

N° d'ordre : 07/2008-E/MT  
REPUBLICQUE ALGERIENNE DEMOCRATIQUE ET POPULAIRE  
MINISTÈRE DE L'ENSEIGNEMENT SUPÉRIEUR ET DE LA  
RECHERCHE SCIENTIFIQUE  
UNIVERSITÉ DES SCIENCES ET DE LA TECHNOLOGIE  
« HOUARI BOUMEDIENNE »  
**FACULTÉ DES MATHÉMATIQUES**



Thèse  
Présentée pour l'obtention du diplôme de Doctorat D'Etat  
**EN : MATHÉMATIQUES**  
Spécialité: Probabilités & Statistiques  
Par : **REBBOUH Amar**  
Sujet

Analyse de l'information apportée par les éléments  
constitutifs d'une structure de tableau de données

Soutenue publiquement le 13 décembre 2008 devant le jury composé de:

Mr	BERRACHEDI	Abdelhafid	Professeur	USTHB	<b>Président</b>
Mr	DJEDOUR	Mohamed	Professeur	USTHB	<b>Directeur de thèse</b>
Mr	AIDER	Meziane	Professeur	USTHB	<b>Examineur</b>
Mr	ANES	Ouali	M/C	INPS, Alger	<b>Examineur</b>
Mr	IBAZIZEN	Mohamed	Professeur	UMTO	<b>Examineur</b>
Mr	HOUACINE	Amrane	Professeur	USTHB	<b>Examineur</b>

# Table des matières

Introduction générale	5
<b>I Structures de juxtaposition de tableaux multiples</b>	<b>8</b>
<b>1 Structure de données à observations répétées</b>	<b>8</b>
1.1 Structure d'une juxtaposition de tableaux de mesure (TofM) . . . . .	8
1.2 Structure de juxtaposition de tableaux de données qualitatives . . . . .	9
<b>2 Quelques exemples pratiques d'une structure de juxtaposition de tableaux multiples</b>	<b>9</b>
<b>3 Quelques travaux antérieurs</b>	<b>12</b>
<b>4 Problèmes posés</b>	<b>15</b>
<b>II Approches factorielles pour analyser les éléments constitutifs d'une juxtaposition de <math>n</math> tableaux à observations multiples</b>	<b>17</b>
<b>1 Analyse après réduction</b>	<b>17</b>
<b>2 <i>Les <math>L</math> observations de chaque individu pour chaque variable qui le décrit sont résumées par une seule valeur</i></b>	<b>18</b>
2.1 Cas continu . . . . .	18
2.2 Ajustement point par point . . . . .	19
2.3 Ajustement du nuage global par une droite . . . . .	19
2.4 L'une des variables est qualitative . . . . .	21
2.5 Ajustement global du nuage . . . . .	22
2.5.1 Tableau disjonctif complet associé à $X_i$ . . . . .	22
2.5.2 Tableau de contingence de Burt associé . . . . .	24
2.5.3 Critère d'ajustement et distance du $\chi^2$ . . . . .	25
2.5.4 Analyse du tableau de Burt: équivalence avec l'analyse du tableau disjonctif complet . . . . .	27
<b>3 Analyse sans réduction</b>	<b>29</b>
3.1 Comparaison globale des tableaux : l'interstructure. . . . .	29
3.1.1 Objet représentatif du tableau $X_k$ décrivant l'individu $\omega_k$ . . . . .	30
3.1.2 Choix de la métrique sur $\mathbb{R}^{L^2}$ . . . . .	31
3.1.3 Recherche du sous espace vectoriel de dimension un. . . . .	33
3.2 Le nuage moyen ou compromis : l'intrastructure . . . . .	38
3.2.1 Construction d'un compromis . . . . .	39
3.2.2 L'expression du compromis . . . . .	39

<b>III</b>	<b>Classification des éléments constitutifs d'une structure de tableaux de mesure</b>	<b>41</b>
<b>1</b>	<b>Introduction</b>	<b>41</b>
<b>2</b>	<b>Structure de juxtaposition de tableaux de données multiples</b>	<b>42</b>
<b>3</b>	<b>Indices de Distances</b>	<b>45</b>
3.1	Indice basé sur la distance de Kullback-Leibler . . . . .	45
3.1.1	Cas discret . . . . .	45
3.1.2	Cas continu . . . . .	46
3.2	Indice de distance basé sur des techniques factorielles . . . . .	47
3.3	Indice basé sur le produit scalaire de Hilbert-Schmidt . . . . .	49
3.3.1	Le produit scalaire de Hilbert-Schmidt . . . . .	50
3.3.2	Norme et distance induite par le produit scalaire de Hilbert-Schmidt	50
<b>4</b>	<b>Critère et problème d'optimisation</b>	<b>51</b>
4.1	La fonction de représentation . . . . .	52
4.2	La fonction d'affectation . . . . .	53
<b>5</b>	<b>L'algorithme</b>	<b>54</b>
<b>6</b>	<b>Structure générale de l'algorithme : Organigramme</b>	<b>54</b>
<b>7</b>	<b>Application</b>	<b>55</b>
7.1	Données en entrée . . . . .	55
7.2	Classification après réduction . . . . .	56
7.3	Classification basé sur l'indice de Hilbert-Schmidt . . . . .	57
7.4	Interprétation de ces résultats à l'aide des courbes représentatives des centres ou noyaux des classes . . . . .	58
<b>8</b>	<b>Conclusion</b>	<b>60</b>
<b>IV</b>	<b>Classification des éléments constitutifs d'une structure de tableaux de données catégorielles</b>	<b>61</b>
<b>1</b>	<b>Introduction</b>	<b>61</b>
<b>2</b>	<b>Formalisme adapté</b>	<b>61</b>
<b>3</b>	<b>Distance entre systèmes physiques aléatoires multiples</b>	<b>63</b>
3.1	Entropie comme mesure de l'incertitude des états des systèmes . . . . .	63
3.2	Entropie d'un système physique multiple aléatoire . . . . .	64

<b>4</b>	<b>Application numérique</b>	<b>67</b>
4.1	Procédure pour estimer la distribution conjointe . . . . .	67
4.2	Exemple numérique . . . . .	68
4.2.1	Procédure pour construire une hiérarchie indicée sur ces objets . . .	69
4.3	<b>Détails sur le déroulement de cet algorithme</b> . . . . .	71
4.4	Application à des données réelles . . . . .	72
<b>5</b>	<b>Conclusion</b>	<b>74</b>
<b>V</b>	<b>Classification des données avec erreurs de mesure</b>	<b>75</b>
<b>1</b>	<b>Introduction</b>	<b>75</b>
<b>2</b>	<b>Estimation des sous ensembles <math>R_i</math></b>	<b>76</b>
2.1	Les $d$ variables sont non corrélées . . . . .	77
2.1.1	Les intervalles de confiance associés . . . . .	77
2.1.2	Premier cas $\sigma_i^j$ sont connus pour $i = 1, \dots, N$ et $j = 1, \dots, d$ . .	77
2.1.3	Deuxième cas $\{\mu_i^j ; j \succeq 1\}$ sont connus et $\sigma_i^j \neq 0 ; \forall j \succeq 1$ . . . . .	78
2.2	Les $d$ variables sont quelconques . . . . .	79
2.2.1	Première approche: Choix du sous espace compromis et estimation des intervalles de confiance . . . . .	79
2.2.2	Seconde approche: Estimation des couples $(\mu_i, \Sigma_i)$ . . . . .	81
<b>3</b>	<b>Choix de l'indice de distance entre éléments représentatifs</b>	<b>81</b>
3.1	Les individus sont décrits par des pavés de $\mathbb{R}^d$ . . . . .	81
3.1.1	Indice de distance entre objets intervalles . . . . .	82
3.1.2	Indice d'aggrégation entre objets et pavés de $\mathbb{R}^d$ . . . . .	82
3.2	Les individus sont décrits par des ellipsoïdes de $\mathbb{R}^d$ . . . . .	82
3.2.1	Distance de Mahannalobis pondérée. . . . .	82
<b>4</b>	<b>Classification d'objets décrits avec erreurs de mesure</b>	<b>83</b>
4.1	Objets décrits par des pavés de $\mathbb{R}^d$ . . . . .	83
4.1.1	<b>Procédure</b> . . . . .	84
4.1.2	Problème d'optimisation . . . . .	85
4.1.3	Algorithme de classification . . . . .	86
4.2	Etude de cas . . . . .	87
4.2.1	Alternative 1: Les erreurs sont cummulées au niveau des noyaux des classes . . . . .	87
4.2.2	Alternative 2: Les erreurs sont cumulées au niveau de la description des individus et les noyaux sont donnés sans erreurs . . . . .	90
4.2.3	Alternative 3: Les erreurs sont cumulées dans la description des individus et dans les noyaux des classes . . . . .	92
4.3	Objets décrits par des ellipsoïdes de $\mathbb{R}^d$ . . . . .	95
4.3.1	Critère et problème d'optimisation . . . . .	95
4.3.2	Caractérisation d'une classe d'objets(résultat très important) . . . .	96

	4
4.3.3 La fonction de représentation . . . . .	98
4.3.4 La fonction d'affectation . . . . .	99
4.3.5 L'algorithme . . . . .	99
4.3.6 Exemples . . . . .	100
<b>VI Conclusion générale</b>	<b>103</b>
<b>VII Bibliographie</b>	<b>105</b>
<b>1 Appendix</b>	<b>113</b>

# Chapitre

## Introduction générale

Cette thèse s'inscrit dans le cadre de la classification automatique et plus précisément dans la classification des éléments constitutifs d'une structure de tableaux multiples. L'analyse des éléments constitutifs d'une juxtaposition de tableaux multiples reste un large domaine d'investigation. L'objectif principal de la classification automatique et de l'analyse des données est de définir des indices de similarité entre individus ou entre variables qui permettent d'associer à un couple d'objets ou de variables une valeur numérique. Ces indices de mesures sont choisis en fonction des données et de la structure de classification désirée.

Nous apportons dans cette thèse quelques nouveaux éclairages qui permettent d'appréhender l'information apportée par les éléments constitutifs d'une structure de juxtaposition de tableaux multiples. Une structure de juxtaposition de tableaux multiples peut être obtenue dès que l'on est conduit à observer plusieurs fois chaque objet pour les paramètres qui le décrivent. Aussi, nous présentons plusieurs approches qui permettent :

- *De comparer globalement des différentes variables*
- *D'étudier les différents individus.*
- *De faire des typologies d'individus suivant des critères adaptés à la structure de données se présentant sous la forme d'une juxtaposition de tableaux à observations multiples.*

La première partie de la thèse est une rétrospective des techniques classiques du type factorielles pouvant s'adapter à ce type de données. Dans la seconde partie, nous définissons trois indices qui permettent sous certaines conditions, de mesurer la similarité entre les éléments constitutifs de la structure. Dans la troisième partie, nous introduisons la notion d'objets aléatoires et de systèmes physiques associés pour modéliser les observations de chaque individu pour les variables qui les décrivent. La quatrième partie est consacrée à la classification d'objets décrits par des données entachées d'erreurs de mesure. Enfin, la dernière partie de la thèse regroupe certaines perspectives, ouvertures et développements futurs.

Si les algorithmes de classification automatique existants se distinguent pour la plupart par le choix de l'indice mesurant la ressemblance ou la dissemblance entre objets. Lermam[62] note que le problème de ce choix est délicat et est loin d'être résolu. Cette remarque reste hélas vraie à l'heure actuelle. Aussi, dans cette thèse, nous avons défini 5 indices qui sont

- L'indice basé sur le produit scalaire de *Hilbert-Schmidt* et sur la définition d'objets représentatifs de chaque élément constitutif de la structure. Cet indice a été adapté à une structure de tableaux de mesure (*TofM*) et un algorithme de type k-means McQueen[65] a été validé pour classifier les clients domestiques d'une entreprise de distribution d'électricité.

