

INFORMATIQUE ORIENTATION LOGICIELS

CLASSIFICATION AUTOMATIQUE

Prof. É. D. Taillard

CLASSIFICATION AUTOMATIQUE

But :

Étant donné un ensemble d'éléments, caractérisés par un certain nombre de mesures, on cherche à trouver des sous-ensembles qui soient *homogènes* et *bien séparés*. Les éléments appartenant à un sous-ensemble doivent donc se ressembler et deux éléments appartenant à des sous-ensembles distincts doivent être différents.

Types de classification :

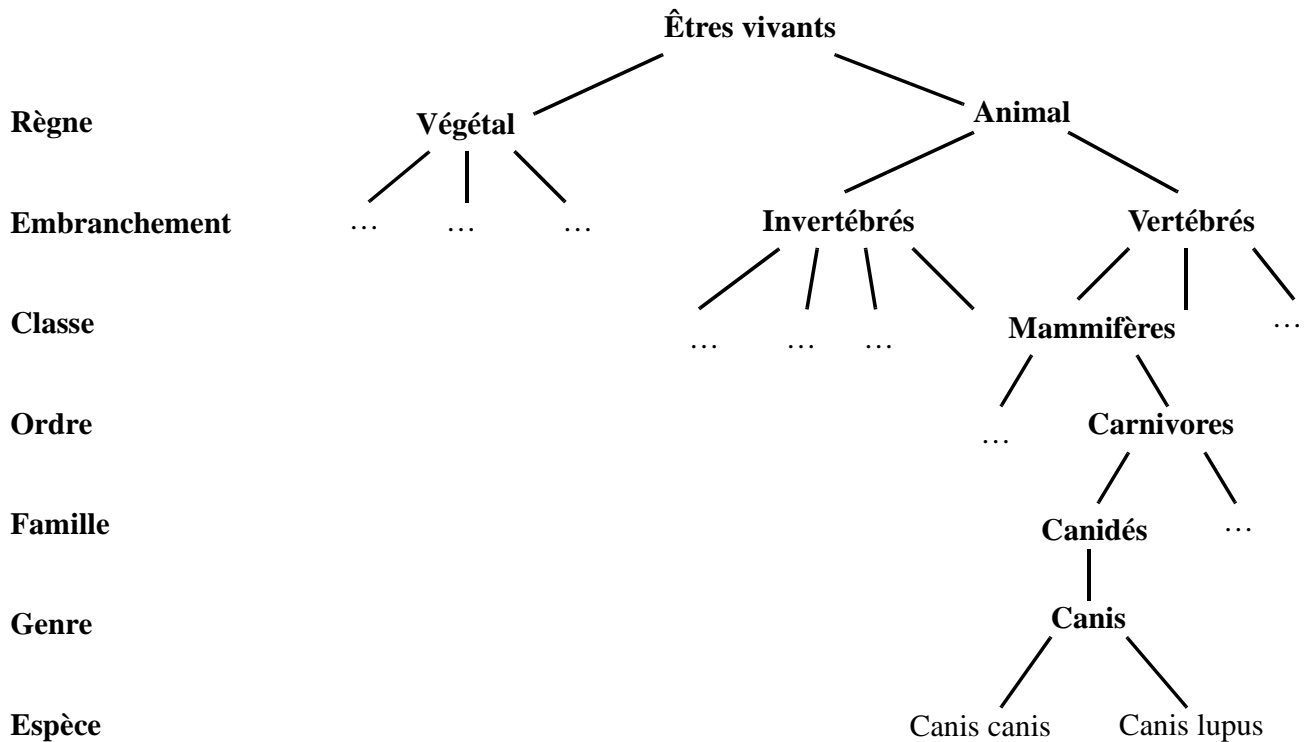
Classification hiérarchique

Les groupes que l'on crée sont emboîtés les uns dans les autres de façon hiérarchique. Exemple : classification des êtres vivants.

Classification non hiérarchique

En classification non hiérarchique, on ne cherche pas à structurer de manière hiérarchique les sous-ensembles entre-eux.

CLASSIFICATION HIÉRARCHIQUE



Méthodes divisives ou agglomératives

En classification *hiérarchique*, on peut procéder par une méthode soit *descendante*, soit *ascendante*.

