

M. BESSON

**À propos des distances entre ensembles de parties**

*Mathématiques et sciences humaines*, tome 42 (1973), p. 17-35

[http://www.numdam.org/item?id=MSH\\_1973\\_\\_42\\_\\_17\\_0](http://www.numdam.org/item?id=MSH_1973__42__17_0)

© Centre d'analyse et de mathématiques sociales de l'EHESS, 1973, tous droits réservés.

L'accès aux archives de la revue « Mathématiques et sciences humaines » (<http://msh.revues.org/>) implique l'accord avec les conditions générales d'utilisation (<http://www.numdam.org/conditions>). Toute utilisation commerciale ou impression systématique est constitutive d'une infraction pénale. Toute copie ou impression de ce fichier doit contenir la présente mention de copyright.

NUMDAM

Article numérisé dans le cadre du programme  
Numérisation de documents anciens mathématiques  
<http://www.numdam.org/>

## A PROPOS DES DISTANCES ENTRE ENSEMBLES DE PARTIES

par  
M. BESSON <sup>1</sup>

### RÉSUMÉ

*Soit E un ensemble fini. Nous nous proposons ici d'établir des indices de distance entre les divers ensembles de parties de E. Nous généraliserons pour cela, à  $\mathcal{P}\mathcal{P}(E)$ , la notion de fonction caractéristique bien connue dans  $\mathcal{P}(E)$ .*

### SUMMARY

*Let E be a finite set. We demonstrate a method to establish distances between different sets of sub-sets of E. We will apply in a generalised form the well-known characteristic function of  $\mathcal{P}(E)$  to  $\mathcal{P}\mathcal{P}(E)$ .*

1. Nous nous proposons d'étudier ici les distances entre ensembles de parties d'un ensemble E fini, que nous convenons d'appeler, pour simplifier l'écriture, les familles de E.

#### 1.1. SOIT E UN ENSEMBLE FINI SOIT R UNE FAMILLE DE E

Un certain nombre de conditions sur R ordonnées par l'implication forment un sup demi-treillis explicité par le schéma 1.

#### 1.2. FONCTIONS CARACTÉRISTIQUES

Nous convenons d'adopter la définition suivante :

Une fonction caractéristique pour un certain type de famille de E est une application de E dans  $N^+^2$  qui caractérise chaque famille de ce type.

Établir une fonction caractéristique pour un certain type de famille de E, c'est donc trouver une injection de l'ensemble des familles de ce type dans l'ensemble des applications de E dans  $N^+$ .

#### 1.3. VECTEURS CARACTÉRISTIQUES

Soient  $x_1, x_2, \dots, x_n$  les éléments de E. Soit f une fonction caractéristique pour une certaine famille R sur E. Nous appellerons « Vecteur caractéristique » associé à (f, R) le vecteur :

$$\vec{V} = \{f(x_1), f(x_2) \dots f(x_n)\}$$

---

1. Ecocentre de MacGregor Comarain. Directeur Nguyen Tien Phuc, 219, rue de Versailles, 92-Ville d'Avray.  
2. Ensemble des entiers positifs ou nuls.

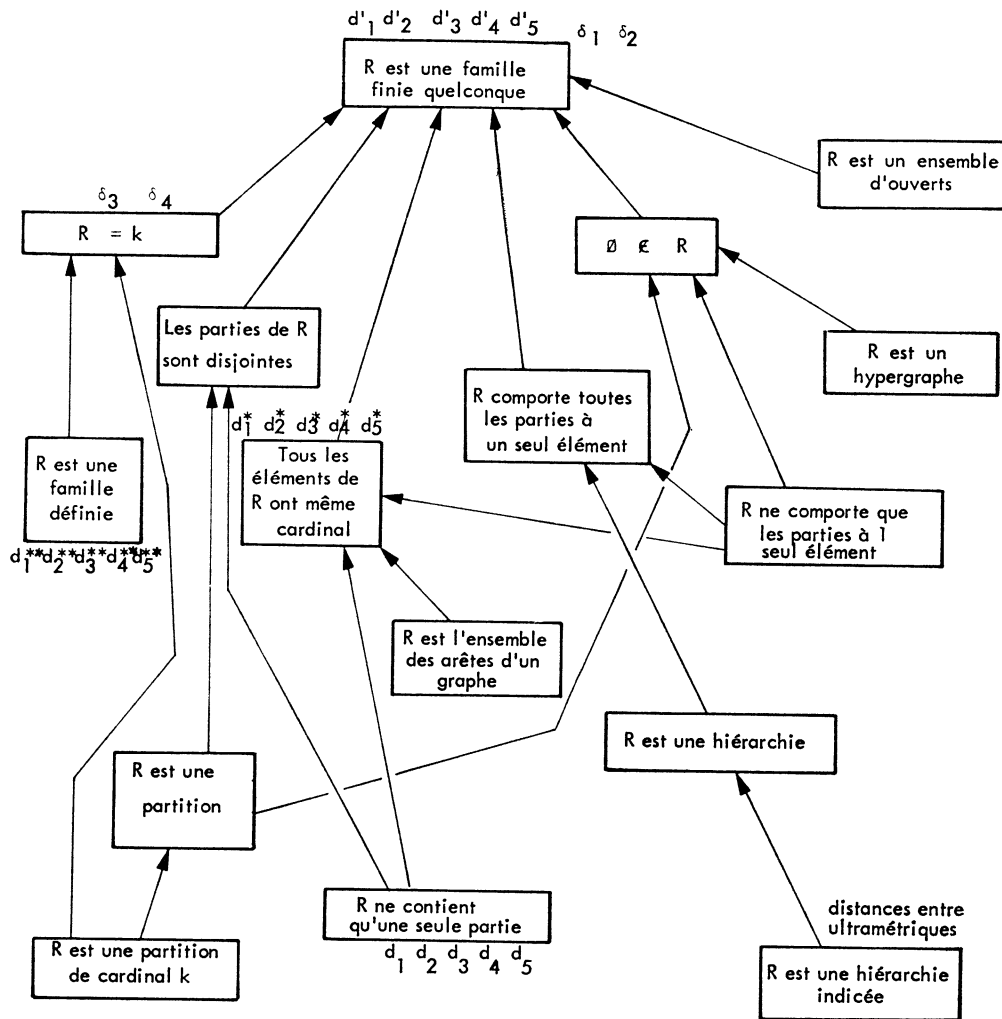


Schéma 1

**Théorème**

Soit T un certain type de familles sur E, et  $f_T$  une fonction caractéristique adaptée à ce type. Toute distance entre les vecteurs caractéristiques construits avec  $f_T$ , induit une distance entre les familles de type T.