- L'indice basé sur des techniques factorielles, adapté dans le cas d'une structure de juxtaposition de tableaux de mesure.
- L'indice basé sur les techniques utilisées en théorie de l'information (système physique, entropie,...) et sur les travaux de Kullback-Leibler[55].
- La distance classique euclidienne entre sommets des pavés qui décrivent les individus dans le cas de la classification de données avec erreurs de mesure et si les variables ayant participé à la description des individus sont non corrélées.
- Si les variables sont corrélées donc quelconques, nous utilisons la distance pondérée de *Mahannalobis*.

Dans la troisième partie, nous introduisons la notion d'objets aléatoires. Les objets aléatoires sont une modélisation des observations de chaque individu pour les variables qui les décrivent. Cette approche utilise la notion d'entropie comme mesure de l'incertitude des états du système associé à chaque élément de la structure. Le concept de système physique est celui proposé par Shannon-Weininger[90]. Cette notion fait encore l'objet de discussions entre scientifiques de différentes spécialités (chimie, théorie du signal, statistique, ...)

Depuis 1824, date à laquelle **Carnot** dans "*réflexions sur la force motrice du feu*" a ébauché les principes fondamentaux d'une nouvelle science : **La thermodynamique**. Cette science est basée sur l'existence d'une grandeur physique fondamentale baptisée : **entropie**.

Le principe que chacun des "systèmes physiques" possède une entropie qui lui est propre et qui représente son degré de "*désordre*" a été admis depuis. Cependant, c'est sur la définition d'un système physique que les discussions entre scientifiques demeurent.

Le fait de considérer qu'une observation est une réalisation d'une variable aléatoire conduit à raisonner sur l'ensemble des états ou valeurs que peut prendre cette variable. Aussi, si la variable est discrète et le nombre d'états fini, chaque état sera mesuré par sa fréquence. Il s'agit d'étudier la probabilité que la variable ou le système associé d'être dans une des configurations possibles ou états. On se ramène donc à la notion de système physique définie par Ludwig Boltzman [1872]

Cette approche qui consiste en la modélisation des éléments constitutifs ou observations répétées de la structure par des systèmes physiques s'étend à des classes d'individus donc à un ensemble beaucoup plus vaste que celui étudié. Nous signalons que ces notions de systèmes physiques et d'entropie ont été développées de manière plus générale par Shanonn[90].

Le formalisme proposé donne une explication satisfaisante de la distance introduite par Kullback -Leibler. De plus, l'indice adapté permet de construire une hiérarchie indicée sur les éléments constitutifs d'une structure de juxtaposition de tableaux multiples et s'étend à de multiples applications et structures dans le cas le plus général.

Lerman [62] a proposé l'indice du lien de vraisemblance maximum qu'il a adapté à une structure de juxtaposition de tables de contingence mais cet indice et les nombreux indices proposés dans la littérature ne s'adaptent pas à notre structure. Parmi ces indices, on peut citer: l'indice de Ward, l'indice de Sokal & Sneath, l'indice d'Anderberg, l'indice de Cover,...Pour plus de précisions voir [7]

En résumé, cette thèse est découpée en 5 chapitres.

Dans le chapitre I nous introduisons les types de structure de données qui nous intéressent et nous présentons l'originalité et l'intérêt pratique.

Dans le chapitre II, nous étudions les différentes approches de type factorielles qui peuvent s'y adapter. Le cas d'une structure de juxtaposition de tableaux de mesure est largement étudié.

Dans le chapitre III, nous définissons 3 indices de ressemblance adaptés aux éléments constitutifs de la structure et détaillons les différentes étapes de l'algorithme des moyennes mobiles ou k-means[40]. L'algorithme utilisant la distance induite par le produit scalaire de *Hilbert-Schmidt* est programmé en MATLAB 5.3 et utilisé pour classifier les clients domestiques d'une entreprise de distribution d'électricité. Ce chapitre a fait l'objet d'une publication dans la revue *communications in statistics*.

Le chapitre IV de cette thèse constitue une nouvelle approche tant au niveau du formalisme que des concepts définis. Ainsi, nous introduisons la notion d'objets aléatoires et adaptons de manière naturelle un indice de distance. Cet indice est apparu, de manière presque inattendue, comme une nouvelle version de la distance de *Kullback-Leibler* définie pour d'autres applications. Les différentes étapes de l'algorithme de classification hiérarchique sont développées, un exemple déroulé manuellement et une application sur données réelles complètent ce chapitre dont une synthèse a fait l'objet d'une publication dans "*Journal of Applied Mathematics and Decision Science*"

Le chapitre V est une extension des résultats établis dans le chapitre III vers une application à un problème classique et important . Il s'agit de la classification d'objets dont les descriptions sont entachées d'erreurs de mesure . Des approches sont proposées et validées sur des exemples.

Une conclusion, certaines perspectives de développement et une bibliographie commentée terminent ce travail.



## Chapitre I

# Structures de juxtaposition de tableaux multiples

Ces structures de données peuvent être obtenues dès que l'on est conduit à effectuer plusieurs observations sur les variables qui décrivent chaque individu. Les individus peuvent ne pas être décrits par le même groupe de variables et le nombre d'observations peut différer d'un individu à l'autre ou d'une variable à l'autre.

## 1 Structure de données à observations répétées

C'est une structure de données particulière de tableaux multiples fréquente en pratique. Elle se schématise comme suit

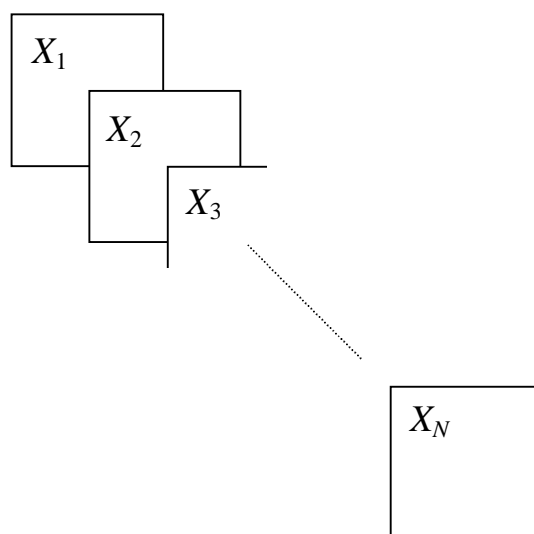


fig 1

### 1.1 Structure d'une juxtaposition de tableaux de mesure (TofM)

Soit  $\Omega$  un ensemble fini de  $N$  individus, le tableau

$$X = [X_1, \dots, X_N],$$

est constitué de la juxtaposition de  $N$  sous tableaux où

$$X_i = [x_{il}^j]_{j=1, d_i=\|P_i\|; l=1, L} ,$$

$X_i$  est le sous tableau de dimension  $L \times d_i$  contenant les observations répétées de l'individu  $\omega_i \in \Omega$  pour le groupe de variables quantitatives (continues)  $P_i$  de cardinal  $d_i$  qui le décrivent  $x_{il}^j$  est la  $l^{\text{ième}}$  observation de l'objet  $\omega_i$  pour la variable  $V^j$ . Une observation  $\theta_i$  est une ligne du tableau  $X$ . À cette observation dite « moyenne » est associée  $N$  individus  $\{\hat{i}; i = 1, N\}$  dits partiels correspondants aux divers sous tableaux  $X_i$ .

## 1.2 Structure de juxtaposition de tableaux de données qualitatives

Soit  $\Omega$  un ensemble fini de  $N$  individus,  $T$  un descripteur de la forme

$$T = [T_1, \dots, T_N],$$

où

$$T_i = [m_{il}^j]_{j=1, d_i=\|P_i\|; l=1, L}$$

$T$  est un tableau constitué de la juxtaposition de  $N$  sous tableaux  $\{T_i; i = 1, n\}$ ,

$T_i$  est le sous tableau de dimension  $L \times d_i$  contenant les observations répétées de l'individu  $\omega_i \in \Omega$  pour le groupe de  $P_i$  variables qualitatives qui le décrivent,  $m_{il}^j$  est la  $l^{\text{ième}}$  modalité prise par l'objet  $\omega_i$  pour la variable  $V^j$ .

Le tableau  $T_i$  regroupe les modalités prises par l'individu  $\omega_i$  pour les variables qui le décrivent.

Cette structure est une structure de données brutes non exploitable sous cette forme pour son analyse, on adopte la démarche classique utilisée dans l'analyse factorielle des correspondances multiples qui consiste à associer à chacun des sous tableaux, le tableau disjonctif complet correspondant. Nous obtenons ainsi, une structure de juxtaposition de tableaux de données catégorielles. L'analyse de cette structure fera l'objet du chapitre IV de cette thèse.

## 2 Quelques exemples pratiques d'une structure de juxtaposition de tableaux multiples

Les exemples de structure de juxtaposition de tableaux de ce type sont nombreux en médecine où  $X_1, \dots, X_N$  sont des dossiers de malades regroupant les bilans effectués en

mesurant les différents paramètres et symptômes qui décrivent la maladie. Les tableaux n'ont pas la même dimension car le choix des paramètres et leurs mesures dépendent de l'état général du malade ou de la maladie.

En écologie végétale, on rencontre souvent cette structure de données particulière. En effet

Si  $R_1, \dots, R_L$  sont  $L$  relevés de la région étudiée,  $e_i^l$  l'espèce végétale  $i$  présente dans le relevé  $R_l$ ,  $v_j^l$  le paramètre  $j$  mesuré sur son cortège floristique.

La structure de données dont on dispose en entrée se présente comme suit

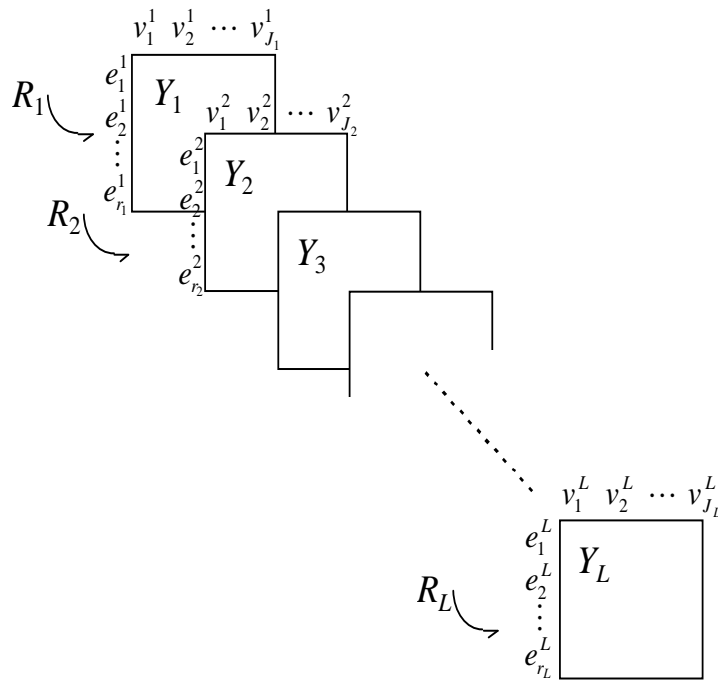


fig 2

Chaque relevé  $R_i$  est décrit par

$$Y_i \wedge V_i$$

où  $Y_i$  est le tableau contenant les mesures faites sur son cortège floristique, est le vecteur contenant les mesures des paramètres édaphiques ou paramètres propres au relevé.

# Un des problèmes posé est l'étude des liens pouvant exister entre le cortège floristique et les paramètres édaphiques donc entre les tableaux  $Y_i$  et les vecteurs  $V_i$  #.

Les paramètres édaphiques sont hétérogènes (qualitatifs et quantitatifs). Dans le cas où l'on pourrait les recoder en un paramètre qualitatif avec plusieurs modalités, ce problème devient un problème type que l'analyse discriminante peut résoudre.

En effet, soit

$$V = \begin{bmatrix} V_1 \\ \vdots \\ V_L \end{bmatrix}$$

le tableau contenant les descriptions des paramètres édaphiques des différents relevés.

