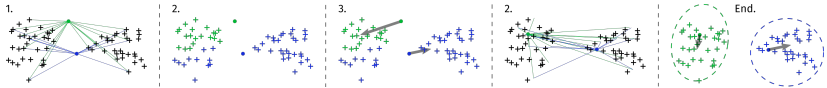
- Le nombre de groupes K est donné.
- Toutes les caractéristiques sont quantitatives.
- La dissimilarité est la distance euclidienne :

$$d(x^i, x^{i'}) = \sum_{j=1}^p (x_j^i - x_j^{i'})^2$$

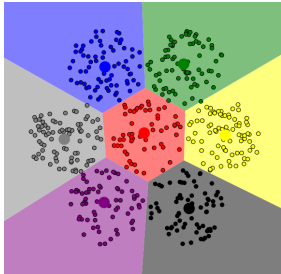
Algorithme des K -means

1. Prendre K individus m_1, \dots, m_K comme centroïdes.
2. Attribuer les individus selon $C(i) = \arg \min_k \{d(x^i, m_k)\}$
3. Redéfinir les centroïdes par $m_k = \frac{\sum_{i=1}^N x^i \mathbf{1}_{C(i)=k}}{\sum_{i=1}^N \mathbf{1}_{C(i)=k}}$
4. Répéter les étapes 2. et 3. tant que $W(C)$ décroît.

Convergence des K -means



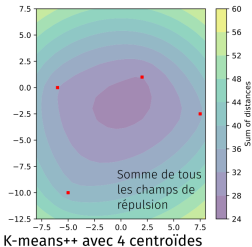
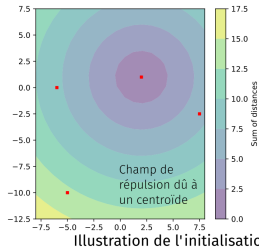
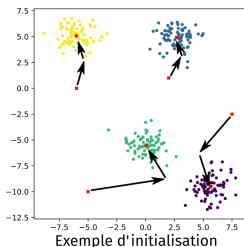
- Assez rapide (quelques dizaines d'itérations suffisent)
- Vers un minimum local de $W(C)$.



Les frontières sont toujours des hyperplans (donc des droites en 2D). Elles forment ce qu'on appelle un **diagramme de Voronoï**.

Initialisation des K -means

- Influence beaucoup le résultat final
- Exemple d'initialisation:
 - aléatoire, répétée plusieurs fois.
 - points extrêmes (stéréotypes caricaturaux)
 - aléatoire guidé (ex: "K-means++")

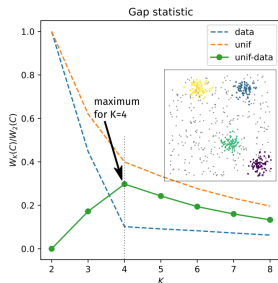
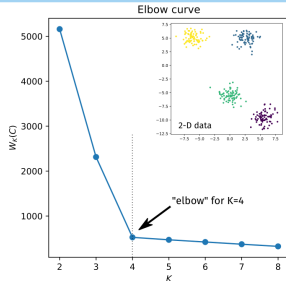


Choix de K

- Peut être imposée par le problème (e.g. 26 lettres dans l'alphabet)
- Estimé automatiquement depuis les données
 - Répéter l'algorithme pour K prenant des valeurs raisonnables et chercher un "coude" dans la courbe de $W_K(C)$.
 - Comparer avec des données suivant une loi uniforme ("gap statistics", Tibshirani et al., 2001).
 - Optimiser un score de classification



Pas de validation croisée possible!



Choix de K : méthode de la silhouette

Soit i un individu affecté à un groupe k . On définit:

- $a(i) = \frac{1}{|C_k|-1} \sum_{i' \in C_k, i' \neq i} d(x^i, x^{i'})$: dissimilarité moyenne entre i et les membres de son groupe.
- $\Delta(i, l) = \frac{1}{|C_l|} \sum_{i' \in C_l} d(x^i, x^{i'})$: dissimilarité moyenne entre i et les membres du groupe l .
- $b(i) = \min_{l \neq k} \Delta(i, l)$: dissimilarité moyenne entre i et les membres du groupe voisin.
- $s(i) = \frac{b(i) - a(i)}{\max\{b(i), a(i)\}}$: **silhouette** du point i . $s(i) \in [-1, 1]$

Quand $s(i) \simeq 1$, $a(i) \ll b(i)$, donc i est bien classé.

