

Russel Rao $s(i_1, i_2) = \frac{2}{2+1+0+1} = \frac{2}{6}$

II / Partition et hierarchie.

II.1 / Partition

Soit E une population à n individus

On note par $\mathcal{P}(E)$, l'ensemble de toutes les parties de E .

et soit $P(E) = \{E_1, \dots, E_k\} \subset \mathcal{P}(E)$ tp :

- $E_i \neq \emptyset \quad \forall i=1, \dots, k$
- $\bigcup_{i=1}^k E_i = E$
- $E_i \cap E_j = \emptyset$ pour $i \neq j$

$P(E)$ est une partition de E .

Exp: $E = \{1, 2, 3, 4\}$. $P(E) = \{\{1, 3\}, \{2, 4\}\}$ partition de E

II.2 / Hierarchisation de parties d'un ensemble

Def 1: Soit E un ensemble finie. et soit H un ensemble de parties de E ($H \subset \mathcal{P}(E)$)

H est une hierarchie sur E si:

- $\forall x_i \in E \quad \{x_i\} \in H$

- $E \in H$

- $\forall h_1$ et $h_2 \in H$ on a soit $h_1 \cap h_2 = \emptyset$ ou $h_1 \subset h_2$ ou $h_2 \subset h_1$

(si l'intersection est vide l'une est incluse dans de l'autre.

• toute classe $h \in H$ est la reunion de classes qui sont contenues de

Exp: $E = \{1, 2, 3, 4, 5, 6\}$.

$H = \{\emptyset, \{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\}, \{1, 2\}, \{1, 2, 3\}, \{4, 5, 6\}, \{1, 2, 3, 4, 5, 6\}\}$

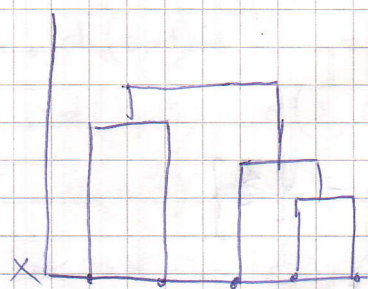
Def 2: Une hierarchie indexee est un couple (H, f) ou H est une hierarchie et $f: H \rightarrow \mathbb{R}_+$ tp pour $h \in H$

- $f(h) = 0 \Leftrightarrow h$ ne contient qu'un seul elt

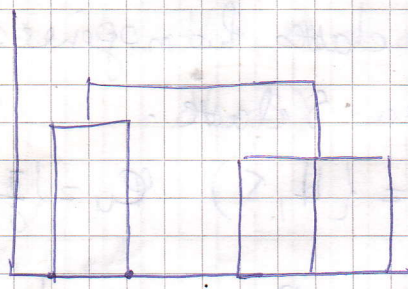
- $\forall h, h' \in H \quad h \subset h' \Rightarrow f(h) \leq f(h')$ " $f \uparrow$ "

Rq: si l'inegalite est large, on parle de l'hierarchie indexee au sens large.

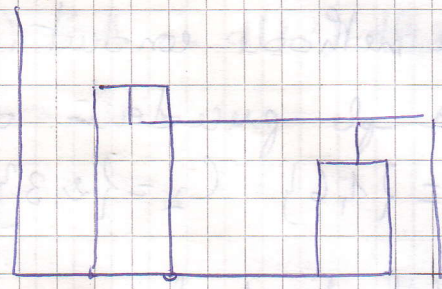
• S'il existe h et h' tq $h \subset h'$ et $f(h) > f(h')$ on dit qu'il y a inversion.



hierarchy indexée



hierarchy indexée
au sens large



hierarchy avec inversion

inverses sont appelés niveau d'agregation. ** voir page sup.

IV / Méthode de classification basée sur la ressemblance.

Soit $E = \{x_1, \dots, x_n\}$ à n inds à classifier. Partir.

Le but : chercher à représenter ses inds (p^b) par des ensembles de pa hiérarchiquement emboîtés

Afin de procéder au regroupement de inds de E , on calcule les distances (ou dissimilarités) entre les inds de E puis 2 à 2.

Plusieurs techniques ont été utilisées, citons la plus importante.

a) Méthode de liaison simple : consiste à regrouper les inds si les inds d'une m classe ont une distance (ou dissimilarité) $d(i, i') \leq \delta$ à un seuil fixé δ .

si par exemple C est une classe $\forall i, i' \in C \Rightarrow d(i, i') \leq \delta$

Rq : La procédure est simple, mais on aboutit généralement à des classes hétérogènes ds le sens où chaque ind n'est pas forcément relié à tous les inds de la classe.