**Remarque :** Ce théorème n'est qu'un cas particulier du théorème suivant :

« Soient E et E' deux ensembles quelconques. Soit  $\Psi$  une injection de E dans E'. Soit d une distance sur E' »

Définissons  $\delta$ , application de  $E \times E$  dans  $R^+$ , par l'égalité suivante :

$$(\forall x \in E) (\forall y \in E) (\delta(x, y) = d(\Psi(x), \Psi(y)))$$

$\delta$  est une distance sur E.»

**II. LES PARTIES**

Nous nous proposons d'étudier maintenant les familles les plus simples : les parties (familles de cardinal 1).

### 2.1. Fonctions caractéristiques

La fonction caractéristique d'une partie X de E est l'application  $f_X$  de E dans [0, 1] telle que :

$$(\forall x \in X) (f_X(x) = 1)$$

$$(\forall x \notin X) (f_X(x) = 0)$$

Le vecteur caractéristique correspondant à une telle partie X est le vecteur :

$$\vec{V}_f(X) = \{f_X(x_1), f_X(x_2) \dots f_X(x_n)\}$$

(vecteur qui ne comporte évidemment que des 0 ou des 1.)

### 2.2. Distances induites par cette fonction caractéristique

Soient X et Y deux parties quelconques de E.

*Distance euclidienne :*

Écrivons la distance euclidienne entre les vecteurs  $V_f(X)$  et  $V_f(Y)$  :

$$\begin{aligned} d(V_f(X), V_f(Y)) &= \sqrt{\sum_{i=1}^n (f_X(X_i) - f_Y(Y_i))^2} \\ &= ||\vec{V}_f(X) - \vec{V}_f(Y)|| \\ &= \sqrt{|X \Delta Y|}^1 \end{aligned}$$

Distance induite correspondante sur l'ensemble des parties de E

$$d_1(X, Y) = \sqrt{|X \Delta Y|}$$

*Distance : somme des différences des coordonnées en valeur absolue*

$$\begin{aligned} d_2(\vec{V}_f(X), \vec{V}_f(Y)) &= \sum_{i=1}^n |f_X(x_i) - f_Y(x_i)| \\ &= |X \Delta Y| \end{aligned}$$

*Distance : sup. des différences des coordonnées en valeur absolue*

$$d(\vec{V}_f(X), \vec{V}_f(Y)) = \text{Sup}_{i=1, n} (|f_X(x_i) - f_Y(x_i)|)$$

La distance induite sur l'ensemble des parties de E sera donc :

$$d_s(X, Y) = \begin{cases} 1 & \text{si } X \neq Y \\ 0 & \text{si } X = Y \end{cases}$$

### 2.3. Autres distances

Soient  $x_1, x_2 \dots x_n$  les éléments de E. Considérons trois parties A, B et C de E :

$$A = \{x_1, x_2, x_3\}$$

---

1.  $\Delta$  représente la différence symétrique.

$$B = \{x_4\}$$

$$C = \{x_1, x_2, x_3, x_4, x_5, x_6, x_7\}$$

Appliquons-leur les distances  $d_1$ ,  $d_2$  et  $d_3$  :

$$|A \Delta B| = 4$$

$$|A \Delta C| = 4$$

$$\text{Donc : } d_1(A, B) = d_1(A, C)$$

$$d_2(A, B) = d_2(A, C)$$

$$d_3(A, B) = d_3(A, C) = 1$$

Et cependant, n'avons-nous pas instinctivement l'impression que la ressemblance entre A et B est bien plus faible que celle entre A et C, qui ont trois éléments en commun, alors que A et B sont disjoints.

Pour pallier à cet inconvénient, et que des parties disjointes correspondent à des distances particulièrement fortes, nous avons cherché à établir d'autres distances (ou indices de distance).

*Distance variant entre 0 et 1 :*

Pour tout couple X et Y de parties de E, non-vides à la fois, posons :

$$d_4(X, Y) = \frac{|X \Delta Y|^1}{|X \cup Y|} \quad \text{et} \quad d_4(\emptyset, \emptyset) = 0$$

$d_4$  vaut 0 si et seulement si  $X = Y$

$d_4$  vaut 1 si et seulement si  $X \cap Y = \emptyset$

sans que X et Y soient simultanément vide.

Ainsi nous aurions dans l'exemple précédent

$$d_4(A, B) = \frac{4}{4} = 1$$

$$d_4(A, C) = \frac{4}{7}$$

Nous pouvons exprimer cette distance d avec les fonctions caractéristiques :

$$d_4(X, Y) = \frac{\sum_{i=1}^n (f_X(x_i) + f_Y(x_i) - 2f_X(x_i) f_Y(x_i))}{\sum_{i=1}^n (f_X(x_i) + f_Y(x_i) - f_X(x_i) f_Y(x_i))}$$

ou avec les vecteurs caractéristiques :

$$d_4(X, Y) = \frac{(\vec{V}_f(X) + \vec{V}_f(Y)) \cdot \vec{1} \cdot \vec{1} - 2 \vec{V}_f(X) \cdot \vec{V}_f(Y)}{(\vec{V}_f(X) + \vec{V}_f(Y)) \cdot \vec{1} - \vec{V}_f(X) \cdot \vec{V}_f(Y)}$$

1. Cette distance a été utilisée également par M. Brissaud voir [6].

La mesure de similarité proposée par Jaccard en 1908 peut ainsi s'écrire :  $f(X, Y) = 1 - d_4(X, Y)$ .

2. Vecteur unité.

$d_4$  vérifie bien en effet les axiomes d'une distance.

i)  $d_4(X, Y) = 0 \Rightarrow X = Y$  (résultat trivial)

ii)  $d_4(X, Y) = d_4(Y, X)$  (résultat trivial)

iii)  $d_4$  vérifie l'inégalité triangulaire <sup>1</sup>

*Indice de distance variant entre 0 et l'infini :*

Pour tout couple X et Y de parties de E non simultanément vides, posons :

$$d_5(X, Y) = \frac{|X \Delta Y|^2}{|X \cap Y|} = \frac{\sum_{i=1}^n (f_X(x_i) + f_Y(x_i) - 2f_X(x_i) \cdot f_Y(x_i))}{\sum_{i=1}^n f_X(x_i) \cdot f_Y(x_i)}$$

$$= \frac{(\vec{V}_f(X) + \vec{V}_f(Y)) \cdot \vec{1} - 2\vec{V}_f(X) \cdot \vec{V}_f(Y)}{\vec{V}_f(X) \cdot \vec{V}_f(Y)}$$

et  $d_5(X, X) = 0$ .