Dans une méthode *descendante* ou *divisive*, on part de l'ensemble de tous les éléments que l'on fractionne en un certain nombre de sous-ensembles. Ces derniers sont eux-même fractionnés récursivement jusqu'à ce que l'on arrive aux éléments individuels (Règne → Embranchement → Classe → ... → Genre → Espèce)

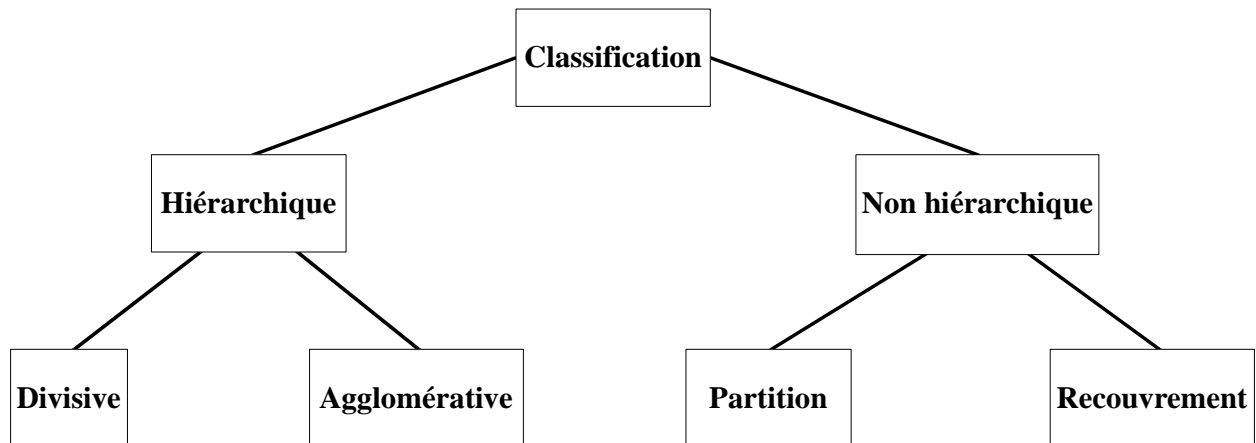
Au contraire, dans un méthode *ascendante* ou *agglomérative*, on part des éléments individuels que l'on regroupe en sous-ensembles avant d'appliquer récursivement les regroupements sur les sous-ensembles jusqu'à obtenir un seul ensemble contenant tous les éléments :

(Espèce → Genre → ... → Embranchement → Règne)

Méthodes recouvrantes et partition

En classification *non hiérarchique*, on peut considérer que chaque élément ne fait partie que d'*un* sous-ensemble, on parle alors de *partition*. On peut aussi considérer que chaque élément en fait partie de *plusieurs*, en attribuant une probabilité d'appartenance à chaque groupe et on parle alors de classification *recouvrante*.

CLASSIFICATION DES MÉTHODES DE CLASSIFICATION



TYPES DE MESURES

Chaque élément est caractérisé par un *ensemble de mesures*. Ces mesures peuvent être de différents types :

Données numériques

Valeur absolue (longueur, température)

Valeur relative (%)

Données nominales

Couleur des yeux (noir, vert, brun, ...).

Pas de comparaisons possibles autres que $A = B$ ou $A \neq B$

Données ordinales

Grade à l'armée ; Petit, moyen, grand. *On peut définir un ordre $A < B < C \dots$*

Données binaires

Cas particulier de variables nominales ou ordinales lorsqu'on a deux valeurs possibles

Sexe, présence ou absence d'une caractéristique

Données conditionnelles, liées à la présence d'autres mesures

Données manquantes

REPRÉSENTATION DES DONNÉES

Les données *brutes* peuvent être représentées par une matrice X à n lignes et p colonnes. Chaque ligne correspond à l'un des n éléments qu'on désire classer, caractérisé par p mesures.

Matrice de données brutes :

		Mesures					
		1	2	...	j	...	p
Éléments	1	x_{11}	x_{12}		x_{1j}		x_{1p}
	2	x_{21}	x_{22}		x_{2j}		x_{2p}
	...						
	i	x_{i1}	x_{i2}		x_{ij}		x_{ip}
	...						
	n	x_{n1}	x_{n2}		x_{nj}		x_{np}

Les colonnes ne sont pas forcément du même type et il peut y avoir des *entrées vides* correspondant aux entrées *manquantes*.