Quand $s(i) \simeq -1$, $a(i) \gg b(i)$, donc i est mal classé.

\Rightarrow la silhouette est un score d'évaluation de la classification.

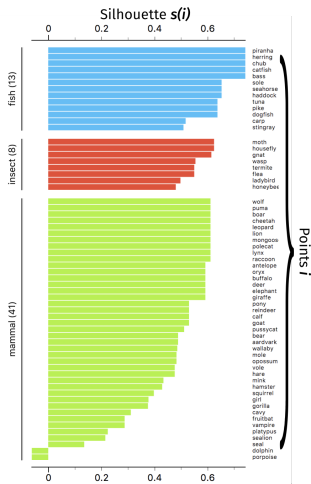
Choix de K : méthode de la silhouette

En pratique, on calcule $s(i)$ pour tous les points i et on trace les $s(i)$ par ordre décroissant sur la verticale.

Les points sont regroupés par cluster.

On voit que les clusters bleu et rouge n'ont que des fortes silhouettes. En revanche, le cluster vert contient des points à silhouette faible.

⇒ Passer de $K = 3$ à $K = 4$?



K -medians

- Toutes les caractéristiques sont quantitatives.
- La dissimilarité est la distance de Manhattan :

$$d(x^i, x^{i'}) = \sum_{j=1}^p |x_j^i - x_j^{i'}|$$

Algorithme des K -medians

1. Prendre K individus m_1, \dots, m_K comme centroïdes.
2. Attribuer les individus selon $C(i) = \arg \min_k \{d(x^i, m_k)\}$
3. Redéfinir les centroïdes par $m_k = \text{median}\{x^i, C(i) = k\}$
4. Répéter les étapes 2. et 3. tant que $W(C)$ décroît.

K -medoids (ou PAM : Partition Around Medoids)

- Caractéristiques sont quelconques.
- La dissimilarité est quelconque.
- Les centroïdes sont des individus de l'échantillon.

Algorithme des K -medoids

1. Prendre K individus m_1, \dots, m_K parmi x^1, \dots, x^N comme centroïdes.
2. Attribuer les individus selon $C(i) = \arg \min_k \{d(x^i, m_k)\}$
3. Redéfinir les centroïdes par $m_k = x^{i_k^*}$ où
$$i_k^* = \arg \min_{i, C(i)=k} \sum_{i', C(i')=k} d(x^i, x^{i'})$$
4. Répéter les étapes 2. et 3. tant que $W(C)$ décroît.

- Il s'agit d'un algorithme très populaire et assez vieux (1967).
- Nombre de groupes K donné a priori.
- **Converge** rapidement, mais vers un minimum local de $W(C)$.
- L'**initialisation** est donc importante, on conseille de répéter l'algorithme plusieurs fois avec différents points de départ.
- Plusieurs astuces existent pour **choisir** K ("coude" dans courbe de $W(C)$, gap statistic...)
- Des **variantes** de l'algorithme permettent de l'adapter à des contraintes particulières (choix de la dissimilarité...)

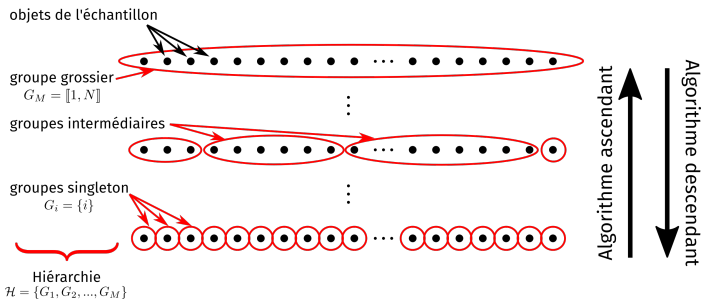
Outline

1. Introduction
2. Mesures de dissimilarité
3. Classification en partitions
4. Classification hiérarchique
 - Classification hiérarchique ascendante
 - Distance cophénétique
 - Dendrogramme
 - Choix de K
5. Pour aller plus loin

Classification hiérarchique ascendante

Une classification hiérarchique donne une **hiérarchie** de groupes (i.e. des groupes imbriqués les uns dans les autres).