Ainsi, nous aurions dans l'exemple précédent :

$d_5(A, B) \rightarrow \infty$

$d_5(A, C) = 4/3$

$d_5$  est bien un indice de distance sur l'ensemble des parties de E.

Mais  $d_5$  n'est pas une distance : l'inégalité triangulaire n'est pas vérifiée.

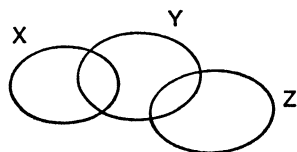
En effet, soient X, Y et Z trois parties d'un ensemble E fini quelconque telles que :

$X \cap Y \neq \emptyset$

$X \cap Z = \emptyset$

$Y \cap Z \neq \emptyset$

Nous avons bien sûr :  $d_5(X, Z) > d_5(X, Y) + d_5(Y, Z)$



#### 2.4. CORRÉLATION ENTRE DEUX VECTEURS CARACTÉRISTIQUES ET INDICE DE DISTANCE CORRESPONDANT

Soit E un ensemble quelconque, à  $n$  éléments. Soient X et Y deux parties de E. Les séries statistiques formées par les diverses coordonnées de leurs vecteurs caractéristiques  $V_f(X)$  et  $V_f(Y)$  correspondent au coefficient de corrélation suivant :

1. Cette démonstration a été faite par G. Morlat, voir [7].

2.  $d_5$  se trouve être l'inverse de la mesure de similarité proposée par Kuleyzinski en 1927, voir [4].

$$r(\vec{V}_f(X), \vec{V}_f(Y)) = \frac{\sum_{i=1}^n (f_X(x_i) - m_X) (f_Y(x_i) - m_Y)}{\sqrt{\sum_{i=1}^n (f_X(x_i) - m_X)^2} \sqrt{\sum_{i=1}^n (f_Y(x_i) - m_Y)^2}}$$

$$\text{avec } m_X = \frac{\sum f_X(x_i)}{n} = \frac{|X|}{n}$$

$$\text{et } m_Y = \frac{\sum f_Y(x_i)}{n} = \frac{|Y|}{n}$$

$$\begin{aligned} r(\vec{V}_f(X), \vec{V}_f(Y)) &= \frac{\sum_{i=1}^n f_X(x_i) (f_Y(x_i) - \frac{|Y|}{n}) - \frac{|X| \cdot |Y|}{n}}{\sqrt{\sum x_i^2 - \frac{|X|^2}{n}} \cdot \sqrt{\sum y_i^2 - \frac{|Y|^2}{n}}} \\ &= \frac{n |X \cap Y| - |X| \cdot |Y|}{n \cdot \sqrt{|Y| - \frac{|Y|^2}{n}} \cdot \sqrt{|X| - \frac{|X|^2}{n}}} \end{aligned}$$

$$r(\vec{V}_f(X), \vec{V}_f(Y)) = \frac{n|X \cup Y| - |X| \cdot |Y|}{\sqrt{|X| \cdot |Y|} \cdot \sqrt{|X| \cdot |Y|}}$$

**Remarque :** Si l'écart-type d'une série statistique est nul, le coefficient de corrélation est considéré comme nul également. De même, ici, par une convention analogue, si une des deux parties X et Y est  $\emptyset$  ou E, r est nul aussi.

**Indice de distance correspondant :**

$$\text{Posons : } d_0(X, Y) = 1 - r(\vec{V}_f(X), \vec{V}_f(Y))$$

$$d_0(X, Y) = \begin{cases} 1 + \frac{|X| \cdot |Y| - n |X \cup Y|}{\sqrt{|X| \cdot |Y|} \cdot \sqrt{|X| \cdot |Y|}} \\ 1 \text{ si } X = \emptyset \text{ ou } E \\ \text{ou} \\ \text{si } Y = \emptyset \text{ ou } E \end{cases}$$

Il s'agit bien d'un indice de distance. En effet :

$$d_0(X, Y) = d_0(Y, X)$$

$$d_0(X, X) = 0$$

$$d_0(X, Y) = 0 \Rightarrow \vec{V}_f(X) \text{ et } \vec{V}_f(Y) \text{ homothétiques sont donc égaux.}$$

Cet indice de distance  $d_0$  varie entre 0 et 2.

— Il vaut 2 si et seulement si :

$$r(\vec{V}_f(X), \vec{V}_f(Y)) = -1$$

$$\text{Donc : } Y = \bar{X}$$

— Il vaut 1 si et seulement si :

$$r(\vec{V}_f(X), \vec{V}_f(Y)) = \frac{|X| \cdot |Y|}{n}$$

(On pourrait dire alors que les parties X et Y sont «indépendantes», par analogie avec des séries statistiques quelconques.)

Mais est-ce une distance ?

C'est-à-dire : l'inégalité triangulaire est-elle vérifiée ?

La question reste ouverte.

## 2.5. STABILITÉ PAR PASSAGE AU COMPLÉMENTAIRE

Une propriété assez intéressante des indices de distance entre parties d'un ensemble E, est leur stabilité par passage au complémentaire.

C'est-à-dire :

$$(\forall X \subset E) (\forall Y \subset E) (d(\bar{X}, \bar{Y})) = d(X, Y)$$

Puisque « $X \Delta Y = \bar{X} \Delta \bar{Y}$ », les distances  $d_1$ ,  $d_2$  et  $d_3$  seront stables. Par contre, les distances  $d_4$ ,  $d_5$  et  $d_0$  ne le seront pas, en général.

—  $d_4$  et  $d_5$  ne vérifieront la propriété de stabilité que si :

$$|X \cap Y| = |\bar{X} \cap \bar{Y}|$$

— Tandis que  $d_0$  ne la vérifiera que si :

$$n |X \cap Y| - |X| \cdot |Y| = n |\bar{X} \cap \bar{Y}| - |\bar{X}| \cdot |\bar{Y}|$$

$$\text{donc : } n |X \cap Y| = |X| \cdot |Y|$$

(parties indépendantes).

## III. LES PARTITIONS

Pour définir des distances dans des ensembles de partitions, nous essaierons d'abord de trouver des fonctions caractéristiques, correspondant à des vecteurs caractéristiques. Toute distance entre les vecteurs induira une distance entre les partitions.



Puis, nous essaierons par une démarche plus directe, de déterminer des indices de distance plus intéressants de divers points de vue.

### 3.1. FONCTIONS CARACTÉRISTIQUES

Nous cherchons ici des fonctions de E dans  $N^+$  qui caractérisent chaque partition.

Soit  $E = \{x_1, x_2, \dots, x_n\}$

Soit T une partition de E.