ÉTAPES DE LA CLASSIFICATION

- 1) *Sélection* d'un échantillon de n éléments parmi lesquels on cherche des groupes
- 2) *Mesure* des éléments. Effectuer p mesures sur chaque élément pour obtenir les *données brutes*.
- 3) *Construction* d'une matrice de *dissimilarités* D , de taille $n \times n$. L'élément d_{rs} mesure la dissimilarité qui existe entre les éléments r et s . Tous les d_{rs} sont du *même type*, généralement un nombre réel avec $d_{rr} = 0$ et $d_{rs} \geq 0$.

Parfois les dissimilarités peuvent être booléennes : $d_{rs} = \begin{cases} 0 & \text{Les éléments } r \text{ et } s \text{ sont similaires} \\ 1 & \text{Les éléments } r \text{ et } s \text{ sont différents} \end{cases}$

- 4) *Choix d'un type de classification* (hiérarchique, partition, ...) avec éventuellement des contraintes additionnelles : nombre de groupes à former, nombre maximum d'éléments dans un groupe, ...
- 5) *Choix d'un critère* pour mesurer l'*homogénéité* ou l'*hétérogénéité* des groupes : attribuer une *mesure de qualité* pour chaque classification possible.
- 6) *Choix* ou conception d'un *algorithme* pour la résolution du problème ainsi défini.

CALCUL DES DISSIMILARITÉS

Selon le type des données, la manière de construire la matrice des dissimilarités varie.

Données binaires

Exemple : 0 1 0 0 1 0 0 1 1 objet *r*

1 1 0 1 0 1 1 0 1 objet *s*

Variable n° 1 2 3 4 5 6 7 8 9

Sommaire des mesures entre *r* et *s*.

		Objet <i>r</i>	
		0	1
Objet	0	1 = <i>a</i>	2 = <i>b</i>
	1	4 = <i>c</i>	2 = <i>d</i>

Quelques possibilités de définition de la dissimilarité :

$$\text{Distance de Hamming : } d_{rs} = \frac{b + c}{a + b + c + d}$$

$$\text{Coefficient de Jaccard : } d_{rs} = \frac{b + c}{a + b + c}. \quad \text{Propriété : } c' \text{ est une distance } (d_{rt} \leq d_{rs} + d_{st} \forall r, s, t)$$

$$\text{Mesures liées aux probabilités : } 1 - d_{rs} = \frac{ad - bc}{(a + b)(a + b + c + d)}, \quad 1 - d_{rs} = \frac{ad - bc}{(a + b)(a + c)}$$

$$\text{Approches géométriques : } 1 - d_{rs} = \frac{a}{\sqrt{(a + b)(a + c)}} \quad \text{covariance}$$

$$1 - d_{rs} = \frac{ad - bc}{\sqrt{(a + b)(a + c)(b + d)(c + d)}} \quad \text{coefficient de corrélation}$$

CALCUL DE DISSIMILARITÉS (2)

Données nominales

Remplacer les *e* états possibles d'une variable d'un type énumératif non ordonné par *e* variables booléennes et se ramener au cas des données binaires.

Données ordinales

Définir une matrice de *coefficients de discordance* (δ_{pq}) entre les états *p* et *q* avec :

$$\delta_{eq} > \delta_{e-1q} > \delta_{e-2q} > \dots > \delta_{qq} = 0 \quad \text{et} \quad \delta_{p1} > \delta_{p2} > \delta_{p3} > \dots > \delta_{pp} = 0$$

		État de l'objet <i>r</i> (<i>j</i> ^{ième} mesure)					
		1	2	...	<i>q</i>	...	<i>e</i>
État de l'objet <i>s</i> (<i>j</i> ^{ième} mesure)	1	0					
	2		0				
	...						
	<i>p</i>	δ_{p1}			δ_{pq}		
	...						
	<i>e</i>	$\delta_{e1}=1$			δ_{eq}		0

La dissimilarité d_{rs} entre les objets *r* et *s* sera la somme (pondérée) des coefficients de discordances pour chaque mesure.