- Groupe : ensemble d'indices $G \subset \llbracket 1, N \rrbracket$.
- Hiérarchie : ensemble de groupes $\mathcal{H} \in \mathcal{P}(\llbracket 1, N \rrbracket)$ tel que
 - $\llbracket 1, N \rrbracket \in \mathcal{H}$
 - Si G et H sont dans \mathcal{H} , $G \cap H \neq \emptyset \Rightarrow (G \subset H \text{ ou } H \subset G)$



Algorithme ascendant (ou agglomératif):

- On commence par considérer tous les individus comme des groupes (singletons) puis on réunit les 2 groupes les plus proches en un nouveau groupe qui remplace les 2 précédents. On réunit à nouveau les 2 groupes les plus proches, et ainsi de suite jusqu'à agglomérer tous les points.

Algorithme descendant (ou divisif):

- On commence par considérer le groupe contenant tous les individus et on le divise de façon à maximiser la dissimilarité intra-groupes. Le groupe avec la plus grande variance est ensuite divisé de la même façon et ainsi de suite jusqu'aux singletons.
- Les algorithmes divisifs sont assez rares

Distance cophénétique

Dans un algorithme ascendant, on commence par considérer tous les points comme des groupes puis on réunit les 2 groupes les plus proches en un groupe père, et ainsi de suite.

- ⇒ Par construction, la distance entre les groupes réunis est croissante : on parle de **distance cophénétique**.
- ⇒ On peut représenter les réunions successives par une structure en arbre : on parle de **dendrogramme**.

Clustering hiérarchique nécessite une distance entre groupes (linkage). On attend de cette distance que

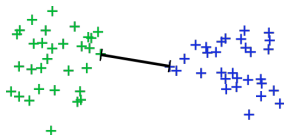
- Les classes qui en résultent soient bonnes :
 - Compacité : classes de petit diamètre
 - Proximité : point d'une classe + proche que ceux de la classe voisine
- Une augmentation du nombre de points lui soit bénéfique.
- Elle ne soit pas affectée par une transformation monotone des distances : utiliser $d_{ii'}$ ou $f(d_{ii'})$ ne doit rien changer si f est monotone.

Soient G et H deux groupes de $\llbracket 1, N \rrbracket$ et $|G|, |H|$ leurs cardinaux respectifs.

$$\text{Saut minimal}^1: D(G, H) = \min_{i \in G, i' \in H} d_{ii'}$$

- Tendence à faire des chaînes de points, les clusters résultants seront étirés, de large diamètre.

⇒ Détection de bandes/strates...



Avantages

- Points d'un même cluster sont proches entre eux
- Invariant par transformation monotone

Inconvénients

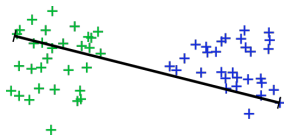
- Clusters peu compacts
- Tend vers 0 quand $N \rightarrow \infty$
- Pas lié à une stat sur la population

¹En anglais: *simple linkage*

Distance cophénétique

$$\text{Saut maximal}^1: D(G, H) = \max_{i \in G, i' \in H} d_{ii'}$$

- Opposé du saut minimal : crée des groupes compacts de petits diamètres mais dont les points au bord peuvent être plus proche d'un autre cluster que d'un point du même cluster.