Une typologie des relevés construite à partir de  $V$  conduit à recoder les paramètres par une variable qualitative dont les modalités sont les classes obtenues. Dans ce cas, on retrouve les conditions d'application de l'analyse factorielle discriminante. Cependant, cette situation n'est pas toujours possible et les résultats sont difficilement interprétables. De plus, dans la pratique, nous n'avons pas les mêmes espèces dans les différents relevés et les relevés sont caractérisés par d'autres paramètres tels que les caractéristiques du sol: calcium, ensoleillement, pentes, etc... De plus, les variables ne sont pas les mêmes pour chaque tableau

$$\{Y_l ; l = 1, L\}$$

A partir de cette structure, on peut construire une structure du type précédent. En effet

Si  $e_1, \dots, e_N$  sont les différentes espèces rencontrées, on construit les tableaux  $X_1, \dots, X_N$  caractérisant respectivement chacune des espèces de la manière suivante.

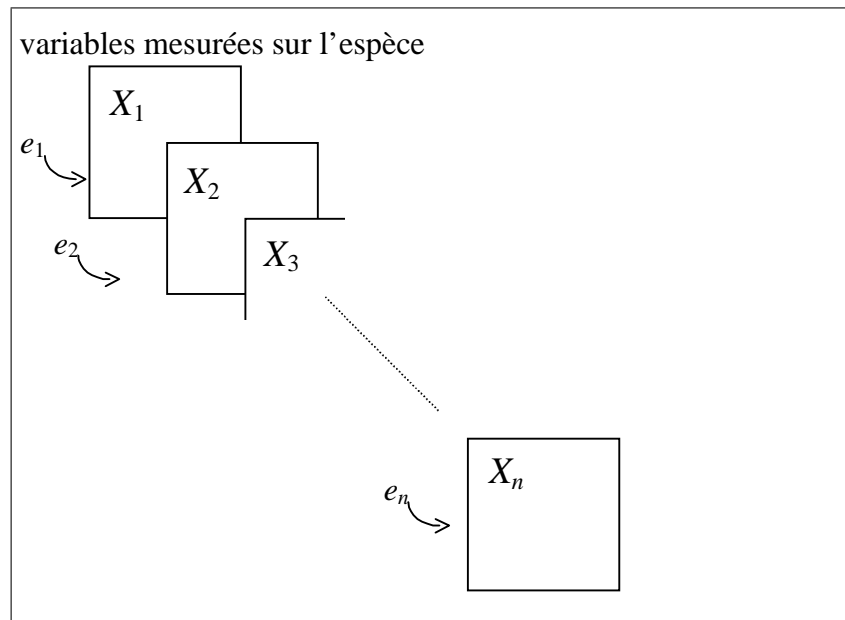


fig 3

Comme les espèces ne sont pas toutes présentes dans les différents relevés, les tableaux  $X_1, \dots, X_N$  n'ont pas la même dimension.

La structure de données obtenue est une structure de juxtaposition de tableaux multiples. Si le problème de l'analyse de tableaux multiples par des techniques factorielles

reste un vaste domaine de recherche Lebart, Morineau et Pion[58], plusieurs modèles et approches sont proposés pour comparer les tableaux, décrire leur structure commune et appréhender leur différence.

### 3 Quelques travaux antérieurs

Le fait que le théorème d'Eckart et Young[35][1935] de décomposition en éléments singuliers permettant de reconstituer partiellement ou totalement la structure de départ n'ait pas de généralisation dans le cas de tableaux multiples, a conduit certains chercheurs à développer des modèles particuliers aux disciplines et à la nature de leurs données.

On peut citer les premiers travaux sur l'analyse de tableaux multiples par les techniques factorielles qui remontent à Tucker[97][1964] et [98][1966], Harschman[45][1970]. On cite dans le même répertoire les travaux de Koomerberg[53][1983]; Carlier[14][1985]; de Von Der Heijder[101] ; Carlier et Coll[15][1988]; ...

Cuttman[26][1941] ; Burt[12][1950] et Hayashi[48][1956] sont les précurseurs de l'ACM " Analyse des correspondances multiples ". Des extensions ont été développées par Escofier et Cordier[1965]; Benzekri[7][1973]; Horst[49][1961]; Kettenring[50][1971];... L'analyse des correspondances multiples a été également développée sous le nom de " *Homogeneity analysis* " et sous le nom de " *Dual Scaling* " par Nishisato[70][1980]. La première application de l'ACM à un tableau disjonctif complet est l'œuvre de JP Nakache[69].

Les résultats et propriétés connus actuellement ont été mis en forme et programmés par Lebart et Tabard[1973]. Enfin, un exposé synthétique de certaines techniques a été réalisé par Tenenhaus et Young[1985]; Lebart, Morineau et Piron[58].

Dans le cas où les tableaux  $X_1, \dots, X_N$ , regrouperaient les mêmes individus décrits par des groupes de variables différentes, l'AFC peut répondre au problème posé.

Escofier et L'hermier des Plantes[63] ont introduit la méthode STATIS " Structure des Tableaux À Trois Indices de la Statistique ". Développée par la suite par Lavit[56]. Cette méthode a pour but d'analyser un ensemble d'individus décrits par plusieurs groupes de variables ou d'analyser un ensemble de variables mesurées sur plusieurs groupes d'individus.

Dans le cas où notre structure se présente sous la forme

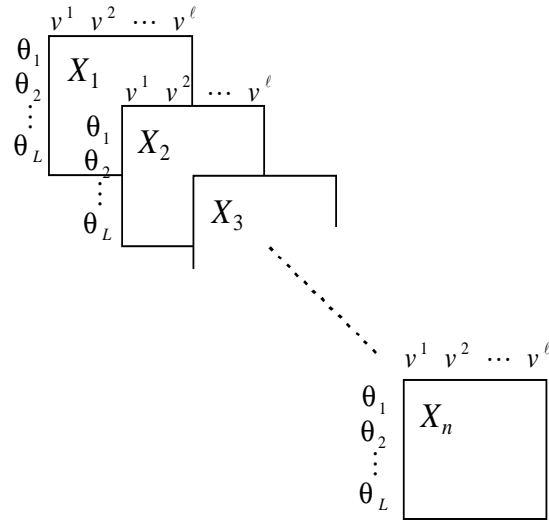


fig 4

où  $\theta_1, \dots, \theta_L$  sont considérés comme un ensemble de  $L$  individus décrits respectivement par les tableaux  $X_1, \dots, X_N$  de même dimension, la méthode STATIS permet effectivement d'étudier l'évolution des différents  $\{\theta_i; i = 1, L\}$ , d'obtenir les liaisons entre les variables  $\{V^j; j = 1, p\}$  et de trouver une structure commune aux différents tableaux.

Le chapitre II de ce travail est consacré à des approches factorielles qui peuvent être adaptées et étendues aux données écologiques si  $\theta_1, \dots, \theta_L$  représentent les relevés  $X_1, \dots, X_N$  les tableaux décrivant les espèces. Ce sera une étude particulière qui nécessite que les mêmes espèces soient présentes dans les différents relevés, ce qui n'est pas la réalité qu'on désire appréhender.

Une autre approche factorielle peut être utilisée dans le cas où dans la structure toutes les variables seraient quantitatives et les tableaux de même dimension.

Cette approche DACP " *Double Analyse en Composantes Principales* " introduite par Bouroche[10] s'adapte parfaitement mais l'interprétation des résultats n'est pas évidente ce qui limite considérablement son application. De plus, elle se repose essentiellement sur la recherche du référentiel commun de représentation, qui demeure, dans l'optique des techniques proposées, un problème assez ouvert.

Cette approche s'articule autour de 3 phases distinctes.

La première étape consiste à l'analyse du phénomène d'évolution globale. Cette évolution est étudiée par une analyse en composantes principales des centres de gravité des nuages associés aux  $L$  tableaux  $Y_1, \dots, Y_N$ . Elle correspond à ce qui est appelé: *l'étude de l'interstructure*.

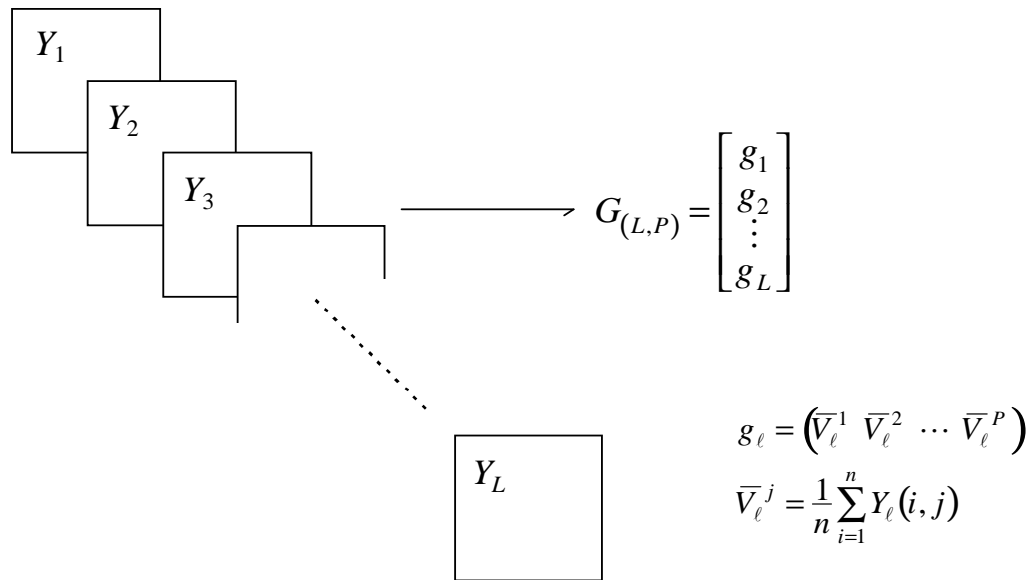


fig 5

Chaque tableau est résumé par le vecteur contenant le centre de gravité du nuage associé. L'ACP du tableau  $G$  permet d'étudier l'évolution des relevés (dans le cas de données écologiques) par l'intermédiaire des centres de gravité.

La seconde étape consiste à étudier la déformation des nuages autour de leur centre de gravité. Pour ce faire, Bouroche[10] propose d'effectuer  $L$  analyses en composantes principales (ACP) des  $L$  nuages de points associés aux tableaux centrés. Ainsi, cela permet d'éliminer le phénomène d'évolution globale.

La troisième phase consiste à rechercher un espace de représentation des espèces commun Bouroche [10] propose 4 procédures heuristiques permettant de déterminer un référentiel commun.

Cette étape correspond à ce qui est appelé : *étude de l'intrastructure*.

Outre les inconvénients propres aux différentes approches pour analyser des tableaux multiples la structure de données dont on dispose n'obéit pas, dans le cas général, aux hypothèses qui sous-tendent ces techniques. Les tableaux juxtaposés n'ont ni le même nombre de lignes, ni le même nombre de colonnes.

Cette structure de données en entrée n'a pas eu non plus, l'attention en classification automatique, à la hauteur de sa fréquence dans les cas pratiques.

Plusieurs structures de classification automatique existent, les plus courantes sont les partitions, hiérarchies, arbres et pyramides. Pour obtenir une partition à  $K$  classes ( $K$  fixé), la pratique courante consiste à optimiser un critère basé sur l'inertie. L'espace où sont définis les objets est muni d'une métrique  $M$ . Une solution est obtenue par un algorithme itératif du type k-means ou centres mobiles, Thordike[96][1953]; Forgy[40][1965]; Ball and Hall[5][1967] ; Macqueen[65]; Diday[31], ...

Il existe plusieurs variantes de ces algorithmes qui se distinguent principalement par certains aspects tels que : les noyaux, centres ou représentants des classes et les métriques sur l'espace où sont définis les objets à classer Benzecri[7][1973] ; Anderberg[2][1973]; Diday[31][1980]; ...

Dans le cas où chaque tableau  $X_i$  de la structure regrouperait les valeurs multiples et non prévisibles à priori, prises par les variables qui décrivent l'objet correspondant et si les objets sont décrits par le même groupe de variables, une modélisation et une approche classification par des partitions dont le critère est basé sur la distance induite par le produit scalaire de *Hilbert-Schmidt* sont proposées dans le chapitre III . Un programme informatique en MATLAB 5.3 est déroulé pour classer les clients d'une entreprise de distribution d'électricité.