Exp : $E = \{1, 2, 3, 4, 5, 6, 7, 8\}$ et la matrice de distances est

$E \setminus E$	1	2	3	4	5	6	7	8
1	0							
2	5	0						
3	3	1	0					
4	2	6	4	0				
5	5	2	2	1	0			
6	1	3	4	3	3	0		
7	6	5	2	2	2	3	0	
8	2	1	5	4	4	3	2	0

si $\delta = 1$ ~~est~~ \Rightarrow 4 classes

$C_1 = \{1, 6\}$ $C_2 = \{2, 3, 8\}$

$C_3 = \{4, 5\}$ $C_4 = \{7\}$

$d(2, 3) = d(8, 2) = 1 \leq \delta$

mais $d(8, 3) = 5 > \delta$

C_2 est donc hétérogène

b) Méthode de liaison complète: Chaque ind. est relié à tous les ind. de sa classe à un niveau $\leq \delta$.

Cette méthode conduit à des classes homogènes.

Exemple précédent on aura 5 classe.

$$C_1 = \{1, 6\} \quad C_2 = \{2, 3\} \quad C_3 = \{4, 5\} \quad C_4 = \{7\} \quad C_5 = \{8\}$$

c) Méthode de liaison moyenne = Regroupement par étapes successives.

Chercher la paire de base constituée d'ind. le plus proches, calculer les distances moyennes (ou dissimilarité moyennes) entre cette paire et les autres ind. On obtient ainsi une matrice de dimension $(n-1, n-1)$. repérer la distance (ou dissimilarité) la plus faible $\leq \delta$ à un seuil donné.

- si cette distance (ou dissimilarité) correspond à une liaison entre deux ind. isolés, ces derniers sont regroupés ds une n^o classe nouvelle matrice.

- si elle correspond à la liaison entre la paire ^{de base} et 1 ind. isolé.

On construit une classe à 3 ind. et ainsi de suite

L'opération s'arrête lorsque les termes de la dernière matrice obtenue soit strictement supérieure au seuil fixe.

Exemple. Soit le tableau de distance précédent. On fixe $\delta = 1$

Soit $\{2, 3\}$ la paire de base - de départ

	1	{2,3}	4	5	6	7	8
1	0						
{2,3}	$\frac{5+3}{2} = 4$	0					
4	2	$\frac{6+4}{2} = 5$	0				
5	5	$\frac{8+2}{2} = 5$	1	0			
6	1	$\frac{3+4}{2} = 3,5$	3	3	0		
7	6	$\frac{5+2}{2} = 3,5$	2	2	3	0	
8	2	$\frac{1+5}{2} = 3$	4	4	3	2	0

La distance la plus faible $\leq \delta = 1$

correspond à la liaison entre 5 et 4

$$d(5,4) = 1 \leq 1$$

On reuni ds une nouvelle matrice (6,6)

$E \setminus E$	1	{2,3}	{4,5}	6	7	8
1	0					
{2,3}	4	0				
{4,5}	$\frac{2+5-3}{2}$	$\frac{5+1}{2}$	0			
6	1	3,5	$\frac{3+3-3}{2}$	0		
7	6	3,5	$\frac{2+6-2}{2}$	3	0	
8	2	3	$\frac{4+4-4}{2}$	3	2	0

$$d(1,6) = 1 \leq 1$$

on réunit 1 et 6 ds 1 m^e classe -
on obtient une matrice (5x5)

	{1,6}	{2,3}	{4,5}	7	8
{1,6}	0				
{2,3}	$\frac{4+3-3}{2}$	0			
{4,5}	$\frac{3+3-3}{2}$	5	0		
7	$\frac{6+3-4}{2}$	3,5	2	0	
8	$\frac{2+3-2}{2}$	3	4	2	0

On arrête puisque toute les
distances sont > au seuil $\delta = 1$
on obtient donc les classes suivantes
{1,6}, {2,3}, {4,5}, {7}, {8}.

V Classification hiérarchique ascendante.

V.1/ Hiérarchie de partie d'un ensemble.

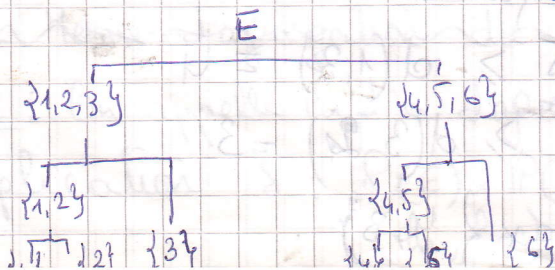
Le but est de regrouper progressivement les inds d'un ensemble E sur la base de leurs ressemblance. Une hiérarchie peut être représentée par un arbre de classification.

Exemple.

Soit $E = \{1, 2, 3, 4, 5, 6\}$.

et $H = \{\emptyset, \{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\}, \{1,2\}, \{1,2,3\}, \{4,5\}, \{4,5,6\}, E\}$
exprèsèdent hiérarchie de parties de E.

On construit un arbre correspondant à H de la manière suivante.



C'est un arbre ascendant.