⇒ Détection de paquets bien isolés

Avantages

- Cluster compacts
- Invariant par transformation monotone

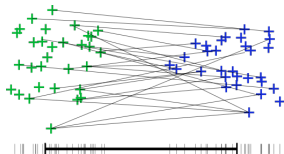
Inconvénients

- Points de clusters différents peuvent être + proche que points d'un même cluster.
- Tend vers $+\infty$ quand $N \rightarrow \infty$
- Pas lié à une stat sur la population

¹En anglais: *complete linkage*

Distance moyenne¹:
$$D(G, H) = \frac{1}{|G||H|} \sum_{i \in G, i' \in H} d_{ii'}$$

- Compromis entre saut minimal et maximal. Elle n'est pas invariante par transformation monotone mais elle est consistante avec la formulation du clustering par mélange de modèle et converge quand $N \rightarrow \infty$.



Avantages

- Compromis compacité/proximité
- Quand $N \rightarrow \infty$, tend vers $\mathbb{E}[d(X, X') | X \in G, X' \in H]$

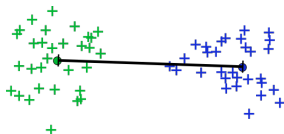
Inconvénients

- Non invariant par transformation monotone

¹En anglais: *group average*

Distance à centroïdes: $D(G, H) = d(\overline{x_G}, \overline{x_H})$
où $\overline{x_G} = \frac{1}{|G|} \sum_{i \in G} x^i$

- Résume un cluster par son point moyen, ce qui a du sens pour les clusters avec des formes assez régulières et de taille équivalente mais qui peut être risqué autrement.



Avantages

- Invariant par transformation monotone
- Stable quand $N \rightarrow \infty$

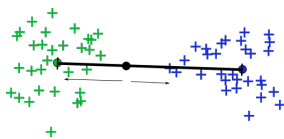
Inconvénients

- Ne tient pas compte de la forme des clusters.
- Petits clusters comptent autant que les gros.

Distance de Ward: $D(G, H) = \frac{w_G w_H}{w_G + w_H} d(\overline{x}_G, \overline{x}_H)$

$$\text{où } \overline{x}_G = \frac{1}{|G|} \sum_{i \in G} x^i \text{ et } w_G = \frac{|G|}{N}$$

- Distance qui maximise l'augmentation d'inertie inter-classe¹. Dans le cas de la distance euclidienne, c'est équivalent à minimiser $W(C)$, ce qui en fait un choix judicieux.



Avantages

- Invariant par transformation monotone
- Stable quand $N \rightarrow \infty$
- Maximise l'augmentation d'inertie inter-classe.

Inconvénients

- Ne tient pas compte de la forme des clusters.
- Inertie n'a de sens que pour la distance euclidienne.

¹Inertie inter-classe : $I_e = \sum_{k=1}^K \frac{|C_k|}{N} d(\overline{x}_{C_k}, \overline{X})^2$ où $\overline{X} = \sum_i x^i / N$

Distance cophénétique : bilan

Nous avons vu 5 distances possibles entre deux groupes G et H :

- Saut minimal: $D(G, H) = \min_{i \in G, i' \in H} d_{ii'}$
- Saut Maximal: $D(G, H) = \max_{i \in G, i' \in H} d_{ii'}$
- Distance moyenne: $D(G, H) = \frac{1}{|G||H|} \sum_{i \in G, i' \in H} d_{ii'}$
- Distance à centroïdes: $D(G, H) = d(\overline{x}_G, \overline{x}_H)$
- Distance de Ward: $D(G, H) = \frac{w_G w_H}{w_G + w_H} d(\overline{x}_G, \overline{x}_H)$

où $\overline{x}_G = \frac{1}{|G|} \sum_{i \in G} x^i$ et $w_G = \frac{|G|}{N}$.

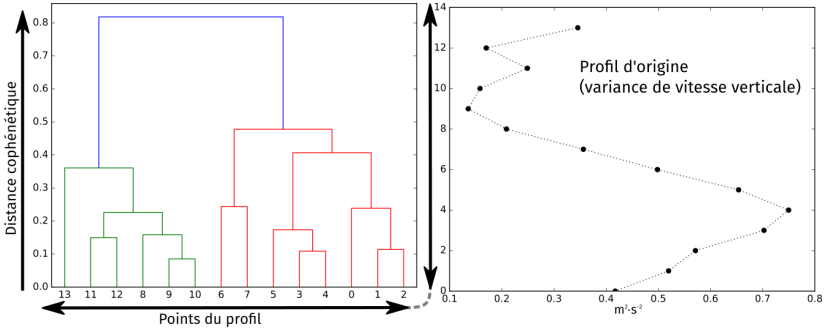
Mais il en existe beaucoup d'autres...