Pour tout élément y de E, posons :

$$f_R(y) = \text{Inf } \{i / (\exists X \in R) (y, x_i) \in X\}$$

C'est-à-dire :  $f_R(y)$  est le plus petit des indices  $i$  tels que  $x_i$  soit réuni avec y dans la partition R.

Soit R' une partition de E, différente de R. Donc :

$$(\exists y \in E) (\exists z \in E) ((y, z) \text{ réunis par } v \text{ et non par } R' \text{ ou } (y, z) \text{ réunis par } R' \text{ et non par } R).$$

Or, si (y, z) est un couple d'éléments de E réunis par R et non pas par R' on a :

$$f_R(y) = f_R(z)$$

$$\text{et } f_{R'}(y) \neq f_{R'}(z) ..$$

Donc les fonctions  $f_R$  et  $f_{R'}$ , sont différentes.

Nous venons donc de définir une injection de l'ensemble des partitions dans l'ensemble des applications de E dans  $N^+$ , c'est-à-dire des fonctions caractéristiques pour les partitions de E.

*Vecteur caractéristique correspondant :*

$$V_f(R) = (f_R(x_1), f_R(x_2) \dots f_R(x_n)) .$$

Les distances euclidiennes, *somme*, et *sup* des différences des coordonnées en valeur absolue, impliquent évidemment des distances dans l'ensemble des partitions. Mais que représentent-elles dans la réalité ?

Il semble que les cas où elles sont bien adaptées sont assez rares.

### 3.2. INDICES DE DISTANCE

Acceptant de ne plus vérifier l'inégalité triangulaire, nous cherchons à établir maintenant de simples indices de distances, (plus significatifs) de la différence réelle entre deux partitions.

Soient Q et R deux partitions quelconques.

Posons :

$$\Psi(Q, R) = \frac{1}{|R|} \left[ \sum_{Z \in R} \text{Inf}_{X \in Q} \frac{|X \Delta Z|}{|X \cup Z|} \right] \text{ Si Q et R non vides à la fois.}$$

C'est-à-dire : pour chaque élément Z de R, nous cherchons l'élément X de Q tel que le rapport des cardinaux de  $X \cap Z$  et de  $X \cup Z$  soit le plus grand possible,  $\Psi$  est la moyenne des quotients ainsi obtenus.

Puis, nous définissons des applications  $\delta_1$  et  $\delta_2$  par les relations :

$$\begin{aligned}\delta_1(\mathbf{R}, \mathbf{Q}) &= \frac{1}{2} [\Psi(\mathbf{R}, \mathbf{Q}) + \Psi(\mathbf{Q}, \mathbf{R})] \\ &= \frac{1}{2|\mathbf{R}|} \left[ \sum_{\mathbf{x} \in \mathbf{R}} \inf_{\mathbf{z} \in \mathbf{Q}} d_4(\mathbf{X}, \mathbf{Z}) \right] + \frac{1}{2|\mathbf{Q}|} \left[ \sum_{\mathbf{z} \in \mathbf{Q}} \inf_{\mathbf{x} \in \mathbf{R}} d_4(\mathbf{X}, \mathbf{Z}) \right] \\ \delta_2(\mathbf{R}, \mathbf{Q}) &= \frac{1}{2} \left[ \frac{1}{1 - \Psi(\mathbf{R}, \mathbf{Q})} + \frac{1}{1 - \Psi(\mathbf{Q}, \mathbf{R})} \right] - 1\end{aligned}$$

Démontrons que  $\delta_1$  et  $\delta_2$  sont des indices de distance sur l'ensemble des partitions :

i)  $\Psi(\mathbf{R}, \mathbf{Q}) \in [0, 1]$

donc :  $\delta_1(\mathbf{R}, \mathbf{Q}) \in [0, 1]$

donc :  $\delta_1(\mathbf{R}, \mathbf{Q}) \geq 0$  et  $\delta_2(\mathbf{R}, \mathbf{Q}) \geq 0$

ii)  $\delta_1$  et  $\delta_2$  sont symétriques par leur définition même

iii)  $\delta_1$  et  $\delta_2$  sont nuls et seulement si  $\mathbf{R} = \mathbf{Q}$ .

En effet :

$$(\mathbf{R} = \mathbf{Q}) \Rightarrow (\Psi(\mathbf{R}, \mathbf{Q}) = \Psi(\mathbf{Q}, \mathbf{R}) = 0)$$

donc :  $(\mathbf{R} = \mathbf{Q}) \Rightarrow (\delta_1(\mathbf{R}, \mathbf{Q}) = 0 \quad \text{et} \quad \delta_2(\mathbf{R}, \mathbf{Q}) = 0)$

D'autre part :

$$(\delta_1(\mathbf{R}, \mathbf{Q}) = 0) \Rightarrow (\Psi(\mathbf{R}, \mathbf{Q}) = \Psi(\mathbf{Q}, \mathbf{R}) = 0)$$

or :  $(\Psi(\mathbf{R}, \mathbf{Q}) = 0) \Rightarrow (\mathbf{Q} \supset \mathbf{R})$

et  $(\Psi(\mathbf{Q}, \mathbf{R}) = 0) \Rightarrow (\mathbf{R} \supset \mathbf{Q})$

donc :  $(\delta_1(\mathbf{R}, \mathbf{Q}) = 0) \Rightarrow (\mathbf{Q} = \mathbf{R})$

(même démonstration pour  $\delta_2$ )—

$\delta_1$  et  $\delta_2$  satisfont donc tous les axiomes des indices de distance.

*Remarque* :  $\delta_1$  vaut 1 si et seulement si :

$$\Psi(\mathbf{Q}, \mathbf{R}) = \Psi(\mathbf{R}, \mathbf{Q}) = 0$$

c'est-à-dire :

$$(\forall \mathbf{X} \in \mathbf{R}) \quad (\forall \mathbf{Z} \in \mathbf{Q}) \quad (\mathbf{X} \cap \mathbf{Z} = \emptyset)$$

or :  $(\bigcup_{\mathbf{z} \in \mathbf{R}} \mathbf{Z}) \cap (\bigcup_{\mathbf{z} \in \mathbf{Q}} \mathbf{Z}) = \emptyset$  (et  $\delta_2$  est alors infini).

Ce qui n'arrive jamais puisque :  $\bigcup_{\mathbf{z} \in \mathbf{R}} \mathbf{Z} = \bigcup_{\mathbf{z} \in \mathbf{Q}} \mathbf{Z} = \mathbf{E}$

*Généralisation*

Il est immédiat de voir que ces indices de distances sont valables non seulement dans l'ensemble des partitions, mais aussi dans l'ensemble des familles quelconques. Leur application dans le cas particulier de l'ensemble des parties donne à nouveau les indices  $d_4$  et  $d_5$ .