On regroupe d'abord deux objets inds le plus proches qui forme un noeud ou sommet, il ne reste donc qu'un seul objet ind.

que (n-1) inds et on répète l'opération jusqu'à regroupement complet.

V.2/ Compatibilité d'une hiérarchie avec une partition.

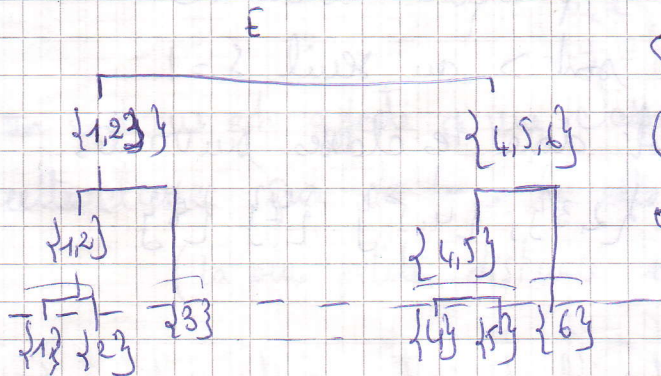
Def: une hiérarchie H est compatible à une partition d'un ensemble E si toute classe de la partition est ds H.

Exemple:

Soit la partition de $E = \{1, 2, 3, 4, 5, 6\}$.

$C_1 = \{1, 2\}$ $C_2 = \{3\}$ $C_3 = \{4, 5\}$ $C_4 = \{6\}$.

ds l'ex^o précédent C_1, C_2, C_3 et C_4 appartiennent à H.



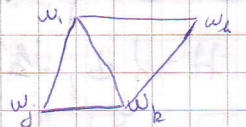
Si on coupe l'arbre par une horizontale (regardez les pointillés) on obtient la partition juste sur l'horizontal.

V.3/ Hiérarchie induite et ultramétrique.

definition: Une ultra métrique est une application $\delta: E \times E \rightarrow \mathbb{R}_+$ tp

- 1 - $\forall w_i, w_j. \delta(w_i, w_j) \geq 0$
- 2 - $\delta(w_i, w_j) = 0 \Rightarrow w_i = w_j$
- 3 - $\delta(w_i, w_j) = \delta(w_j, w_i)$
- 4 - $\forall i, j, k. \delta(w_i, w_j) \leq \max \{ \delta(w_i, w_k), \delta(w_k, w_j) \}$

Proposition: Une condition suffisante et nécessaire pour que une distance soit une ultramétrique est que tous les triangles de E soient isocèles avec la base inférieure aux côtés



Exemple.

$\frac{E}{E}$	1	2	3	4
1	0	5	4	5
2	5	0	3	1
3	4	3	0	2
4	5	1	2	0

c'est un tab de distance mais pas ultramétrique.

$$d(1, 2) = 5 \geq d(1, 3) = 4$$

$$\geq d(3, 4) = 3$$

$$\Rightarrow d(1, 2) \geq \max \{ 4, 3 \}$$

Le tab suivant est un tab d'ultramétrique

$\begin{matrix} E \\ E \end{matrix}$	1	2	3	4
1	0	6	0	6
2	6	0	5	2
3	0	5	0	5
4	6	2	5	0

Théorème : Il existe une bijection entre une hiérarchie induite et ultramétrique.

En d'autres termes, toute hiérarchie induite sur E correspond une ultramétrique et inversement.

On propose deux familles de méthodes pour obtenir une hiérarchie induite :

- Les regroupements progressifs
- Les passages directs à l'ultramétrique

V.4 / Les regroupements progressifs.

On a plusieurs étapes successives.

Sur le tableau des distances ou dissimilarité, on repère la plus petite valeur.

Soit A la classe $\{(w_i, w_{i'}) \mid d(w_i, w_{i'}) = \text{plus petite valeur}\}$.

On calcule la distance ou dissimilarité de chaque ind à A , on obtient un tableau de dim $(n-1, n-1)$.

On repère du nouveau tableau la dist ou diss la plus petite.

On répète la procédure jusqu'à obtenir un regroupement complet.

R_p si la + petite valeur correspond à une dist ou diss entre deux autres ind il seront regroupés ds une m classe, si cela correspond à dist ou diss entre un ind et A on forme une classe B de 3 ind's.

On a 3 manières de calculer les distances par utilisation d'indice au ~~titre~~ d'agrégation du ~~lien~~ minimum. d'agrégation

Def : On appelle indice d'agrégation entre deux groupes d'ind i et i' une

$$\text{application } \mathcal{F} : \mathcal{P}(E) \times \mathcal{P}(E) \longrightarrow \mathbb{R}^+ \\ (h_1, h_2) \longmapsto \mathcal{F}(h_1, h_2) \text{ tq } \mathcal{F}(h_1, h_2) = \mathcal{F}(h_2, h_1)$$