Choix ? Arbitrage avantages/inconvénients par un expert du problème à résoudre.

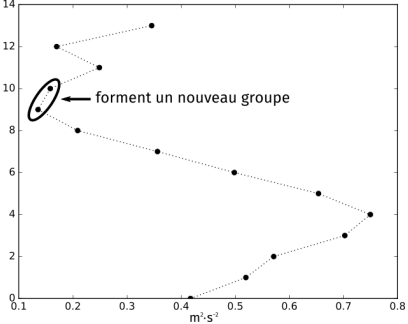
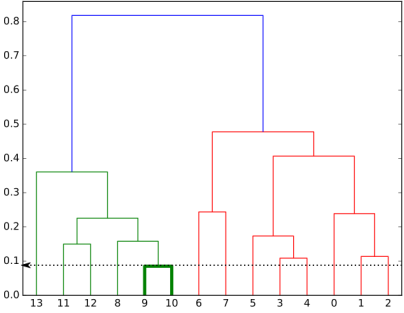
Dans un algorithme ascendant, on commence par considérer tous les individus comme des groupes puis on réunit les 2 groupes les plus proches en un groupe père, et ainsi de suite.

- ⇒ Par construction, la distance entre les groupes réunis est croissante : on parle de **distance cophénétique**.
- ⇒ On peut représenter les réunions successives par une structure en arbre : on parle de **dendrogramme**.

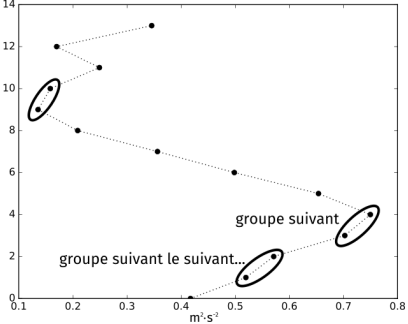
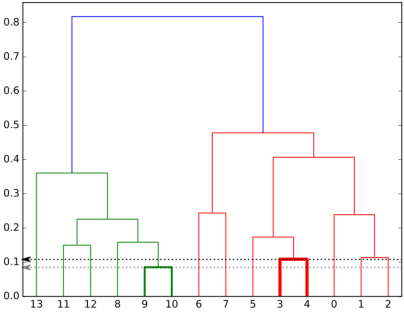
Dendrogramme



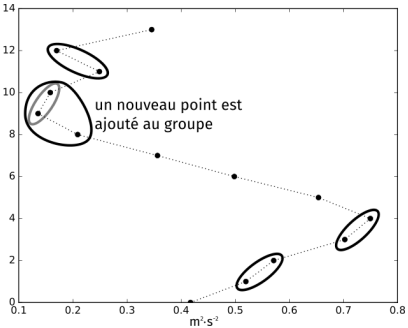
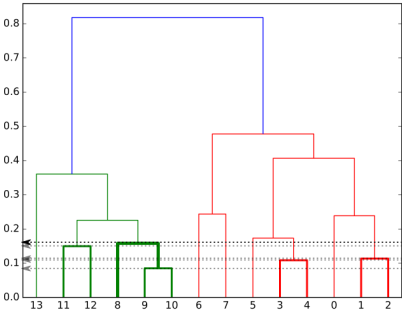
Dendrogramme