CALCUL DE DISSIMILARITÉS (3)

Données numériques

$$\text{Distance de Manhattan : } d_{rs} = \sum_{j=1}^p |x_{rj} - x_{sj}| \quad \text{Distance euclidienne : } d_{rs} = \sqrt{\sum_{j=1}^p (x_{rj} - x_{sj})^2}$$

$$\text{Distance de Minkowski : } d_{rs} = \sqrt[\lambda]{\sum_{j=1}^p |x_{rj} - x_{sj}|^\lambda}$$

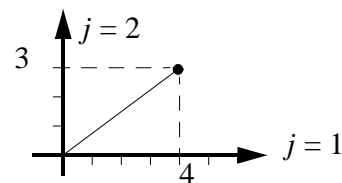
$\lambda = 1$: distance de Manhattan,

$\lambda = 2$: distance euclidienne, $\lambda = \infty$: maximum de la différence entre deux coordonnées.

$$\lambda = 1 : d_{rs} = 4 + 3 = 7$$

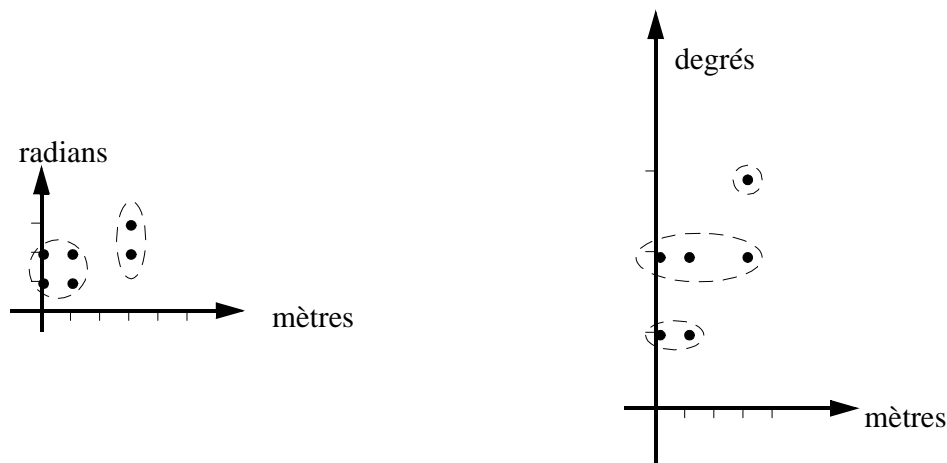
$$\lambda = 2 : d_{rs} = \sqrt{16 + 9} = 5$$

$$\lambda = \infty : d_{rs} = \max(4, 3) = 4$$



STANDARDISATION DES DONNÉES

Dans les calculs de dissimilarités présentés plus haut, on a fait l'hypothèse que toutes les mesures avaient la même importance et étaient basées sur la même unité, ce qui n'est pas forcément vrai. La structure des groupes peut être fortement influencée par les unités dans lesquelles les mesures ont été effectuées :



STANDARDISATION DES DONNÉES (2)

Pour éviter qu'une mesure prenne beaucoup plus d'importance que les autres, il convient de *normaliser* les unités.

Dans le cas d'une *distance de Manhattan*, il convient de diviser chaque composante par la différence *maximale* entre deux valeurs de cette composante. Dans le cas d'une *distance euclidienne*, il convient de diviser par l'*écart-type* des valeurs de chaque composante. De plus, il se peut que l'on désire donner explicitement un *poids* plus grand à une mesure qu'à une autre. On associera donc un poids w_j à chaque mesure, poids qui prendra en considération l'importance subjective que l'on veut donner à la $j^{\text{ème}}$ mesure divisée par un facteur de normalisation.

Exemple : on veut donner 4 fois plus d'importance à la 1^{ère} mesure qu'à la 2^e et utiliser des distances euclidiennes.

$$\text{On aura : } w_1 = \frac{4}{\frac{1}{n} \cdot \left(\sum (x_{i1})^2 - \frac{(\sum x_{i1})^2}{n} \right)}, \quad w_2 = \frac{1}{\frac{1}{n} \cdot \left(\sum (x_{i2})^2 - \frac{(\sum x_{i2})^2}{n} \right)}, \quad d_{rs} = \sqrt{\sum_{j=1}^p w_j (x_{rj} - x_{sj})^2}$$

MESURES DE DIFFÉRENTS TYPES

Très souvent, les données brutes sont de *plusieurs types*. Par exemple, un sol peut être caractérisé par son acidité (Ph, numérique), sa couleur (ordinal), le type de pierres présentes (nominal) et la présence de vers (binaires).