La structure de données se présente dans ce cas

$$\begin{array}{l}
 \text{A} \\
 \text{l'objet} \\
 \omega_i \\
 \text{correspond} \\
 \text{le} \\
 \text{tableau} \\
 X_i
 \end{array}
 \rightarrow
 \begin{array}{c}
 \theta_1 \\
 \theta_2 \\
 \vdots \\
 \theta_{L_i}
 \end{array}
 \left| \begin{array}{cccc}
 V^1 & V^2 & \dots & V^P \\
 \times & \times & \dots & \times \\
 \times & \times & \dots & \times \\
 & & & \\
 \times & \times & & \times
 \end{array} \right.
 \begin{array}{l}
 \rightarrow E_1^i \\
 \rightarrow E_2^i \\
 \\
 \rightarrow E_P^i
 \end{array}$$

tab 1

Cette modélisation exige que les variables soient observées le même nombre de fois pour chaque objet qu'elles décrivent. Les tableaux  $X_1, \dots, X_n$  ont donc le même nombre de lignes.

Ce cas particulier, même s'il ne s'étend pas au cas général, est proche de la réalité car les contraintes exigées (même groupe de variables pour les objets) sont fréquentes dans toutes les techniques classiques d'estimation. Les lignes de chacun des tableaux, correspondent aux observations des objets pour les variables qui les décrivent. Ces observations doivent être prises dans les mêmes conditions si l'on veut que les comparaisons aient un sens.

Enfin, la distance de Kullback-Leibler s'adapte au cas d'une structure où les tableaux n'ont pas le même nombre de lignes mais doivent avoir le même nombre de colonnes càd qu'il est nécessaire que les individus soient décrits par les mêmes variables.

## 4 Problèmes posés

Lerman[60] [62] a proposé l'algorithme du lien de vraisemblance qu'il a adapté pour traiter les éléments constitutifs d'une structure de données sous forme d'une juxtaposition de tables de contingence mais la structure dont on dispose, nouvelle par sa forme ne pas être transformée pour être appréhendée par l'approche de Lerman[62].

Il s'agit donc de définir une autre manière de décrire les objets. La description devrait tenir compte



1. *De la variabilité des observations.*
  2. *Du caractère aléatoire des observations.*
  3. *Du fait que les objets ne sont pas décrits par le même groupe de variables.*
  4. *Et enfin du fait que les objets ou variables ne sont pas observés le même nombre de fois.*
- Il faut de plus, que cette description offre des possibilités graphiques et de simulation.

## Chapitre II

# Approches factorielles pour analyser les éléments constitutifs d'une juxtaposition de $n$ tableaux à observations multiples

## 1 Analyse après réduction

On suppose dans ce paragraphe que les objets sont décrits par le même groupe de variables et observés le même nombre de fois. Soit  $X$  la structure de données suivante

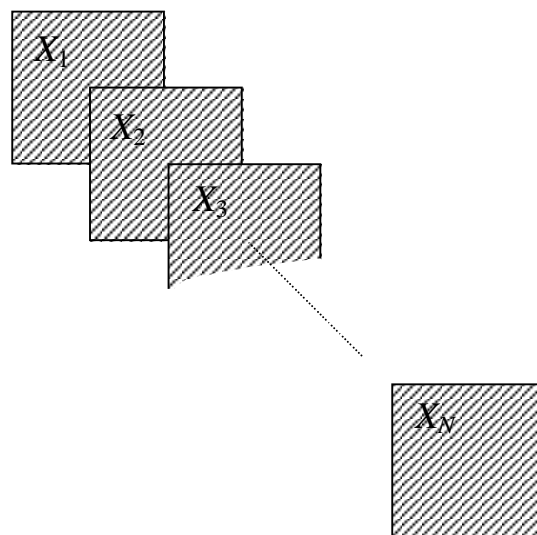


fig 6

où les  $n$  blocs de sous tableaux ont la même dimension.

Pour tout  $i = 1, \dots, n$ ;  $X_i$  est un tableau de dimension  $(L \times d)$  tel que



Soit  $x_i^j$  la  $j$ ème colonne de cette matrice, on désire résumer cette colonne par une seule valeur. Deux approches sont possibles pour répondre à ce problème

- Faire un ajustement colonne par colonne pour chaque tableau.
- Faire un ajustement global du nuage associé aux  $d$  colonnes par un sous-espace affine ou un sous-espace vectoriel de dimension 1 si le nuage est centré.

Cette dernière technique est la première étape de la méthode *DACP* (Double Analyse en Composantes Principales) introduite par Bouroche[10] et consiste à l'analyse du phénomène d'évolution globale des différents tableaux.

## 2.2 Ajustement point par point

Dans ce cas,  $X_i^j$  est une série statistique de  $L$  observations de la variable  $V^j$ . Cette série peut être représentée géométriquement par un vecteur de  $\mathbb{R}^L$ . La technique d'usage pour résumer la série par un seul point nécessite le choix d'une métrique sur et la résolution du problème d'optimisation lié à cette métrique

$$\text{Min} \{ D^2 (X_i^j, t.J_L), t \in \mathbb{R} \}$$

$J_L$  est le vecteur de  $\mathbb{R}^L$  dont toutes les composantes valent 1. Si  $D$  est la métrique euclidienne classique, l'optimum  $t_{ij}^*$  est atteint pour la moyenne empirique de la série.

D'autres métriques conduisent à la médiane, au mode de la série comme solution de ce problème. En général, ce problème n'admet pas toujours de solution.

Le choix de la métrique repose sur la forme de la distribution ou de la fonction de répartition représentée graphiquement par un histogramme ou par un diagramme en bâtons. La découpe en classes de la série pose le problème fondamental qui est "*l'effet partition*": un changement de la partition de l'espace des observations peut modifier le graphe de l'historgramme, gommer ou faire apparaître des symétries ou des modes... De plus, il est difficile de mesurer la qualité de l'ajustement. Enfin, cet ajustement ne permet pas de reconstituer même partiellement la série (les données de départ) car l'application qui associe à la série sa valeur centrale n'est pas injective.

## 2.3 Ajustement du nuage global par une droite

Le tableau  $X_i$  peut être considéré comme un nuage de  $d$  points de  $\mathbb{R}$

$$X_i = [X_i^1, \dots, X_i^d]$$

.Sans restreindre la généralité, supposons que les variables sont centrées et que  $\mathbb{R}^L$  est muni de la métrique euclidienne classique. On a

$$\bar{X}_i^j = \frac{1}{L} \sum_{t=1}^L x_{it}^j = 0$$

Le problème posé est un problème de réduction numérique. Autrement dit, un problème de compression des données. Parmi les critères d'ajustement d'un nuage de  $d$  points par un sous-espace vectoriel, celui retenu et qui conduit probablement aux calculs analytiques les plus simples est le critère classique des moindres carrés. Il consiste à rechercher dans notre cas, la droite d'allongement maximum du nuage de points et donc à rendre minimale la somme des carrés des écarts

$$\sum_{j=1}^d M_j H_j,$$

$M_j$  et  $H_j$  sont deux points  $\mathbb{R}^L$  de tels que

$$\overrightarrow{OM_j} = X_i^j,$$

$H_j$  est la projection orthogonale de  $M_j$  sur  $\Delta_u$  # l'axe engendré par le vecteur unitaire  $\vec{u}$  # et passant par l'origine centre de gravité du nuage "*condition nécessaire pour qu'une droite satisfasse le critère des moindres carrés.*"

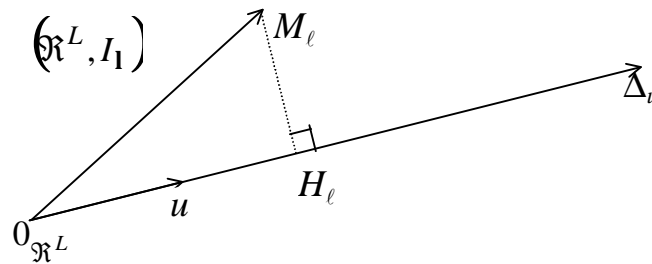


fig 7

La droite  $\Delta_u$  qui réalise le meilleur ajustement du nuage des  $d$  points de  $\mathbb{R}^L$  au sens des moindres carrés est la droite passant par l'origine de  $\mathbb{R}^L$  engendrée par le vecteur propre unitaire  $U_1^*$  associé à la plus grande valeur propre  $\lambda_1$  de la matrice carrée d'ordre  $d$ ,

$$X_i' \cdot X_i.$$

La qualité d'ajustement se mesure par

$$T_1 = \frac{\lambda_1}{d} = \frac{\lambda_1}{\sum_{t=1}^d \lambda_t} = \frac{\lambda_1}{\text{tr}(X_i' \cdot X_i)}$$

Dans ce cas, chaque vecteur  $X_i^j$  pour tout  $j = 1, d$  est résumé par sa projection sur la droite

$$X_i^j \longrightarrow y_i^j = (X_i^j)' \cdot u \in \mathbb{R}.$$

Plusieurs travaux ont été menés pour établir la loi de la matrice

$$S = X_i' \cdot X_i$$

dans le but d'expliciter la loi de ses valeurs propres. Sous certaines hypothèses, il a été montré que  $S$  suit la loi de *Wirshart*. Fisher[39][1939] a établi la densité de probabilité des valeurs propres issues d'une matrice de *Wirshart*. Dans Anderson[3] sont données la densité de probabilité de *Wirshart* et de certaines lois dérivées.

Il est apparu que l'utilisation du taux d'inertie  $T_1$  comme outil d'évaluation globale de la qualité de représentation est très délicat.  $T_1$  représente une part de la variance brute initiale qui n'est pas une mesure de référence adéquate. Benzecri[7] utilise un taux d'inertie corrigé dans le cas de l'analyse des correspondances d'un tableau disjonctif complet. Ce taux est donné sous la forme

$$\tau(\lambda) = \left( \frac{d}{d-1} \right) \left( \lambda - \frac{1}{d} \right),$$

pour  $\lambda \succ \frac{1}{d}$  où  $d$  est le nombre de variables actives.

L'ajustement global du nuage par une droite fait appel à une forme quadratique définie positive et à un ajustement par minimisation d'un critère lié à cette forme.

D'autres approches sont possibles en modifiant le type de distance, la nature du sous espace vectoriel ou les deux en même temps. D'autres critères peuvent être aussi utilisés; à la méthode des moindres carrés, basée sur la norme  $L^2$ , on peut substituer celles des moindres valeurs absolues basée sur la norme  $L^1$  Van Cutsen[99][1994]; ...

## 2.4 L'une des variables est qualitative

Si  $V^j$  est une variable qualitative à  $r_j$  modalités  $\{m_l; l = 1, r_j\}$  décrivant l'individu  $\omega_i$ . La colonne  $j$  du tableau  $X_i$  est de la forme

$$X_i = [x_i^1, \dots, x_i^d],$$

où  $x_i^j$  est l'une des modalités.

Le tableau  $X_i$  est un tableau hétérogène donné sous forme de codage brut. Il est inexploitable sous cette forme.

Devant un tel tableau, nous adoptons l'attitude des factorialistes qui consiste à transformer les variables continues en variables qualitatives dont les modalités sont les classes de l'histogramme correspondant. Les techniques qui peuvent être utilisées pour ce recodage sont nombreuses. Suivant la forme de la dispersion des points sur l'axe des réels, on peut adapter, si le nombre de classes est fixé à priori, la méthode hiérarchique ascendante à une variable et qui consiste à

1. *Choisir un indice d'agrégation*
2. *Dérouler un algorithme du type hiérarchique jusqu'à agréger les  $L$  valeurs en  $k$  classes.*

$$X^j = \begin{bmatrix} x_{i1}^j \\ \vdots \\ x_{iL}^j \end{bmatrix} \mapsto \widetilde{X}^j = \begin{bmatrix} y_{i1}^j \\ \vdots \\ y_{ik}^j \end{bmatrix} \text{ où } y_{il}^j = m_t^j \text{ si } x_{il}^j \in \text{à la classe } t.$$

Ainsi le tableau  $X_i$  peut être transformé en un tableau disjonctif complet  $Z_i$  ou en un tableau de Burt  $B_i$ .