### 3.3. EXEMPLE

Soit  $E = \{x_1, x_2, x_3, x_4, x_5, x_6, x_7\}$

Soit  $R = [\{x_1, x_4, x_5, x_6\}, \{x_2, x_3\}, \{x_7\}] \quad |R| = 3$

Soit  $Q = [\{x_2, x_4, x_5, x_6, x_7\}, \{x_1, x_3\}] \quad |Q| = 2$

*Distance induite*

Vecteur caractéristique :

$$\vec{V}_f(R) = [1, 2, 2, 1, 1, 1, 7]$$

$$\vec{V}_f(Q) = [1, 2, 1, 2, 2, 2, 2]$$

Toute distance entre ces vecteurs à 7 dimensions induit une distance dans l'ensemble des partitions dont l'intérêt est en général faible.

*Indice de distance*

$$\begin{aligned} \Psi(R, Q) &= \frac{1}{3} \left[ \frac{3}{6} + \frac{2}{3} + \frac{4}{5} \right] \\ &= \frac{1}{3} \left[ \frac{15 + 20 + 24}{30} \right] = \frac{59}{90} \end{aligned}$$

$$\begin{aligned} \Psi(Q, R) &= \frac{1}{2} \left[ \frac{3}{6} + \frac{2}{3} \right] \\ &= \frac{1}{2} \left[ \frac{3 + 4}{6} \right] = \frac{7}{12} \end{aligned}$$

$$\delta_1(R, Q) = \frac{1}{2} \left[ \frac{59}{90} + \frac{7}{12} \right] = 0,62$$

$$\delta_2(R, Q) = \frac{1}{2} \left[ \frac{90}{31} + \frac{12}{5} \right] - 1 = 1,2 + 1,5 = 2,7.$$

### 3.4. AUTRES EXEMPLES

Soit  $E = \{x_1, x_2, x_3, x_4, x_5, x_6, x_7\}$

Soit  $R = [E]$

Soit  $Q = [\{x_1\}, \{x_2\}, \{x_3\}, \dots, \{x_7\}]$

R et Q sont encore deux partitions de E.

$$\begin{aligned} \Psi(Q, R) &= \sum_{z \in R} \inf_{x \in Q} \left( \frac{|X \Delta Z|}{|X \cup Z|} \right) \\ &= \frac{6}{7} \end{aligned}$$

$$\Psi(R, Q) = \frac{1}{7} \left[ \sum_{x \in Q} \inf_{z \in R} \left( \frac{|X \Delta Z|}{|X \cup Z|} \right) \right]$$

$$\Psi(Q, R) = \frac{1}{7} \left[ 7 \cdot \frac{6}{7} \right] = \frac{6}{7}$$

$$\text{Donc : } \delta_1(R, Q) = \frac{1}{2} \left[ \frac{6}{7} + \frac{6}{7} \right] = \frac{6}{7} = \frac{|E| - 1}{|E|}$$

$$\delta_2(R, Q) = \frac{1}{2} [7 + 7] - 1 = 6 = |E| - 1$$

(voir annexe p. 34).

Considérons maintenant un autre exemple :

Soit  $E = \{x_1, x_2, x_3, x_4\}$

Soit  $R = [\{x_1, x_2\}, \{x_3, x_4\}]$

Soit  $Q = [\{x_1, x_3\}, \{x_2, x_4\}]$

$$\begin{aligned} \Psi(R, Q) &= \frac{1}{2} \left[ \sum_{z \in R} \inf_{x \in Q} \left( \frac{|X \Delta Z|}{|X \cup Z|} \right) \right] \\ &= \frac{1}{2} \times 2 \left[ \frac{2}{3} \right] = \frac{2}{3} < \frac{|E| - 1}{|E|} \end{aligned}$$

$$\Psi(Q, R) = \frac{2}{3}$$

$$\text{Donc : } \delta_1(R, Q) = \frac{2}{3}$$

$$\delta_2(R, Q) = \frac{1}{2} [3 + 3] - 1 = 2 < |E| - 1$$

#### IV. LES FAMILLES DE CARDINAL K

4.0. Un ensemble de  $k$  parties de  $E$  forme une famille de cardinal  $k$ . Comme pour les partitions nous déterminerons d'abord des distances à l'aide de fonction et vecteurs caractéristiques. Puis, déçus par leur manque de signification réelle, nous chercherons à établir des indices de distance plus représentatifs.

4.1. Soit  $R$  une famille de cardinal  $k$ . Ordonnons ces  $k$  parties suivant l'ordre lexicographique des éléments de  $E$  contenus :

$$R = \{A_1, A_2, \dots, A_k\}$$

Soit  $\alpha_i$  la fonction caractéristique habituelle associée à la partie  $A_i$ .

A tout élément  $x$  de  $E$ , par ces fonctions caractéristiques  $\alpha_i$ , correspondent les chiffres (0 ou 1) :

$$\alpha_1(x), \alpha_2(x), \dots, \alpha_k(x)$$

Considérons alors le nombre :

$$f(x) = \alpha_1(x) \alpha_2(x) \dots \alpha_i(x) \dots \alpha_k(x)$$

(ce nombre ne comporte que des zéros et des 1. On peut le considérer comme étant en système binaire et le traduire en système décimal. Mais cela n'a rien d'obligatoire.)

$f(\mathbf{x})$  définit une fonction caractéristique, à laquelle correspond un vecteur caractéristique.

Mais il n'y a pas de raison de privilégier ainsi  $A_i$  par rapport à  $A_j$  si  $i > j$ .

C'est pourquoi les distances induites par celles entre vecteurs caractéristiques ne seront guère intéressantes.

Cependant, l'une de ces distances induites pallie à cet inconvénient. C'est la distance : nombre de différences entre les mêmes coordonnées de deux vecteurs.