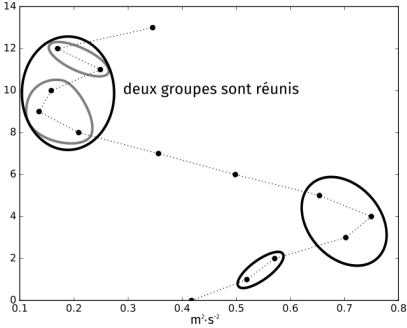
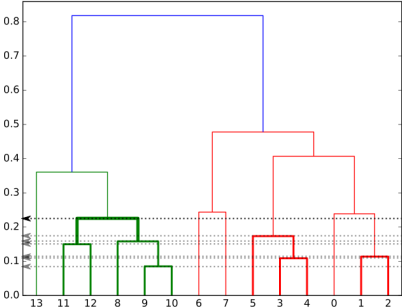
Dendrogramme



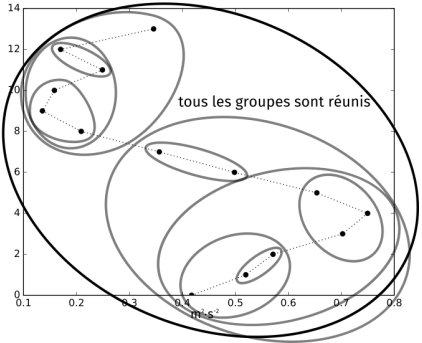
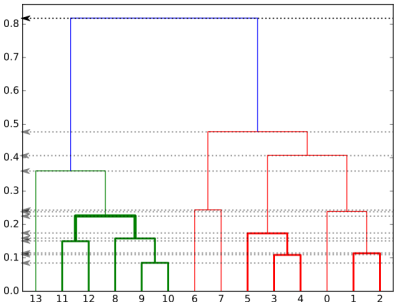
Dendrogramme



Dendrogramme



Dendrogramme



Avantages

- Résumé visuel et complet des données
- Permet d'identifier des groupes à plusieurs échelles
- Les groupes sont imbriqués (cohérence entre les échelles)
- Une fois le nombre de groupes fixés, les groupes sont déjà calculés.

Inconvénients

- Dépend du choix de distance cophénétique.
- C'est le résultat d'un algorithme et non les données brutes.
- Donne toujours une structure hiérarchique, même lorsqu'il n'y en a pas.
- De petites modifications des données peuvent produire un dendrogramme assez différent.

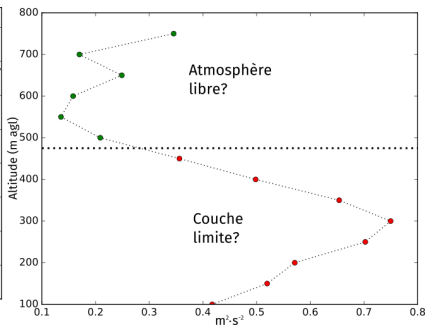
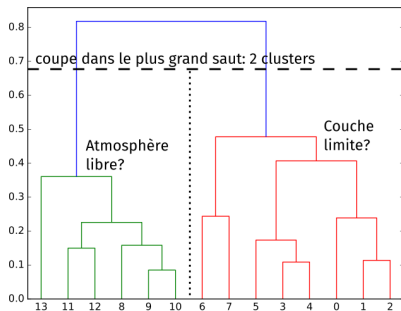
Le clustering hiérarchique ne donne pas une partition mais une hiérarchie, visualisable avec un dendrogramme. Pour obtenir une partition, il faut choisir un niveau dans la hiérarchie.

⇒ **Où couper dans le dendrogramme ?**

- Chercher les liens qui suivent de grands sauts
- De façon automatique : calculer les coefficients d'inconsistance

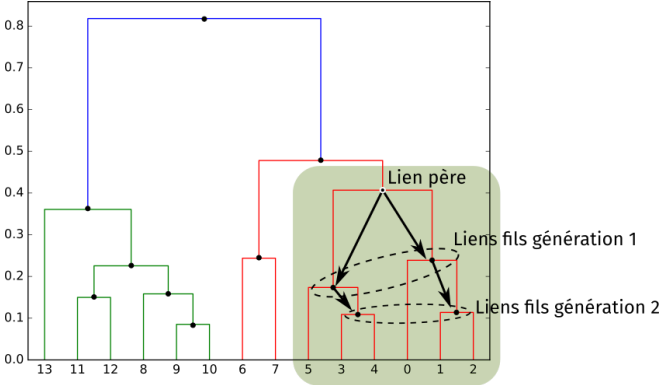
Choix de K : saut le plus grand

Détecter le saut le plus grand (par examen visuel ou en dérivant la distance cophénétique suivant les générations).