L'ajustement de  $X_i$  se fait à partir de  $Z_i$  ou  $B_i$ .

## 2.5 Ajustement global du nuage

### 2.5.1 Tableau disjonctif complet associé à $X_i$

Soit  $L$  le nombre d'observations répétées sur chaque variable décrivant les individus  $\{\omega_i; i\}$ . Posant

$$P = \sum_{j=1}^d r_j,$$

$r_j$  nombre de modalité ou classes de la variable  $V^j$ .

On construit à partir de  $X_i$  le tableau  $Z_i$  à  $L$  lignes et  $P$  colonnes décrivant les réponses des  $L$  observations pour les  $d$  variables par un codage binaire. Le tableau  $Z_i$  est la juxtaposition de  $d$  sous-tableaux

$$Z_i = [Z_i^1, \dots, Z_i^d],$$





L'effectif total du tableau disjonctif complet est

$$z_i = \sum_{t=1}^L \sum_{k=1}^P (t_{lk}^r) = L \times d.$$

### 2.5.2 Tableau de contingence de Burt associé

On construit à partir de  $Z_i$ , le tableau symétrique  $B_i$  d'ordre  $P$  qui rassemble les croisements deux à deux de toutes les variables qui décrivent l'individu  $\omega_i$ .

$$B_i = (Z_i)' Z_i.$$

Le terme général de  $B_i$  s'écrit

$$b_{jj'}^i = \sum_{l=1}^L z_{jj'}^l.$$

$B_i$  est une juxtaposition de tableaux de contingence.

- Les marges sont pour tout  $j \leq P$  (identiques pour les lignes et les colonnes de par la symétrie de  $B_i$ ) données par

$$b_{j''}^i = \sum_{j=1}^P b_{jj'}^i = d. (Z_i)_{.j'}$$

- L'effectif total  $b_i$  de  $B_i$  vaut

$$\begin{aligned} b_i &= \sum_{j=1}^P \left( \sum_{j'=1}^P b_{jj'}^i \right) = d. \sum_{j'=1}^P (Z_i)_{.j'} = d \left( \sum_{j'=1}^P \sum_{k=1}^L Z_{kj}^i \right) \\ &= d \left( \sum_{k=1}^L \underbrace{\sum_{j=1}^P Z_{kj}^i}_d \right) = d \left( \sum_{k=1}^L d \right) = d^2 . L \end{aligned}$$

$$Z_i = [[Z_i^1] \dots [Z_i^d]] \Rightarrow B_i = (Z_i)' Z_i,$$

le tableau  $B_i$  est formé de  $d^2$  blocs. Le bloc  $(Z_i^q)' \cdot Z_i^{q'}$  indicé par  $q$  et  $q'$  est une matrice rectangulaire de dimension  $r_q \cdot r_{q'}$  qui est le tableau de contingence croisant les variables  $V^q$  et  $V^{q'}$ . Si  $q = q'$ ;  $(Z_i^q)' \cdot Z_i^q$  est un bloc carré de dimension  $r_q$  obtenu par le croisement de la variable  $V^q$  par elle-même.

C'est donc une matrice diagonale de dimension  $r_q$  qu'on note  $W_i^q$

$$W_i^q = (Z_i^q)' \cdot Z_i^q,$$

Les éléments diagonaux sont donnés par

$$(W_i^q)_l = (Z_i)_{.l}.$$

Soit  $D^i$  la matrice diagonale d'ordre  $P$  ayant les mêmes éléments diagonaux que  $B_i$

$$d_{ll}^i = (D_i)_{ll} = b_{ll}^i = (B_i)_{ll} = (Z_i)_{.l}.$$

La matrice  $D^i$  peut être considérée comme une matrice à  $d^2$  blocs où seules les  $d$  matrices diagonales  $\{D_q^i; q = 1, d\}$  ne sont pas nulles. On a

$$D_q^i = (Z_i^q)' \cdot Z_i^q = W_i^q.$$

### 2.5.3 Critère d'ajustement et distance du $\chi^2$ .

Les observations sont toutes affectées d'une masse  $m_l = \frac{1}{L}$ . Chacune des modalités  $m_j$  est affectée du poids  $\frac{1}{dL} (Z_i)_{.j}$ . La distance de  $\chi^2$  appliquée à un tableau disjonctif complet conserve un sens. En effet, dans  $\mathbb{R}^L$ , la distance entre modalités s'écrit

$$d^2(j, j') = \sum_{l=1}^L L \left[ \frac{Z_{lj}^i}{(Z_i)_{.j}} - \frac{Z_{lj'}^i}{(Z_i)_{.j'}} \right]^2,$$

ainsi, deux modalités obtenues par les mêmes observations coïncident. Par ailleurs les modalités de faible effectif sont éloignées des autres modalités. La distance entre deux observations  $\theta_l$  et  $\theta_{l'}$  s'exprime par

$$d^2(l, l') = \frac{1}{d_i} \sum_{l=1}^L \frac{L}{(Z_i)_{.j}} [z_{lj}^i - z_{l'j}^i]^2,$$

ainsi deux observations sont proches si les résultats sont les mêmes modalités, elles sont éloignées sinon.

Soient

- $F_i = \frac{1}{dL} Z_i$  de terme général  $f_{lj}^i = \frac{Z_{lj}^i}{dL}$ ,
- $D_P^i = \frac{1}{dL} D^i$  de terme général  $f_{.j}^i = \delta_{lj} \frac{(Z_i)_{.j}}{d}$
- $D_L^i = \frac{1}{L} \cdot I_L$  ( $I_L$  est la matrice identité d'ordre  $L$ ) de terme général  $f_{.j}^i = \frac{\delta_{lj}}{L}$

Pour trouver les axes factoriels, on diagonalise la matrice

$$S_i = (F_i)' (D_L^i) F_i \cdot D_P^i = \frac{1}{d} (Z_i)' Z_i (D^i)^{-1},$$

Le terme général de  $S_i$  est donné par

$$(S^i)_{j\bar{j}i} = \frac{1}{d(Z_i)_{.j'}} \sum_{l=1}^L z_{lj}^i \cdot z_{l'j}^i,$$

Dans  $\mathbb{R}^P$ , l'équation du  $\alpha^{i\text{ème}}$  axe factoriel  $U_\alpha$  est donnée par

$$\frac{1}{d} (Z_i)' \cdot Z_i \cdot D_i^{-1} \cdot U_\alpha.$$

L'équation du  $\alpha^{i\text{ème}}$  facteur  $\varphi_\alpha = D_i^{-1} \cdot U_\alpha$  vérifie

$$\frac{1}{d} (Z_i)' \cdot Z_i \cdot \varphi_\alpha = \lambda_\alpha.$$

De même que l'équation du  $\alpha^{i\text{ème}}$  facteur  $\Psi_\alpha$  dans  $\mathbb{R}^L$  s'écrit

$$\frac{1}{d} Z_i D_i^{-1} \cdot (Z_i)' \cdot \Psi_\alpha = \lambda_\alpha \Psi_\alpha.$$

Les facteurs  $\varphi_\alpha$  et  $\Psi_\alpha$  représentent les coordonnées des points lignes et les points colonnes sur l'axe factoriel  $\Delta_{U_\alpha}$ . Les relations de transition entre facteurs sont

$$\begin{cases} \Psi_\alpha = \frac{1}{d\sqrt{\lambda_\alpha}} Z_i \cdot \varphi_\alpha \\ \varphi_\alpha = \frac{1}{\sqrt{\lambda_\alpha}} D_i^{-1} (Z_i)' \cdot \Psi_\alpha \end{cases}$$

Les coordonnées de l'observation  $\theta_l$  sur l'axe factoriel  $\Delta_{U_\alpha}$  sont données par

$$(\Psi_\alpha)_l = \frac{1}{\sqrt{\lambda_\alpha}} \sum_{j=1}^P \frac{Z_{lj}^i}{(Z_i)_{.j}} (\varphi_\alpha)_j = \frac{1}{d\sqrt{\lambda_\alpha}} \sum_{j \in P(l)}^P (\varphi_\alpha)_j$$

où  $P(l)$  désigne l'ensemble des modalités obtenues à l'observation  $\theta_l$ .  
La coordonnée de la modalité  $j$  sur l'axe factoriel  $\Delta_{U_\alpha}$  est donnée par

$$(\varphi_\alpha)_j = \frac{1}{\sqrt{\lambda_\alpha}} \sum_{l=1}^L \frac{Z_{lj}^i}{(Z_i)_{.j}} (\Psi_\alpha)_l = \frac{1}{(Z_i)_{.j} \sqrt{\lambda_\alpha}} \sum_{j \in I(j)}^P (\Psi_\alpha)_l$$

où  $I(j)$  désigne l'ensemble des observations ayant été la modalité  $j$ .

#### 2.5.4 Analyse du tableau de Burt: équivalence avec l'analyse du tableau disjonctif complet

L'analyse des correspondances appliquée à un tableau disjonctif complet  $Z_i$  est équivalente à l'analyse du tableau de Burt [12]  $B_i$  correspondant et produit les mêmes facteurs. En effet

L'analyse des correspondances du tableau  $B_i$  (symétrique d'ordre  $P \times P$ ) se ramène à l'analyse d'un nuage de  $P$  points-modalités dans  $\mathbb{R}^P$ . Les marges de ce tableau, en ligne et comme en colonne sont les éléments diagonaux de la matrice  $D_i$ .

Compte tenu de l'équation donnant le  $\alpha^{ième}$  facteur  $\varphi_\alpha$  de l'analyse du tableau disjonctif complet  $Z_i$

$$\frac{1}{d} D_i^{-1} (Z_i)' \cdot Z_i \cdot \varphi_\alpha = \lambda_\alpha \varphi_\alpha (*).$$

La matrice à diagonaliser dans le cas l'AC de  $Z_i$  est

$$S_i = \frac{1}{d} D_i^{-1} (Z_i)' \cdot Z_i = \frac{1}{d} D_i^{-1} B_i.$$

Pour l'analyse du tableau  $B_i$  associé à  $Z_i$ , le tableau des fréquences relatives  $F_i$  s'écrit

$$F_i = \frac{1}{L \cdot d} B_i; D_P^{-1} = D_L^{-1} = \frac{1}{L \cdot d} D_i^{-1}.$$

On diagonalise la matrice

$$S_i^* = \frac{1}{d_i} D_i^{-1} B_i \cdot D_i^{-1} B_i (*) \Rightarrow S_i^* = (S_i)^2.$$

En multipliant les deux membres de(\*) par  $\frac{1}{d_i} D_i^{-1} B_i$  on obtient

$$\frac{1}{d_i^2} D_i^{-1} B_i \cdot D_i^{-1} B_i \varphi_\alpha = \lambda_\alpha^2 \varphi_\alpha.$$

Les facteurs des deux analyses sont donc colinéaires dans  $\mathbb{R}^P$  mais les valeurs propres associées sont différentes. On a la relation

$$\lambda_{B_i} = \lambda_{Z_i}^2,$$

au sens où les facteurs issus de l'analyse de  $Z_i$ , représentent les coordonnées factorielles des modalités ont pour norme  $\lambda_{Z_i}$ . Alors le facteur correspondant de l'analyse de  $B_i$ , noté  $\varphi_{B_i}$  a pour norme

$$\lambda_{B_i} = \lambda_{Z_i}^2.$$

### Relations liant les deux systèmes de coordonnées factorielles

$$\varphi_{B_i \alpha} = \varphi_\alpha \sqrt{\lambda_\alpha}.$$

Pour résumer le tableau  $X_i$  par une droite, on choisira la droite qui donne la plus grande inertie, c'est à dire celle engendrée par le premier axe factoriel dans ce cas

$$\begin{array}{c} Z_i \\ \nearrow \quad \searrow \\ X_i \longrightarrow Y_1 \in \mathbb{R}^P \end{array}$$

dont les coordonnées sont données par

$$(Y_1)_j = (\varphi_1)_j = \frac{1}{(Z_i)_{.j} \sqrt{\lambda_\alpha}} \sum_{j \in I(j)}^P (\Psi_1)_l$$

où  $\Psi_1$  est le vecteur propre associé à la plus grande valeur  $\lambda_1$  non triviale de la matrice

$$\frac{1}{d} (Z_i)' \cdot Z_i \cdot D_i^{-1}$$

Cette étude préliminaire à l'analyse globale d'une juxtaposition de tableaux à observations multiples se résume ainsi