#### 4.2. EXEMPLE

Soit  $E = \{x_1, x_2, x_3, x_4, x_5, x_6, x_7\}$

Soit  $R = [\{x_1, x_2, x_3, x_4, x_5, x_7\}, \{x_6, x_7\}]$

Soit  $Q = [\{x_2, x_4, x_5, x_6, x_7\}, \{x_1, x_4, x_3\}]$

$|R| = |Q| = 2$

*Distance induite :*

Vecteur caractéristique :

$$V_f(R) = \begin{array}{c|c|c} 10 & & 2 \\ 10 & & 2 \\ 10 & & 2 \\ 10 & & 2 \\ 10 & & 2 \\ 01 & & 1 \\ 11 & & 3 \\ \hline \end{array} = \begin{array}{c} 2 \\ 2 \\ 2 \\ 2 \\ 2 \\ 1 \\ 3 \end{array}$$

↑  
système binaire

$$V_f(Q) = \begin{array}{c|c|c} 10 & & 2 \\ 01 & & 1 \\ 10 & & 2 \\ 11 & & 3 \\ 01 & & 1 \\ 01 & & 1 \\ 01 & & 1 \\ \hline \end{array} = \begin{array}{c} 2 \\ 1 \\ 2 \\ 3 \\ 1 \\ 1 \\ 1 \end{array}$$

↑  
système binaire

Toute distance entre ces vecteurs à 7 dimensions induit une distance dans l'ensemble des familles de cardinal  $k$ .

Calculons la distance : nombre de différences entre les mêmes coordonnées :

$$d_6(R, Q) = 0 + 1 + 0 + 1 + 1 + 0 + 1 \\ = 4$$

Mais cette distance présente encore un inconvénient : celui de tenir compte d'une correspondance implicite entre les parties de  $R$  et les parties de  $Q$ .

C'est pourquoi nous la conseillerons particulièrement dans le cas des familles « définies » (paragraphe 5-3).

#### 4.3. INDICES DE DISTANCE

Soient  $Q$  et  $R$  deux familles quelconques à  $k$  éléments de  $E$ . Les indices  $\delta_1$  et  $\delta_2$  (étudiés en 3.2) tenaient compte de l'erreur minimale que l'on faisait en associant à chaque élément de  $R$  un élément de  $Q$ , et à chaque élément de  $Q$ , un élément de  $R$ . Nous allons définir maintenant, avec une optique légèrement différente, des indices  $\delta_3$  et  $\delta_4$  qui se rapporteront à l'erreur minimale commise en associant bijectivement les éléments de  $R$  aux éléments de  $Q$  (ceci est possible puisque  $Q$  et  $R$  ont même cardinal  $k$ ).

Soit  $B(Q, R)$  l'ensemble des bijections de  $Q$  dans  $R$ .

Posons :

$$\Psi'(Q, R) = \inf_{g \in B(Q, R)} \left[ \sum_{x \in Q} \frac{|X \Delta g(X)|}{|X \cup g(X)|} \right] \times \frac{1}{k}$$

Il en ressort que :

(i)  $\Psi'(Q, R) = \Psi'(R, Q)$

(ii)  $\Psi'$  varie entre 0 et 1

Posons alors :

$$\delta_3(Q, R) = \Psi'(Q, R)$$

$\delta_3$  est un indice de distance pour l'ensemble des familles de cardinal  $k$ , variant entre 0 et 1.

En effet :

(i)  $\delta_3 \geq 0$

(ii)  $\delta_3 = 0 \Rightarrow \Psi'(Q, R) = 1 \Rightarrow Q \subset R$

Mais comme  $|Q| = |R| = k$

$$\delta_3 = 0 \Rightarrow Q = R$$

(iii)  $\delta_3$  est évidemment symétrique, puisque  $\Psi'$  l'est.

*Remarque :*  $\delta_3$  vaut 1 si et seulement si

$$\left( \bigcup_{x \in Q} X \right) \cap \left( \bigcup_{y \in R} Y \right) = \emptyset$$

Si au lieu de désirer un indice de distance variant entre 0 et 1, nous en voulions un variant entre 0 et l'infini, il faudrait poser :

$$\delta_4(R, Q) = \frac{1}{1 - \Psi'(R, Q)} - 1$$

(il est trivial de vérifier qu'il s'agit bien d'un indice de distance).

*Remarque :*  $\delta_4$  est infini si et seulement si

$$\left( \bigcup_{z \in Q} Z \right) \cap \left( \bigcup_{x \in R} X \right) = \emptyset$$

*Cas particulier :* si  $k = 1$ , nous nous trouvons à nouveau dans l'ensemble des parties de  $E$ , où  $\delta_4$  et  $\delta_5$  se confondent avec  $d_4$  et  $d_5$ .

4.4. Reprenons l'exemple précédent (paragraphe 4.2.)

*Indice de distance :*

$$\begin{aligned} \Psi'(Q, R) &= \frac{1}{2} \inf \left[ \left( \frac{2}{7} + \frac{5}{5} \right), \left( \frac{3}{6} + \frac{3}{5} \right) \right] \\ &= \frac{1}{2} \inf \left[ \left( \frac{9}{7} \right), \left( \frac{15 + 18}{30} \right) \right] = \frac{1}{2} \inf \left[ \frac{9}{7}, \frac{33}{30} \right] \end{aligned}$$

$$= \frac{1}{2} \frac{33}{30} = \frac{33}{60} = 0,5$$

$$\delta_3(Q, R) = \Psi'(Q, R) = 0,55$$

$$\delta_4(Q, R) = \frac{60}{27} - 1 = \frac{33}{27} = 1,2.$$

## V. AUTRES FAMILLES

### 5.1. HIÉRARCHIES INDICÉES

Nous rappelons qu'une hiérarchie indicée est le couple formé par une hiérarchie  $\mathcal{H}^1$  et une application  $\varphi$  de cette hiérarchie dans  $\mathbb{R}^+$  telle que :

$$(i) \quad (\forall X \in \mathcal{H}) \quad (|X| = 1 \Rightarrow \varphi(X) = 0)$$

$$(ii) \quad (\forall X \in \mathcal{H}) \quad (\forall Y \in \mathcal{H}) \quad (Y \supset X \Rightarrow \varphi(X) \leq \varphi(Y))$$

Nous savons d'autre part qu'une hiérarchie indicée est équivalente à une ultramétrie <sup>2</sup>.

On peut donc la représenter par un vecteur à  $n^2$  dimensions (les diverses valeurs de l'ultramétrie pour les  $n^2$  couples d'éléments de  $E$ ). Et toute distance entre ces vecteurs définira une distance dans une hiérarchie indicée.

### 5.2. FAMILLES DONT TOUS LES ÉLÉMENTS ONT POUR CARDINAL $p$

(ou bien : ensemble de  $p$ -uplets de  $E$ )

Soit  $R$  une telle famille :

$$R = \{X_1, X_2 \dots X_k\}$$

avec :

$$(\forall i \leq k) \quad (|X_i| = p)$$

Plaçons-nous dans  $E^p$  (ensemble de  $p$ -uplets d'éléments de  $E$ ). Appelons  $Y_i$  l'ensemble des permutations des éléments de  $X_i$  ( $Y_i$  est une partie de  $E^p$ ).