Si la *majorité des mesures sont d'un type*, on peut essayer de ramener toutes les données à ce type :

Numérique	→	Booléen	acide, basique
Booléen	→	Numérique	0, 1
Ordinal	→	Numérique	rang
Numérique	→	Nominal	classes non ordonnées
Numérique	→	Ordinal	classes ordonnées

S'il n'y a *pas de type prépondérant*, il faut passer par la définition d'un *poids à attribuer à chaque critère*.

VALEURS MANQUANTES

Si une mesure k est *manquante* pour l'élément r (et/ou s) on peut l'*ignorer* dans le calcul des dissimilarités :

$$d_{rs} = \sqrt{\sum_{j=1, j \neq k}^p w_j (x_{rj} - x_{sj})^2}$$

Une autre technique consiste à *remplacer* la valeur manquante par la *moyenne des autres valeurs* ou bien par la valeur de l'élément le plus similaire.

CRITÈRES DE SÉPARATION

On suppose dès maintenant que l'on dispose d'une matrice carrée dont les entrées mesurent la dissimilarité qu'il y a entre chaque paire d'éléments. La *séparation* d'un groupe G par rapport aux autres groupes peut être mesuré de diverses manières :

$$\text{séparation}(G) = \min_{r \in G, s \notin G} (d_{rs})$$

$$\text{coupe}(G) = \sum_{r \in G} \sum_{s \notin G} d_{rs} ; \quad \text{coupe normalisée}(G) = \frac{\sum_{r \in G} \sum_{s \notin G} d_{rs}}{|G|(n - |G|)}$$

Un *critère global* de mesure de *qualité* d'une classification, basé sur les séparations des groupes, peut être défini de diverses manières, selon que l'on se focalise sur le groupe le moins bien séparé ou sur la séparation moyenne de groupes :

$$\max_G(\text{séparation}(G)), \quad \max_G(\text{coupe}(G))$$

$$\max_G(\sum \text{séparation}(G)), \quad \max_G(\sum \text{coupe}(G))$$

CRITÈRES D'HOMOGENÉITÉ

L'homogénéité des éléments à l'intérieur d'un groupe G peut également être définie de diverses manières :

$$\text{diamètre}(G) = \max_{r, s \in G} (d_{rs})$$

$$\text{rayon}(G) = \min_{r \in G} (\max_{s \in G} (d_{rs}))$$

$$\text{étoile}(G) = \min_{r \in G} \left(\sum_{s \in G} d_{rs} \right) \quad \text{étoile normalisée}(G) = \frac{\min_{r \in G} \left(\sum_{s \in G} d_{rs} \right)}{|G| - 1}$$

$$\text{clique}(G) = \sum_{r \in G} \sum_{s \in G} d_{rs} \quad \text{clique normalisée}(G) = \frac{\sum_{r \in G} \sum_{s \in G} d_{rs}}{|G|(|G| - 1)}$$

Le critère *global* de qualité d'une classification, basé sur l'homogénéité consistera à *minimiser le plus grand diamètre* (rayon, étoile, ...) ou la *somme des diamètres* (rayon, ...) pour l'ensemble des groupes.

CRITÈRES BASÉS SUR DES CENTROÏDES

Dans le cas de mesures dans un *espace euclidien*, on peut mesurer l'homogénéité de G par rapport à un centre \bar{c} ne faisant pas partie des éléments. \bar{c} peut être défini comme le *centre de gravité* des éléments de G :

$$\bar{c} = \frac{1}{|G|} \sum_{i \in G} (x_{i1}, x_{i1}, \dots, x_{ip}) = \frac{1}{|G|} \sum_{i \in G} \vec{x}_i$$

Le critère d'homogénéité basé sur la *somme des carrés* des distances au centre \bar{c} est très répandu :

$$\text{SC}(G) = \sum_{i \in G} \left\| \vec{x}_i - \bar{c} \right\|^2$$

Pour ce critère, le centre de gravité \bar{c} est le point minimisant la valeur de $\text{SC}(G)$. Pour le critère basé sur la somme des distances :

$$\text{SD}(G) = \sum_{i \in G} \left\| \vec{x}_i - \tilde{c} \right\|$$

le centre \tilde{c} minimisant la valeur de la somme des distances ne peut plus s'exprimer analytiquement comme pour \bar{c} ; il est nécessaire d'*estimer numériquement* la position optimale de \tilde{c} pour un groupe G .