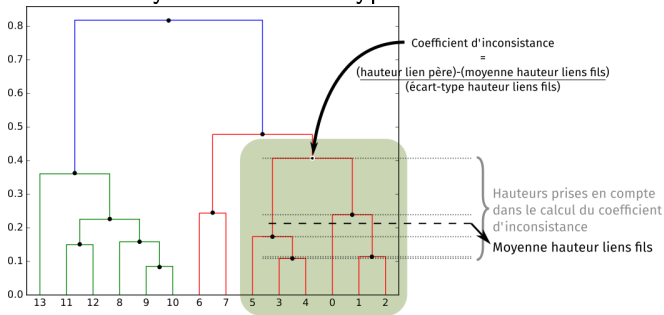
Choix de K : coefficient d'inconsistance

Pour chaque lien, on identifie ses liens fils sur d générations (souvent $d = 2$).



Choix de K : coefficient d'inconsistance

On calcule la moyenne et l'écart type de la hauteur de ces liens.



L'écart normalisé avec la hauteur du lien père est le coefficient d'inconsistance.

On coupe le lien avec le **plus fort coefficient d'inconsistance**.

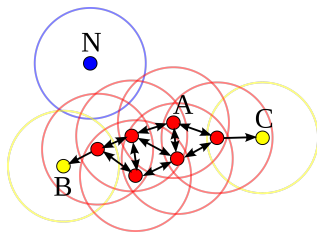
Outline

1. Introduction
2. Mesures de dissimilarité
3. Classification en partitions
4. Classification hiérarchique
5. Pour aller plus loin
 - L'algorithme DBSCAN
 - L'algorithme Expectation-Maximisation
 - Autres scores de classifications
 - Bibliographie

Density-Based Spatial Clustering of Applications with Noise (Ester et al., 1996)

Pseudo-algorithme:

1. Pour tout point x^i , trouver son ϵ -voisinage.
2. Les points avec plus de m voisins appelé *points cœur* (A). Les points cœur connecté entre eux sont mis dans le même cluster.
3. Les points non-cœur sont soit *au bord* (B,C) s'ils sont voisins d'un point cœur, soit *aberrant* (N) sinon.

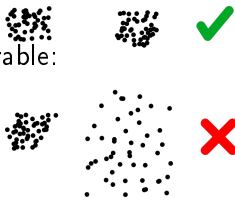


Avantages

- Trouve automatiquement le nombre "naturel" de classes.
- Les groupes peuvent être de forme quelconque (circulaire, entrelacés...).
- Résilient aux données aberrantes et peut même les identifier.

Inconvénients

- Choix des paramètres ϵ et m .
- Les points au bord relié à plus d'un groupe peuvent changer d'affectation suivant l'ordre des points.
- Les groupes doivent être de densité comparable:



L'algorithme Expectation-Maximisation

- Le nombre de groupes K est donné.
- On suppose que chaque groupe est distribué selon une loi gaussienne $\mathcal{N}(\mu_k, \Sigma_k)$.
- On cherche alors à estimer les paramètres μ_k et Σ_k ainsi que la responsabilité de chaque gaussienne.

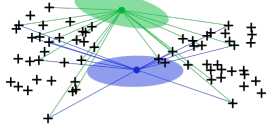
Algorithme EM

1. Initialiser $(\hat{\mu}_k, \hat{\sigma}_k, \hat{\pi}_k)_{k \in \llbracket 1, K \rrbracket}$ (aléatoirement).
2. Expectation: attribuer une responsabilité à chaque gaussienne pour la position des points
3. Maximisation: recalculer les paramètres qui maximisent la vraisemblance
4. Répéter les étapes 2. et 3. tant que $W(C)$ décroît.