### Cas continu

$$X'_{(M=n \times d, L)} = \begin{bmatrix} X'_1 \\ \vdots \\ X'_n \end{bmatrix} \xrightarrow{\text{après réduction}} Y'_{(n, d)} = \begin{bmatrix} Y''_1 \\ \vdots \\ Y'_n \end{bmatrix}$$

← analyse par des méthodes classiques →

### Cas qualitatif

$$\begin{array}{ccc} & Z'_{(nP, L)} = \begin{bmatrix} Z'_1 \\ \vdots \\ Z'_n \end{bmatrix} & \\ \nearrow & \downarrow & \searrow \\ X'_{(M=vn \times d, vL)} = \begin{bmatrix} X'_1 \\ \vdots \\ X'_n \end{bmatrix} & \downarrow Y'_{(n, d)} = \begin{bmatrix} Y'_1 \\ \vdots \\ Y'_n \end{bmatrix} & \\ & B'_{(nP, L)} = \begin{bmatrix} B'_1 \\ \vdots \\ B'_n \end{bmatrix} \nearrow & \end{array}$$

Cette approche n'est efficace que si l'inertie expliquée par le premier axe est suffisante. Ce qui n'est pas toujours le cas en pratique. Comme elle procède par double ajustements, elle peut servir comme étape préliminaire à l'analyse de grands tableaux.

## 3 Analyse sans réduction

### 3.1 Comparaison globale des tableaux : l'interstructure.

Dans cette partie, nous étudions les liens pouvant exister entre les tableaux décrivant les individus et les comparer globalement. On est conduit pour ce faire à

- Choisir des objets représentatifs des tableaux
- Choisir une métrique entre objets représentatifs
- Trouver une image euclidienne des objets représentatifs.

### 3.1.1 Objet représentatif du tableau $X_k$ décrivant l'individu $\omega_k$

Soit  $W_k$  la matrice des produits scalaires internes au tableau  $X_k$  qui décrit l'individu  $\omega_k$ , on a

$$W_k = \underbrace{(X_k) \cdot (X_k)'}_{(l, r_k) (r_k, L)} \text{ ou } W_k = X_k \cdot Q_k (X_k)'$$

$W_k$  est une matrice carrée d'ordre  $L$ ,  $Q_k$  est une métrique interne aux observations de l'individu  $\omega_k$  pour le groupe de variables qui le décrit  $Q_k$  peut être la métrique des poids affectés à chacune des observations.

L'objet  $W_k$  est représentatif de l'individu  $\omega_k$  car il regroupe les produits scalaires entre les observations de l'individu pour les  $P_k$  variables qui le décrivent.  $W_k$  contient donc tous les liens inter observations.  $W_k$  est une matrice symétrique.

$W_k$  peut être considérée comme un point de  $\mathbb{R}^{L^2}$  en empilant ses  $L$  colonnes. Si  $V^{j^1}, \dots, V^{j^{r_k}}$  sont les  $r_k$  variables qui décrivent l'individu  $\omega_k$ . On a

$$X_k = [x_{t \ js}^k]_{t=1, L; \ s=1, r_k}$$

Soient  $W_k^1, \dots, W_k^L$  les  $L$  colonnes de la matrice  $W_k$ .

**Proposition 1**

La  $t^{\text{ième}}$  coordonnée du vecteur  $W_k^l$  est donnée par

$$(W_k^l)_t = \sum_{m=1}^{r_k} (x_{t \ jm}^k) (x_{t \ jm}^k)$$

Si  $\widehat{W}_k$  est le vecteur de  $\mathbb{R}^{L^2}$  obtenu en empilant les  $L$  colonnes de la matrice  $W_k$

$$\widehat{W}_k = \begin{bmatrix} [W_k^1] \\ \vdots \\ [W_k^L] \end{bmatrix} \text{ où } W_k^l = \begin{bmatrix} W_{1i}^k \\ \vdots \\ W_{Li}^k \end{bmatrix} \in \mathbb{R}^L.$$

Le tableau

$$W_{L^2} = [\widehat{W}_1, \dots, \widehat{W}_n],$$

est le tableau contenant les  $n$  vecteurs de  $\mathbb{R}^{L^2}$  en colonnes,  $W_{L^2}$  est un tableau de dimension  $L^2 \times n$ . On définit ainsi, un nuage de  $n$  points tableaux représentés par des vecteurs de  $\mathbb{R}^{L^2}$  et représentatifs des tableaux caractérisant les  $n$  individus.

### 3.1.2 Choix de la métrique sur $\mathbb{R}^{L^2}$

Afin de représenter graphiquement les  $n$  points tableaux représentatifs des  $n$  objets  $\{\omega_1, \dots, \omega_n\}$ , il est nécessaire de définir une distance entre éléments représentatifs.

#### Produit scalaire de Hilbert-Schmidt Définition 1

Le produit scalaire de Hilbert-Schmidt entre deux objets représentatifs  $\widehat{W}_k$  et  $\widehat{W}_{k'}$  est défini par

$$\langle \widehat{W}_k, \widehat{W}_{k'} \rangle = \text{tr}(D.W_k, DW_{k'}),$$

où  $\widehat{W}_k$  et  $\widehat{W}_{k'} \in \mathbb{R}^{L^2}$ ;  $W_k$  et  $W_{k'}$  sont des matrices symétriques d'ordre  $L$ ,  $D$  : matrice diagonale d'ordre  $L$  contenant les poids affectés aux observations.  $\text{tr}A$  est la trace de la matrice  $A$ .

**Norme et distance induites par le produit scalaire de Hilbert-Schmidt.** La norme d'un élément représentatif est donnée par

$$\left\| \widehat{W}_k \right\|_{HS} = \langle \widehat{W}_k, \widehat{W}_k \rangle_{HS} = \text{tr}(D.W_k, DW_k)$$

#### Proposition 2

Si  $\lambda_k^{(l)}$  est la  $l^{\text{ième}}$  valeur propre de la matrice  $\mathbf{W}_k.D$ . Alors

$$\left\| \widehat{W}_k \right\|_{HS}^2 = \sum_{l=1}^L \lambda_k^{(l)},$$

**Preuve:** Soient  $\lambda_1, \dots, \lambda_p$  les  $p$  valeurs propres de la matrice  $A$  diagonalisable de dimension  $P$ . on a

$$A = Q^{-1}D_\lambda Q,$$

où  $Q = [U_1, \dots, U_p]$  avec  $U_i$  vecteur propre associé à  $\lambda_i$  de  $A$ ,  $D_\lambda$  matrice diagonale dont les éléments sont les valeurs propres  $\{\lambda_i; i = 1, p\}$ . On a

$$\text{tr}(A) = \sum_{l=1}^p \lambda_l,$$



d'une part et d'autre part

$$A^2 = A \times A = (Q^{-1}D_\lambda Q) (Q^{-1}D_\lambda Q) = Q^{-1} (D_\lambda)^2 Q$$

or

$$(D_\lambda)^2 = \begin{pmatrix} \lambda_1^2 & & \\ 0 & \cdot & 0 \\ & & \lambda_p^2 \end{pmatrix},$$

donc  $A_2$  est aussi diagonalisable et

$$tr(A^2) = \sum_{l=1}^p (\lambda_l)^2$$

La distance entre deux éléments représentatifs  $\widehat{W}_1$  et  $\widehat{W}_2$  des individus  $\omega_1$  et  $\omega_2$  est définie à partir du produit scalaire de Hilbert-Schmidt de la manière suivante

$$d_{HS}(\widehat{W}_1, \widehat{W}_2) = \|\widehat{W}_1 - \widehat{W}_2\|_{HS} = \sqrt{\|\widehat{W}_1\|^2 + \|\widehat{W}_2\|^2 - 2 \prec \widehat{W}_1, \widehat{W}_2 \succ_{HS}}$$

### Remarque 1

En pratique, il est fréquent que les éléments  $\widehat{W}_k$  aient des normes très différentes les unes des autres. C'est une situation très délicate quant à l'interprétation des résultats car il est évident que les tableaux de norme assez élevée influent sur les résultats. On est donc conduit à normaliser les éléments représentatifs  $\widehat{W}_k$ .

$$\widehat{W}_k \rightarrow \frac{\widehat{W}_k}{\|\widehat{W}_k\|_{HS}}$$

L'analyse générale du tableau  $W_{L^2}$  en se plaçant sur  $\mathbb{R}^n$  où l'espace vectoriel  $\mathbb{R}^{L^2}$  est muni de la métrique  $M_{HS}$  induite par le produit scalaire de Hilbert-Schmidt et  $\mathbb{R}^n$  est muni de la métrique euclidienne  $I_n$  conduit à la diagonalisation de la matrice carrée  $S$  d'ordre  $n$ . Le terme général de  $S$  est défini par

$$s_{kk'} = tr(W_k \cdot W_{k'}).$$

Si les éléments représentatifs  $\{\widehat{W}_k; k = 1, n\}$  ne sont pas normés et

$$s_{kk'} = \frac{tr(W_k \cdot W_{k'})}{\|\widehat{W}_k\|_{HS} \|\widehat{W}_{k'}\|_{HS}},$$

Si ces éléments sont normés

On se place dans  $(\mathbb{R}^n; I_n)$ , au tableau  $W_{L^2}$  correspond  $n$  points tableaux représentés dans  $\mathbb{R}^{L^2}$  par les colonnes et  $L^2$  points vecteurs de  $\mathbb{R}^n$  qui sont les lignes du tableau. On cherche à ajuster le nuage des  $L^2$  points de  $\mathbb{R}^n$  par des sous espaces vectoriels de dimension  $1, \dots, r$  qui satisfassent au critère des moindres carrés.

### Remarque 2

Si les variables qui décrivent les différents individus sont centrées. Cette hypothèse implique les éléments représentatifs sont centrés et justifie l'ajustement du nuage par des sous espaces vectoriels au lieu de sous espaces affines. En effet

$V^j$  centrée pour tout  $j \in Q_i$  et pour tout  $i = 1, n$ , On a

$$\frac{1}{L} \sum_{h=1}^L (x_{hj}^i) = \bar{x}_j = 0,$$

Or

$$\overline{W_t^i} = \frac{1}{L} \sum_{h=1}^L (W_{ht}^i) = \frac{1}{L} \sum_{h=1}^L \left( \sum_{m=1}^{r_t} (x_{hjm}^i) (x_{tjm}^i) \right) \Rightarrow$$

$$\overline{W_t^i} = \sum_{m=1}^{r_t} x_{tjm}^i \left( \frac{1}{L} \sum_{h=1}^L (x_{hjm}^i) \right) = \sum_{m=1}^{r_t} x_{tjm}^i \cdot 0 = 0$$

### 3.1.3 Recherche du sous espace vectoriel de dimension un.

On dispose d'un nuage de  $L^2$  points de  $(\mathbb{R}^n; I_n)$  qui représentent les lignes du tableau  $W_{L^2}$ .