CLASSIFICATION HIÉRARCHIQUE

Méthode agglomérative

Idée générale : créer tout d'abord n groupes ne contenant qu'un élément. À chaque étape, fusionner les deux groupes les plus proches et recalculer une dissimilarité entre ce nouveau groupe et les autres.

Algorithme général

```

I = {1 .. n}          -- indices des groupes
Gi = {i}, ni = 1    -- groupes et nombre d'éléments par groupe
Calculer dij pour tout i, j ∈ I
Pour k = 1 .. n - 1 répéter :
    Trouver un dij minimum (i, j ∈ I)

    -- Fusionner les groupes i et j.
    Poser Gi ← Gi ∪ Gj; ni ← ni + nj

    -- Recalculer les distances
    dki ← αidik + αjdjk + βdij + δ|dik - djk| pour tout k ∈ I \ {i, j}
    dik ← dki pour tout k ∈ I \ {i, j}
    -- Mettre à jour l'ensemble des indices
    I ← I \ {j}
    
```



● k

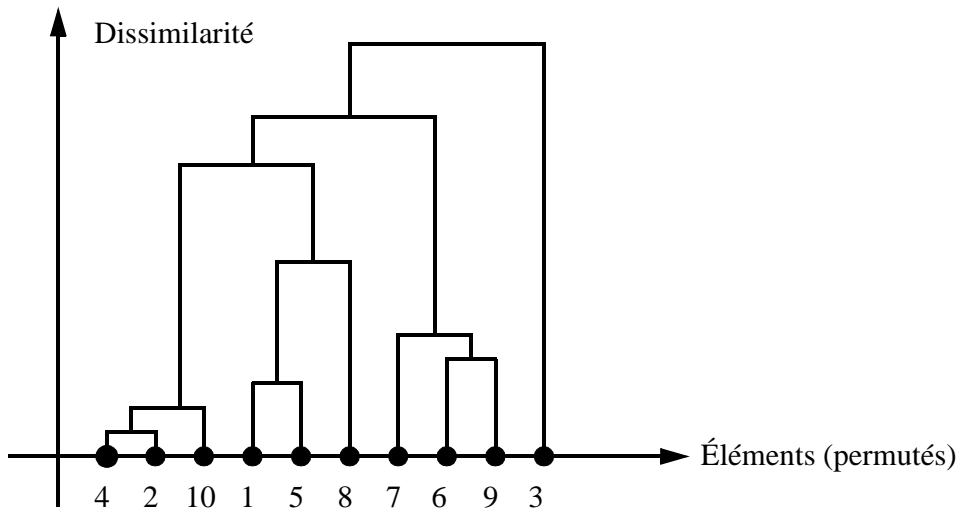
STRATÉGIES DE RECALCUL DES DISTANCES

Dans l'algorithme général précédent, plusieurs *stratégies* ont été proposées pour le calcul des distances (ligne : $d_{ki} \leftarrow \alpha_i d_{ik} + \alpha_j d_{jk} + \beta d_{ij} + \delta |d_{ik} - d_{jk}|$) une fois que deux groupes ont été fusionnés. Ces stratégies peuvent être caractérisées par diverses manières de calculer α_i , α_j , β , et δ :

	α_i	α_j	β	δ
<i>Lien simple</i> (plus proche voisin)	$\frac{1}{2}$	$\frac{1}{2}$	0	$-\frac{1}{2}$
<i>Lien complet</i> (voisin le plus éloigné)	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{1}{2}$
<i>Lien moyen</i>	$\frac{n_i}{n_i + n_j}$	$\frac{n_j}{n_i + n_j}$	0	0
<i>Centre</i>	$\frac{n_i}{n_i + n_j}$	$\frac{n_j}{n_i + n_j}$	$-\frac{n_i n_j}{n_i + n_j}$	0
<i>Somme des carrés incrémentale</i> (Ward)	$\frac{n_i + n_k}{n_i + n_j + n_k}$	$\frac{n_j + n_k}{n_i + n_j + n_k}$	$-\frac{n_k}{n_i + n_j + n_k}$	0

DENDROGRAMMES

Une fois qu'une classification hiérarchique a été trouvée, il est possible de la représenter *graphiquement* au moyen de *dendrogrammes* (diagramme en arbre) :



Durant la construction, il faut *réordonner les éléments* de sorte qu'il n'y ait pas de croisement entre ligne horizontale et verticale.