Posons alors :

$$Y = \bigcup_{i=1}^{i=k} Y_i; \quad Y_i = \bigcup_{i=1}^{i=k} (\text{permutations des éléments de } X_i).$$

$R$  correspond de façon injective à  $Y$ , qui est une partie de  $E^p$ .

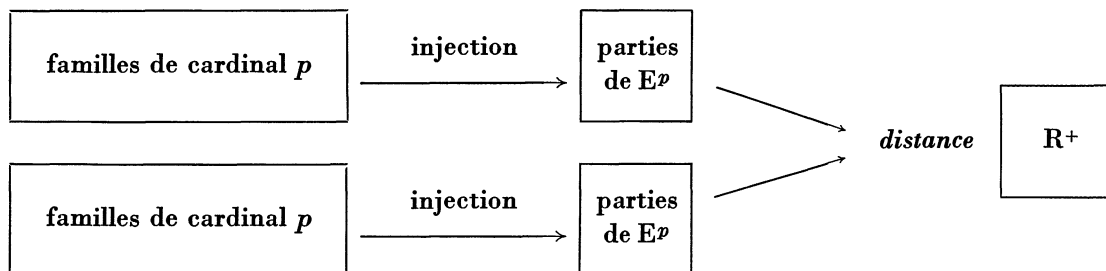
Et nous connaissons d'autre part diverses distances entre les parties ( $d_1, d_2, d_3, d_4$  et l'indice  $d_5$ ). Chacune de ces distances induit par conséquent une distance dans l'ensemble des familles de cardinal  $k$ .

---

1. Voir [8].

2. Voir [7] annexe B.

*Schéma explicatif*



*Étudions plus précisément ces distances*

Soient  $R$  et  $R'$  deux familles dont les éléments ont pour cardinal  $p$ .

Soient  $Y$  et  $Y'$  les parties de  $E^p$  correspondantes de la façon que l'on vient de voir.

$$\begin{aligned} d_1^* (R, R') &= d_1 (Y, Y') = \sqrt{|Y \Delta Y'|} \\ &= \sqrt{k! \cdot |R \Delta R'|} \end{aligned}$$

$$\begin{aligned} d_2^* (R, R') &= d_2 (Y, Y') = |Y \Delta Y'| \\ &= k! \cdot |R \Delta R'| \end{aligned}$$

$$d_3^* (R, R') = d_3 (Y, Y') = \begin{cases} 1 & \text{si } R \neq R' \\ 0 & \text{si } R = R' \end{cases}$$

$$\begin{aligned} d_4^* (R, R') &= d_4 (Y, Y') = \frac{k! |R \Delta R'|}{k! |R \cup R'|} \\ &= \frac{|R \Delta R'|}{|R \cup R'|} \end{aligned}$$

$$\begin{aligned} d_5^* (R, R') &= d_5 (Y, Y') = \frac{k! |R \Delta R'|}{k! |R \cap R'|} \\ &= \frac{|R \Delta R'|}{|R \cap R'|} \end{aligned}$$

*Remarque :* si  $k = 1$ , nous trouvons les mêmes distances entre  $R$  et  $R'$  qu'entre les parties de  $E$  formées par la réunion des éléments de  $R$ , et celle de ceux de  $R'$ .

**5.3. FAMILLES DÉFINIES**

Nous appelons «famille définie d'ordre  $k$ », le couple formé par une certaine famille  $R$  de cardinal  $k$  et une bijection  $\varphi$  de  $R$  dans  $(1, 2, \dots, k)$ .

Posons alors, pour tout couple  $(R, \varphi)$  et  $(R', \varphi')$  de familles d'ordre  $k$

$$d_1^{**} (R, R') = \sum_{i=1}^k d_1 (\varphi^{-1}(i), \varphi'^{-1}(i))$$

$$d_2^{**} (R, R') = \sum_{i=1}^k d_2 (\varphi^{-1}(i), \varphi'^{-1}(i))$$



$$d_3^{**} (R, R') = \sum_{i=1}^k d_3 (\varphi^{-1} (i), \varphi'^{-1} (i))$$

$$d_4^{**} (R, R') = \sum_{i=1}^k d_4 (\varphi^{-1} (i), \varphi'^{-1} (i))$$

$$d_5^{**} (R, R') = \sum_{i=1}^k d_5 (\varphi^{-1} (i), \varphi'^{-1} (i))$$

$d_1, d_2, d_3$  et  $d_4$  étant des distances  $d_1^{**}, d_2^{**}, d_3^{**}$  et  $d_4^{**}$  le sont aussi (les axiomes de la distance se vérifient, il est trivial de le démontrer). Enfin,  $d_5$  étant un indice de distance,  $d_5^{**}$  l'est aussi.

*Remarque* : une autre distance spécialement bien adaptée à ce cas précis est la distance  $d_6$  (voir 4.2.).

#### 5.4. FAMILLES FINIES QUELCONQUES

*Distances* : chaque famille sur  $E$  étant une partie de l'ensemble des parties de  $E$ , nous pouvons appliquer entre les familles les distances  $d_1, d_2, d_3, d_4$  ou  $d_5$  définies entre les parties.

Nous aurons ainsi, pour tout couple  $(Q, R)$  de polyparties :

$$d'_1 (Q, R) = \sqrt{|Q \Delta R|}$$

$$d'_2 (Q, R) = |Q \Delta R|$$

$$d'_3 (Q, R) = \begin{cases} 1 & \text{si } Q = R \\ 0 & \text{si } Q \neq R \end{cases}$$

$$d'_4 (Q, R) = \frac{|Q \Delta R|}{|Q \cup R|}$$

$$d'_5 (Q, R) = \frac{|Q \Delta R|}{|Q \cap R|}$$

Mais il y a  $2^n$  parties dans un ensemble de  $n$  éléments. Et, par conséquent, des familles de cardinal faible, auront très peu de chances d'avoir des éléments communs. ( $Q \cap R = \emptyset$ ).

Et nous aurons très souvent :

$$d'_1 (Q, R) = \sqrt{|Q| + |R|}$$

$$d'_2 (Q, R) = |Q| + |R|$$

$$d'_3 (Q, R) = 1$$

$$d'_4 (Q, R) = 1$$

$$d'_5 (Q, R) = \infty$$

ce qui n'est guère intéressant.