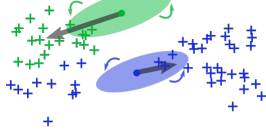
L'algorithme EM: lien avec les K-means ?

Expectation-Maximisation

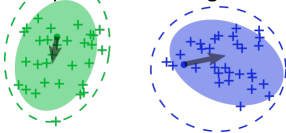
Attribute a responsibility for each Gaussian



Update Gaussians' parameters

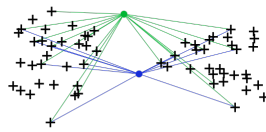


Repeat until convergence

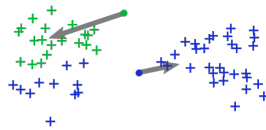


K-means

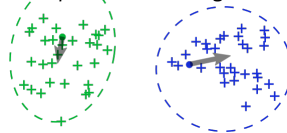
Attribute to closest centroid



Update centroids



Repeat until convergence



Se base sur

- les points moyens de chaque groupe $\mu_k = \frac{1}{|I_k|} \sum_{i \in I_k} x^i$
- la distance moyenne entre un point et le centre de son groupe $\bar{\delta}_k = \frac{1}{|I_k|} \sum_{i \in I_k} d(x^i, \mu_k)$

$$S_{DB} = \frac{1}{K} \sum_{k=1}^K \max_{k' \neq k} \left(\frac{\bar{\delta}_k + \bar{\delta}_{k'}}{d(\mu_k, \mu_{k'})} \right)$$

Interprétation:

- S_{DB} petit: bonne classification
- S_{DB} grand: mauvaise classification
- $S_{DB} = 0$: cas dégénéré où les points d'un groupes sont confondus avec son centre

Se base sur

- les points moyens de chaque groupe $\mu_k = \frac{1}{|I_k|} \sum_{i \in I_k} x^i$
- les diamètres de chaque groupe $\Delta_k = \max_{i, i' \in I_k} d(x^i, x^{i'})$

$$S_D = \frac{\min_{1 \leq k < k' \leq K} d(\mu_k, \mu_{k'})}{\max_{1 \leq k \leq K} \Delta_k}$$

Interprétation:

- S_D petit: mauvaise classification
- S_D grand: bonne classification
- $S_D = 0$: lorsqu'au moins deux groupes ont le même centre (clusters circulaires concentriques)

Se base sur

- la variance inter-groupes $B = \sum_{k=1}^K |I_k| \|\mu_k - \mu\|^2$
- les variances intra-groupe $W_k = \frac{1}{|I_k|} \sum_{i \in I_k} \|x^i - \mu_k\|^2$

$$S_{CH} = \frac{(N - K)B}{(K - 1) \sum_{k=1}^K W_k}$$

Interprétation:

- S_{CH} petit: mauvaise classification
- S_{CH} grand: bonne classification
- $S_{CH} = 0$: cas dégénéré où tous les points sont confondus
- Croît linéairement avec N : ordre de grandeur très variable suivant données

Sources:

- Friedman, Hastie, & Tibshirani (2001). *The elements of statistical learning*. Springer Edition. 2nd Edition. Chapitre 14.
- Saporta (2006). *Probabilités, analyse des données et statistique*.
- Jain, Murty, & Flynn (1999). *Data clustering: a review*. ACM computing surveys (CSUR).
- Rousseeuw (1987). *Silhouettes: a graphical aid to the interpretation and validation of cluster analysis*. Journal of computational and applied mathematics.
- Desgraupes (2017). *Clustering indices*. Comprehensive R Archive Network.

Références:

- (Ester et al., 1996) : Ester, Kriegel, Sander & Xu (1996, August). *A density-based algorithm for discovering clusters in large spatial databases with noise*.
- (Tibshirani, 2001) Tibshirani, Walther, & Hastie (2001). *Estimating the number of clusters in a data set via the gap statistic*. Journal of the Royal Statistical Society: Series B (Statistical Methodology).

Merci pour votre attention
Questions ?

INP-ENM & Météo-France

thomas.rieutord@meteo.fr



www.enm.meteo.fr

www.umr-cnrm.fr