Soit  $W^i$  la  $i^{\text{ème}}$  ligne du tableau  $W_{L^2}$ ;  $W^i \in (\mathbb{R}^n; I_n)$

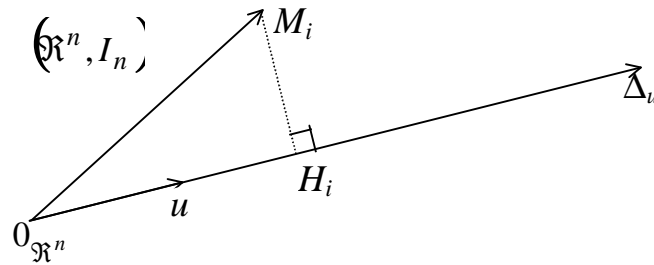


fig 8

$H_i$  projection orthogonale de  $M_i$  sur  $\Delta_u$ ;  $\overrightarrow{OM} = (W^i)'$ . On a

$$\overline{OH_i} = W^i \cdot u$$

En notant par

$$\mathbf{U} = \begin{bmatrix} \overline{OH_1} \\ \vdots \\ \overline{OH_n} \end{bmatrix} = \begin{bmatrix} W^1 \cdot u \\ \vdots \\ W^n \cdot u \end{bmatrix} = W_{L^2} \cdot u \in (\mathbb{R}^{L^2}, M_{HS})$$

Il s'agit de rechercher  $u^*$  tel que l'on ait

$$\underset{u, \|u\|=1}{Max} \|U\|_{HS},$$

avec  $\|U\|_{HS} = U' \cdot M_{HS} \cdot U$

Posant  $M_{HS} = M$  pour faciliter les notations. On doit résoudre le problème d'optimisation suivant :

$$\underset{u, \|u\|=1}{Max} (u' \cdot (W_{L^2})' \cdot M \cdot (W_{L^2}) \cdot u),$$

le Lagrangien s'écrit

$$\varphi(u) = u' \cdot (W_{L^2})' \cdot M \cdot W_{L^2} \cdot u - \lambda (u' \cdot u),$$

condition nécessaire

$$\frac{\partial \varphi}{\partial u} = 0 \Leftrightarrow (W_{L^2})' . M . W_{L^2} . u = \lambda . u,$$

$\lambda$  est donc valeur propre de la matrice  $S = (W_{L^2})' . M . W_{L^2} . u$ . Le maximum est obtenu en choisissant la plus grande valeur propre de  $S$ . L'analyse générale du tableau  $W_{L^2}$  conduit à la diagonalisation de la matrice  $S$ . Le terme général de la matrice à diagonaliser  $S$  est

$$(S)_{kk'} = s_{kk'} = (\widehat{W}_k)' . M . (\widehat{W}_{k'}) = \prec \widehat{W}_k, \widehat{W}_{k'} \succ_{HS} = tr (W_k . W_{k'}) ,$$

en effet

$$W_{L^2} = [\widehat{W}_1, \dots, \widehat{W}_n] \Rightarrow (W_{L^2})' = \begin{bmatrix} (\widehat{W}_1)' \\ \vdots \\ (\widehat{W}_n)' \end{bmatrix} \Rightarrow$$

$$(W_{L^2})' . M . (W_{L^2}) = \left[ (\widehat{W}_k)' . M . \widehat{W}_{k'} \right]_{k=1, L; k'=1, L} \Rightarrow s_{kk'} = tr (W_k . W_{k'})$$

Si l'on norme des éléments représentatifs alors

$$s_{kk'} = \frac{tr (W_k . W_{k'})}{\| \widehat{W}_k \|_{HS} \| \widehat{W}_{k'} \|_{HS}} (*)$$

(\*) est le coefficient *Rv* de [Robert et Escofier, 1976].

Cette analyse du tableau  $W_{L^2}$  permet de représenter les  $n$  points tableaux dans un espace de plus faible dimension et de comparer globalement les tableaux entre eux et par conséquent les individus qu'ils représentent.

Si tous les tableaux sont voisins, ils seront concentrés autour d'un point dans l'espace et le premier axe joindra l'origine à ce point. On pourrait au contraire voir les tableaux s'échelonner le long de cet axe et mesurer ainsi sur l'axe une sorte d'adéquation du tableau au modèle moyen.

**Propriété et interprétation des coefficients Rv** Distance entre deux éléments représentatifs normés  $\widehat{W}_k$  et  $\widehat{W}_{k'}$ .

**Proposition 3**

$$d_{HS} \left( \widehat{W}_k, \widehat{W}_{k'} \right) = \sqrt{2(1 - Rv(k, k'))}$$

où

$$Rv(k, k') = s_{kk'} = \frac{tr(W_k \cdot W_{k'})}{\left\| \widehat{W}_k \right\|_{HS} \left\| \widehat{W}_{k'} \right\|_{HS}}$$

En effet (distances et produits scalaires calculés à partir de la définition d'Hilbert-Schmidt)

$$Rv(k, k') = \prec \frac{\widehat{W}_k}{\left\| \widehat{W}_k \right\|_{HS}}, \frac{\widehat{W}_{k'}}{\left\| \widehat{W}_{k'} \right\|_{HS}} \succ_{HS}$$

et

$$d_{HS}^2 \left( \widehat{W}_k, \widehat{W}_{k'} \right) = \prec \frac{\widehat{W}_k}{\left\| \widehat{W}_k \right\|_{HS}} - \frac{\widehat{W}_{k'}}{\left\| \widehat{W}_{k'} \right\|_{HS}}, \frac{\widehat{W}_k}{\left\| \widehat{W}_k \right\|_{HS}} - \frac{\widehat{W}_{k'}}{\left\| \widehat{W}_{k'} \right\|_{HS}} \succ$$

En utilisant la bilinéarité du produit scalaire, on obtient

$$d_{HS} \left( \widehat{W}_k, \widehat{W}_{k'} \right) = \sqrt{2(1 - Rv(k, k'))}$$

**Propriétés remarquables**

$$Si Rv(k, k') = 1 \Rightarrow d_{HS} \left( \widehat{W}_k, \widehat{W}_{k'} \right) = 0 \Rightarrow \frac{\widehat{W}_k}{\left\| \widehat{W}_k \right\|_{HS}} = \frac{\widehat{W}_{k'}}{\left\| \widehat{W}_{k'} \right\|_{HS}} = c_k$$

Cette égalité se traduit en pratique par le fait que les tableaux  $X_k$  et  $X_{k'}$  sont équivalents dans le sens où les observations de l'objet  $\omega_{k'}$  se déduisent des observations de l'individu  $\omega_k$  par une homothétie de rapport  $c_k$ .

-Si  $Rv(k, k') = 0 \Rightarrow$  les variables qui décrivent l'objet  $\omega_k$  sont non corrélées avec les variables qui décrivent  $\omega_{k'}$

En effet,

pour l'individu  $\omega_{k'}$ , la variable est identifiée par le vecteur  $X_j^{k'}$  de  $\mathbb{R}^L$  où

$$X_{k'} = \left[ \dots \begin{bmatrix} X_{1i}^{k'} \\ \vdots \\ X_{Li}^{k'} \end{bmatrix} \dots \right]$$

Pour l'individu  $\omega_k$ , la variable est identifiée par le vecteur  $X_{.j}^k$  de  $\mathbb{R}^L$  où

$$X_k = \left[ \dots \begin{bmatrix} X_{1i}^k \\ \vdots \\ X_{Li}^k \end{bmatrix} \dots \right]$$

On suppose que  $D = I_L$  (matrice identité d'ordre  $L$ ), on ne privilégie donc aucune observation.

On a

$$W_k = [W_{tj}^k]_{t=1, L; j=1, L}; W_{k'} = [W_{tj}^{k'}]_{t=1, L; j=1, L}$$

$$\Rightarrow tr(W_k \cdot W_{k'}) = \sum_{u=1}^L \left( \sum_{v=1}^L W_{uv}^k W_{uv}^{k'} \right)$$

or

$$W_{uv}^k = \sum_{m=1}^{r_k} (X_{u j_m}^{k'}) (X_{v j_m}^{k'}) \text{ et } W_{uv}^{k'} = \sum_{t=1}^{r_{k'}} (X_{u j_t}^{k''}) (X_{v j_t}^{k'}) \Rightarrow$$

$$tr(W_k \cdot W_{k'}) = \sum_{u=1}^L \left( \sum_{v=1}^L \left( \sum_{m=1}^{r_k} (X_{u j_m}^{k'}) (X_{v j_m}^{k'}) \times \sum_{t=1}^{r_{k'}} (X_{u j_t}^{k''}) (X_{v j_t}^{k'}) \right) \right)$$

$$\Rightarrow tr(W_k \cdot W_{k'}) = \sum_{t=1}^{r_{k'}} \sum_{m=1}^{r_k} \underbrace{\sum_{u=1}^L \sum_{v=1}^L \left( (X_{u j_m}^{k'}) (X_{v j_m}^{k'}) \times (X_{u j_m}^{k''}) (X_{v j_m}^{k'}) \right)}_{*}$$

$$(*) = \sum_{u=1}^L \left( (X_{u j_m}^{k'}) (X_{u j_m}^{k'}) \times \sum_{v=1}^L (X_{v j_m}^{k''}) (X_{v j_m}^{k'}) \right) = \left( \prec X_{..j_m}^k, X_{.j_m}^{k'} \succ \right)^2 \Rightarrow$$

$$tr(W_k \cdot W_{k'}) = \sum_{t=1}^{r_{k'}} \sum_{m=1}^{r_k} \left( \prec X_{..j_m}^k, X_{.j_m}^{k'} \succ \right)^2 \Rightarrow Si tr(W_k \cdot W_{k'}) = 0 \Rightarrow$$

$$\prec X_{j_m}^k, X_{j_m}^{k'} \succ = 0 \Rightarrow$$

les variables décrivant l'individu  $\omega_k$  sont non corrélées avec les variables décrivant l'individu  $\omega_{k'}$ .

Interprétation géométrique des coefficients  $Rv$  On a

$(\mathbb{R}^2, M_{HS})$

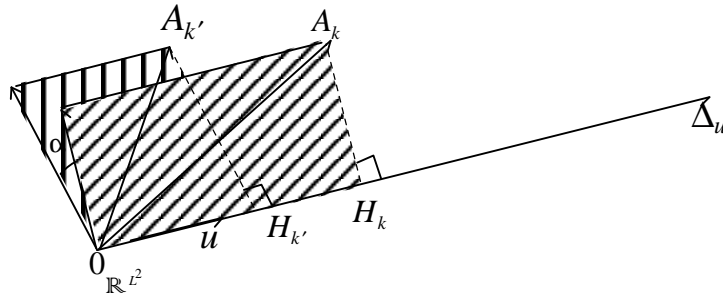


fig 9

$$\overrightarrow{OA_k} = \widehat{W}_k ; \overrightarrow{OA_{k'}} = \widehat{W}_{k'} \text{ et } (\overrightarrow{OA_k}; \overrightarrow{OA_{k'}}) = \alpha$$

$H_k$  et  $H_{k'}$  sont respectivement les projections  $M_{HS}$  orthogonales de  $\overrightarrow{OA_k}$  et  $\overrightarrow{OA_{k'}}$  sur l'axe  $\Delta_u$ . On a

$$\begin{aligned} \angle(\widehat{W}_k, \widehat{W}_{k'}) &= \|\widehat{W}_k\| \|\widehat{W}_{k'}\| \cos \alpha \\ \Rightarrow \cos \alpha &= \frac{1}{\|\widehat{W}_k\| \|\widehat{W}_{k'}\|} \angle(\widehat{W}_k, \widehat{W}_{k'}) = Rv(k, k') \end{aligned}$$

ainsi les coefficients  $Rv$  représentent les cosinus des angles entre les éléments représentatifs de l'image euclidienne.

### 3.2 Le nuage moyen ou compromis : l'intrastructure

L'analyse de l'interstructure a mis en évidence, sans les expliquer, les ressemblances (ou dissemblances) entre les différents tableaux étudiés.

Il est utile et intéressant de construire un nuage moyen qui soit un compromis entre les différents nuages associés aux tableaux  $\{X_k; k = 1, n\}$  par leurs éléments représentatifs  $\{\widehat{W}_k; k = 1, n\}$ . Ce compromis doit être de même nature que les  $\widehat{W}_k; k = 1, n$ .

### 3.2.1 Construction d'un compromis

Le compromis peut être défini de différentes façons en fonction de la nature des données et des connaissances dont on dispose. Le bon sens veut que l'on définisse comme la moyenne pondérée des éléments représentatifs  $\widehat{W}_k$  ou  $\frac{\widehat{W}_k}{\|\widehat{W}_k\|_{HS}}$ . Soit  $c$  le compromis en question.