*Indices de distance* : abandonnons encore l'inégalité triangulaire et appliquons les indices de distance  $\delta_1$  et  $\delta_2$  définis dans le cas particulier des partitions par les formules :

$$\Psi (Q, R) = \frac{1}{|R|} \left[ \sum_{z \in R} \inf_{x \in Q} \frac{|x \Delta z|}{|x \cup z|} \right]$$

$$\Psi (Q, \emptyset) = 1$$

$$\delta_1 (R, Q) = \frac{1}{2} [\Psi (R, Q) + \Psi (Q, R)]$$

$$\delta_2 (R, Q) = \frac{1}{2} \left[ \frac{1}{1 - \Psi (R, Q)} + \frac{1}{1 - \Psi (Q, R)} \right] - 1$$

Ils s'appliquent sans aucune modification sur l'ensemble des familles finies quelconques.

*Exemple :*

Soit  $E = \{x_1' x_2' x_3' x_4' x_5' x_6'\}$

et quatre familles sur  $E$  :

$$Q = [\{x_1' x_2' x_3' x_4'\}, \{x_5, x_6\}]$$

$$R = [\{x_1, x_2, x_3\}, \{x_3, x_4, x_5, x_6\}]$$

$$Q' = [\emptyset]$$

$$R' = [E]$$

$$Q \cap R = \emptyset$$

$$d'_1 (Q, R) = \sqrt{2 + 2} = 2$$

$$d'_2 (Q, R) = 2 + 2 = 4$$

$$d'_3 (Q, R) = 1$$

$$d'_4 (Q, R) = 1$$

$$d'_5 (Q, R) \rightarrow \infty$$

$$\Psi (Q, R) = \frac{1}{2} \left[ \frac{1}{4} + \frac{2}{4} \right] = \frac{3}{8}$$

$$\Psi (R, Q) = \frac{1}{2} \left[ \frac{1}{4} + \frac{2}{4} \right] = \frac{3}{8}$$

$$\delta_1 (Q, R) = \frac{3}{8}$$

$$\delta_2 (Q, R) = \frac{1}{5} \left[ \frac{8}{5} + \frac{8}{5} \right] - 1$$

$$\delta_2 (Q, R) = \frac{3}{5}$$

D'autre part :  $Q' \cap R' = \emptyset$

Donc :

$$d'_1 (Q', R') = \sqrt{2}$$

$$d'_2 (Q', R') = 2$$

$$d'_3 (Q', R') = 1$$

$$d'_4 (Q', R') = 1$$

$$d'_5(Q', R') \rightarrow \infty$$

D'autre part :

$$\Psi(Q', R') = 1$$

$$\Psi(R', Q') = 1$$

Donc :

$$\delta_1(Q', R') = 1$$

$$\delta_2(Q', R') \rightarrow \infty$$

les 5 premières distances donnant donc des résultats plus grands pour Q et R que pour Q' et R' (qui paraissent cependant plus «éloignés»). Tandis que  $\delta_1$  et  $\delta_2$  donnent des résultats plus conformes à ce que nous «sentons» être la réalité de l'éloignement.

## VI. REMARQUE GÉNÉRALE

Reprenons le schéma du début de cette note et considérons les différents types de famille. Supposons que pour une famille donnée, le fait d'être du type T' implique celui d'être du type T.



Alors l'ensemble des familles de type T' est inclus dans l'ensemble des familles de type T.

Et toute distance (ou indice de distance) pour les familles de type T est aussi une distance (ou indice de distance) pour les familles de type T'.

## CONCLUSION

Nous reprenons ici, le schéma initial complété avec l'indication des distances (ou indices de distance) qui s'y appliquent (voir page 18).

Le lecteur pourra ainsi aisément voir quels sont les distances et indices de distance qui s'appliquent pour chaque type de familles (en tenant compte de la remarque précédente).

Tel était le but que nous nous proposons ici.

## ANNEXE

*Remarque générale sur les indices  $\delta_1$  et  $\delta_2$*

Soit E un ensemble quelconque à n éléments. Soient R et Q deux familles quelconques de E.

*Rappel*

$$\text{Posons : } \Psi(Q, R) = \frac{1}{|R|} \sum_{z \in R} \text{Inf}_{x \in Q} \frac{|X \Delta Z|}{|X \cup Z|}$$

On a :

$$\delta_1 (R, Q) = \frac{1}{2} [\Psi (Q, R) + \Psi (R, Q)]$$
$$\delta_2 (R, Q) = \frac{1}{2} \left[ \frac{1}{1 - \Psi (R, Q)} + \frac{1}{1 - \Psi (Q, R)} \right] - 1$$

$\delta_1$  est un indice de distance qui varie entre 0 et 1 ( $\delta_1 (\{\emptyset\}, \{E\}) = 1$ ).

Si Q et R sont des partitions, il ne pourra pas prendre la valeur 1.

Il nous semble que le maximum de  $\delta_1$  sera atteint

si  $R = \{E\}$

et  $Q = [\{x_1\}, \{x_2\} \dots \{x_n\}]$ .

On aura alors :  $\Psi (Q, R) = \frac{n-1}{n}$  ;  $\Psi (R, Q) = \frac{n-1}{n}$

$$\text{et } \delta_1 (R, Q) = \frac{n-1}{n}$$

De même l'indice de distance  $\delta_2$  qui varie entre 0 et l'infini pour des familles quelconques, reste fini entre des partitions. Il nous semble que le maximum de  $\delta_2$  pour les partitions sera atteint dans le même cas particulier que le maximum de  $\delta_1$ . Et nous aurons alors :

$$\delta_2 (Q, R) = n - 1.$$

Nous proposons donc ces deux conjectures.

Suivant ces indices, on pourrait dire que deux partitions ne peuvent être aussi « éloignées » que deux familles quelconques.

## BIBLIOGRAPHIE

- [1] BARBUT, M., MONJARDET, B., *Ordre et classification, algèbre et combinatoire*, Paris, Presses Universitaires de France, 1970, 2 tomes.
- [2] BERGE, C., *Graphes et hypergraphes*, Paris, Dunod, 1970.
- [3] CHOQUET, M., *Topologie*, Paris, Masson, 1964.
- [4] LERMAN, I.C., *Les bases de la classification automatique*, Paris, Gauthier-Villars, 1970.
- [5] ROY, B., *Algèbre moderne et théorie des graphes*, Paris, Dunod, 1969.
- [6] BRISSAUD, M., *Étude des populations et structures relationnelles*, Polycopié pour les sociologues, Faculté de Lyon, avril 1971.
- [7] BESSON, M., *Quelques méthodes d'analyses multicritères*, Thèse de 3<sup>e</sup> cycle, mars 1972, Université Paris.
- [8] ROUX, M., *Algorithme pour construire une hiérarchie particulière*, Thèse de 3<sup>e</sup> cycle, décembre 1968, Université Paris.