CLASSIFICATION GLOBALE

Problème de la p -médiane

Le problème de la p -médiane est très souvent utilisé pour procéder à une classification globale. Il peut être présenté ainsi : On cherche p centres parmi les n éléments, tels que si l'on rattache chaque élément à son centre le plus proche, la somme totale des dissimilarités est minimale. Le critère global d'optimisation est donc la minimisation de la somme des étoiles.

Exemple graphique :



Deux classifications possibles avec $p = 2$, selon les centres choisis.

FORMULATION MATHÉMATIQUE :

1) Étant donné une matrice D ($n \times n$) de dissimilarités, trouver p indices $J = \{j_1, j_2, \dots, j_p\}$

$$\text{minimisant : } \sum_{i=1}^n \min_{j \in J} (d_{ij})$$

2) Introduire des variables indicatrices $y_j = \begin{cases} 1 & j \text{ est sélectionné} \\ 0 & j \text{ non sélectionné} \end{cases}$ et $x_{ij} = \begin{cases} 1 & i \text{ est rattaché à } j \\ 0 & i \text{ non rattaché à } j \end{cases}$.

$$\text{minimiser } \sum_{i=1}^n \sum_{j=1}^p d_{ij} x_{ij} \quad \text{Somme des étoiles}$$

$$\text{sous } \sum_{j=1}^p x_{ij} = 1 \quad \forall i \quad \text{On doit rattacher chaque élément à un centre}$$

$$\text{contraintes } x_{ij} \leq y_j \quad \forall i, j \quad \text{On ne peut pas allouer un élément à un centre non sélectionné}$$

$$\sum_{j=1}^p y_j = p \quad \text{On doit ouvrir } p \text{ centres}$$

$$x_{ij} \in \{0, 1\} \quad \forall i, j$$

$$y_j \in \{0, 1\} \quad \forall j$$

COMPLEXITÉ

Le problème de la p -médiane est *difficile* : on ne connaît pas d'algorithme polynomial pour le résoudre (c'est-à-dire en $O(n^c)$ avec une constante c ne dépendant ni de n , ni de p). Par contre, si p est fixé, il existe un algorithme polynomial évident, mais qui ne peut être mis en œuvre que pour des valeurs de p très petites :

Pour $j_1 = 1 \dots n - p + 1$

Pour $j_2 = j_1 \dots n - p + 2$

 ...

Pour $j_p = j_{p-1} \dots n$

 Ouvrir les centres j_1, j_2, \dots, j_p

 Trouver le rattachement optimum des éléments

 Calculer la valeur de l'objectif global

Si la solution la meilleure trouvée jusqu'ici, la mémoriser

ALGORITHMES HEURISTIQUES DE RÉOLUTION

1) **K-Means** (dans ce nom, le K fait référence au nombre de centre i.e. $\equiv p$)

Demander à un expert de placer les p centres. En l'absence d'expert, utiliser une procédure constructive simple (par exemple, placer les centres aléatoirement)

Répéter

Allouer les éléments au centre le plus proche

Évaluer la qualité de la solution

Pour chaque groupe G d'éléments alloués au même centre, trouver le meilleur centre possible parmi les éléments de G

Tant que l'on réussi à améliorer la qualité de la solution.

2) **K-Means itéré**

Répéter I fois

Placer les centres aléatoirement

Appliquer K -Means

Mémoriser la solution si elle est meilleure

Tant que l'on réussi à améliorer la qualité de la solution.

Retourner la meilleure solution trouvée.

3) **K-Means avec repositionnement**

Placer les p centres aléatoirement (= *solution_courante*);

Poser *meilleure_solution* := *solution_courante*;

Répéter

Pour $i = 1..p$, **répéter**

Pour $j = 1..n$, **s'il n'y a pas de centre en j , répéter**
solution_courante := *meilleure_solution*

Déplacer le centre i en j dans *solution_courante*

Répéter -- appliquer K -Means à *solution_courante*

Allouer les éléments au centre le plus proche

Évaluer la qualité de *solution_courante*

Pour chaque groupe G d'éléments alloués au même centre :

Placer le centre optimalement dans *solution_courante*

Tant que l'on a réussi à améliorer la qualité *solution_courante*

Si *solution_courante* est meilleure que *meilleure_solution*

meilleure_solution := *solution_courante*

Fin_si

Fin_pour j

Fin_pour i

Tant que l'on a réussi à améliorer la qualité de *meilleure_solution*

Retourner *meilleure_solution*.