Alors

$$c = \sum_{t=1}^n \alpha_t \widehat{W}_k \text{ ou } \sum_{t=1}^n \beta_t \frac{\widehat{W}_k}{\|\widehat{W}_k\|_{HS}},$$

Le problème est donc la détermination des coefficients  $\{\alpha_t; t\}$  ou  $\{\beta_t; t\}$ .

Sous les conditions : (c1) et (c2) où

(c1) : Le compromis  $c$  est l'objet représentatif le plus corrélé au sens de Hilbert-Schmidt avec les objets  $W_t$ .

(c2):  $c$  est de même nature que les  $\widehat{W}_k; k = 1, n$ .

(c2) se traduit par

$$\|c\|_{HS} = \sum_{t=1}^n |\alpha_t| \|\widehat{W}_k\|_{HS}$$

Si les éléments représentatifs  $\widehat{W}_k$  sont normés.

$$\|c\|_{HS} = \sum_{t=1}^n |\alpha_t|$$

### 3.2.2 L'expression du compromis

Soit  $u^1$  le vecteur propre de la matrice  $S$  associé à la plus grande valeur propre  $\lambda_1$ . ( $u^1$  est choisi de telle sorte que toutes ses composantes soient positives) ce qui est possible grâce au théorème de Frobenius \*

$$u^1 = \begin{bmatrix} \gamma_1^1 \\ \vdots \\ \gamma_1^n \end{bmatrix}$$

#### Proposition 4

Les coefficients  $\alpha_t$  et  $\beta_t$  sont donnés par



$$\alpha_t = \frac{1}{\sqrt{\lambda_1}} \left[ \sum_{r=1}^L \lambda_r \sqrt{s_{rr}} \right] \lambda_t \gamma_1^t$$

$$\beta_t = \frac{1}{\sqrt{\lambda_1}} \lambda_t \gamma_1^t$$

Une démonstration est proposée par Larvit[56]

L'expression du compromis est

$$\widehat{W}_1 = \sum_{t=1}^L \left( \left[ \frac{1}{\sqrt{\lambda_1}} \left[ \sum_{r=1}^L \lambda_r \sqrt{s_{rr}} \right] \lambda_t \gamma_1^t \right] W_t \right),$$

dans le cas où les éléments représentatifs sont  $\{W_t; t = 1, L\}$

$$\widehat{W}_2 = \sum_{t=1}^L \left( \left[ \frac{1}{\sqrt{\lambda_1}} \lambda_t \gamma_1^t \right] \frac{\widehat{W}_t}{\|\widehat{W}_t\|_{HS}} \right),$$

dans le cas normé.

L'analyse du compromis revient ensuite à effectuer l'analyse en composantes principales ou l'analyse générale de  $\widehat{W}_1$  ou de  $\widehat{W}_2$ . Cela permet de dégager la structure du nuage des observations communes aux tableaux.

# Chapitre III

## Classification des éléments constitutifs d'une structure de tableaux de mesure

### Summary

This chapter deals with the clustering of the elements of a structure of juxtaposition of data measuring tables. One of the main issues in such problems is the selection of a one-dimensional quantity to represent the information included in the repeated observations of each variable. We propose the use of 3 different indexes to measure the distance between elements of a structure and uses the last one based on the Hilbert-Schmidt inner product for clustering purposes through an algorithmic procedure. The algorithm proposed is applied for clustering the customers of an electric company where each customer is described by a curve of load.

**Keywords:** Similarity, constitutive element, structure of multiple table, kullback-Leibler distance, factorial analysis, clustering, criteria, inner product of Hilbert-Schmidt

## 1 Introduction

L'analyse des éléments constitutifs d'une juxtaposition de tableaux multiples reste un large domaine d'investigation. L'objectif principal de la classification automatique et de l'analyse des données est de définir des indices de similarité entre individus ou entre variables qui permettent d'associer à un couple d'objets ou de variables une valeur numérique. Ces indices de mesures sont choisis en fonction des données et de la structure de classification désirée.

Dans le cas d'une structure de juxtaposition de tableaux de mesure ( $TofM$ ), nous définissons 3 indices de similarité. Le premier indice est basé sur la distance de Kullback-Leibler, le second indice sur les techniques de l'analyse factorielle et le troisième sur le produit scalaire de Hilbert-Schmidt. Nous utilisons le troisième indice pour construire un algorithme du type k-means[65].

L'algorithme est utilisé pour classifier les clients domestiques d'une entreprise de distribution d'électricité où chaque client est décrit par des courbes de charges ou d'appels de puissance.

Cette structure de données particulière ( $TofM$ ) est obtenue dès que l'on est conduit à effectuer plusieurs observations pour chacune des variables qui décrit chaque individu.

Des exemples de ce type de structure sont nombreux en pratique, on peut citer

**En médecine:** Dans cas d'une maladie incurable, on observe les patients sur une longue période par des bilans réguliers des paramètres qui donnent une idée sur l'évolution de la maladie.

**En écologie végétale :** Dans le cas où l'on s'intéresse à un ensemble fini d'espèces végétales sur un échantillon de terre. Chaque espèce est caractérisée par un certain nombre

de paramètres écologiques pour chaque relevé.

**En industrie:** Si les objets sont identifiés par des courbes représentatives des pics de consommation sur une période donnée à des instants réguliers.

Dans le paragraphe qui suit (section 2), nous définissons la structure de juxtaposition de tables de données de mesure et introduisons les trois indices de dissimilarité entre éléments de la structure. Dans la section 3 et la section 4 nous construisons un algorithme, écrit en Matlab 5.3, de type k-means utilisant l'indice de distance basé sur le produit scalaire de Hilbert Schmidt. Cet algorithme, section 5 de cette partie, a permis de classifier les clients d'une entreprise de distribution d'électricité où chaque client est décrit par ses courbes de charges construites à partir d'observations régulières.

Une conclusion et certaines remarques constituent la section 6 de ce chapitre.

## 2 Structure de juxtaposition de tableaux de données multiples

Soit  $(\Omega, \Delta)$  une base de connaissance où  $\Omega$  est un ensemble fini d'individus ou d'objets et  $\Delta$  un ensemble de descripteurs de la forme

$$\Delta = [X_1, \dots, X_n],$$

$X_i$  est une matrice d'ordre  $L_i \times d$ , qui contient soit les observations du système des  $d$  variables  $\mathbf{V} = \{V^1, \dots, V^d\}$ . Ces observations sont soit des réalisations du vecteur  $\mathbf{V}$  considéré comme aléatoire soit des observations répétées de l'individu  $\omega_i$  pour chacune des  $d$  variables qui le décrit,  $L_i$  est le nombre de ces observations.

$$X_i = [X_i^j; j = 1, d],$$

$X_i^j$  est un vecteur avec  $L_i$  composantes contenant les observations de l'individu  $\omega_i$  pour la variable  $V^j$ .

$$(X_i^j)' = [x_{i1}^j, \dots, x_{iL_i}^j],$$

$x_{ij}^j$  est la lième observation de l'individu  $\omega_i$  pour la variable  $V^j$ .

Si pour tout  $j = 1, \dots, d$  et pour tout  $i = 1, \dots, n$ ;  $X_i^j \in \mathbb{R}^{L_i}$  alors  $\Delta$  est une structure de juxtaposition de tableaux de mesure (*TofM*). Dans les autres cas,  $\Delta$  peut être considérée comme une juxtaposition de tableaux de données catégorielles (chapitre IV) (*TofC*). C'est une structure de données particulière fréquente en pratique et qui se présente sous la forme

$$\omega_1 \leftrightarrow [X_1]$$

$$\omega_2 \rightarrow [X_2]$$

$$\omega_n \hookrightarrow [X_n]$$

Pour ce type de structure, habituellement nous résumons chaque colonne de chaque tableau par un nombre. La structure est transformée en un tableau classique définissant la description d'un nombre fini d'individus par un nombre fini de variables.

Cependant, cette synthèse nécessite plusieurs hypothèses quant au choix de la valeur qui résume les différentes observations de l'individu pour la variable. Cette valeur peut être obtenue par lissage typologique en utilisant le critère des moindres carrés et la norme  $L^2$  (Anderson[2](1974); Kodell and Chen[52](2001); ...) ou le critère des moindres carrés et la norme  $L^1$  (Arabie and Hubert,[4](1992); Van Cutsen[99](1994); ...).

De plus, ces techniques opèrent de manière ponctuelle, les ajustements obtenus ne permettent pas de reconstituer même partiellement les données de départ car les applications qui associent à chaque colonne une valeur centrale (moyenne, mode, médiane, quantiles, ...) ne sont pas injectives.

Cette réduction se présente comme suit

$$(X)' = [(X_1)', \dots, (X_n)'] \mapsto (Y)' = [(Y_1)', \dots, (Y_n)']$$

$X$  est de dimension  $(M, d)$  avec  $M = \sum_{i=1}^n L_i$ , tandis que  $Y$  est de dimension  $(n, d)$ .

Dans le cas où chaque tableau regroupe les réalisations indépendantes du vecteur aléatoire ou du système des  $d$  variables aléatoires  $\{V^1, \dots, V^d\}$ , on a besoin de rechercher pour tout  $i = 1, \dots, n$  les estimateurs numériques de ce système qui sont les espérances mathématiques  $\{\mu_i^j; j = 1, d\}$  et les composantes des matrices de variances  $\Sigma_i$ .

$$\Sigma_i = \begin{pmatrix} \Sigma_{11}^i & & & \\ & \Sigma_{22}^i & & \\ & & \Sigma_{kk}^i \dots \Sigma_{jk}^i & \\ & & & \Sigma_{nn}^i \end{pmatrix},$$

Où  $\Sigma_{ll}^i = \text{var}(V^l)$  et  $\Sigma_{lk}^i = \text{cov}(V^j, V^k)$  pour  $j \neq k$ .

**Proposition 5**

Nous obtenons les estimateurs suivants

$$\forall j = 1, d; \mu_i^j = \frac{1}{L_i} \sum_{l=1}^{L_i} x_{il}^j; (\sigma_i^j)^2 = \Sigma_{jj}^i = \frac{1}{L_i - 1} \sum_{l=1}^{L_i} (x_{il}^j - \mu_i^j)^2,$$

$$j \neq k; \text{cov}(V^j, V^k) = \Sigma_{lk}^i = \frac{1}{L_i - 1} \sum_{l=1}^{L_i} (x_{il}^j - \mu_i^j) (x_{il}^k - \mu_i^k),$$

ces estimateurs sont sans biais, consistants et convergents et ne dépendent pas du nombre d'observations.

Si les  $d$  variables sont non corrélées même si la taille des échantillons ou le nombre d'observations est relativement petit [10 ou 20], on peut construire des intervalles de confiance correspondants au seuil fixé  $\beta$ .

- Cas où les  $\sigma_i^j$  sont connues.

L'intervalle de confiance au seuil fixé  $\beta$  est donné par

$$I_i^j(\beta) = ]\mu_i^j - \varepsilon_i^j(\beta) ; \mu_i^j + \varepsilon_i^j(\beta)[ ,$$

où

$$\varepsilon_i^j(\beta) = \sqrt{\frac{\sigma_i^j}{L_i - 1}} \arg \Phi^* \left( \frac{1 + \beta}{2} \right),$$

$\Phi^*(y)$  est la valeur telle que la *fdr* de la distribution normale est  $y$ .

- Cas où  $\{\mu_i^j ; j \succeq 1\}$  sont connues et  $\sigma_i^j = 0 ; \forall j \succeq 1$  et les tableaux  $X_i$  sont centrés.

De la même manière, les intervalles de confiance approximatifs de la variance au seuil fixé  $\beta$  sont donnés par

$$I_i^j(\beta) = ]\sigma_i^j - t_\beta.D_i^j ; \sigma_i^j + t_\beta.D_i^j[$$

où

$$D_i^j = \sqrt{\frac{0.8.L_i + 1/2}{L_i(L_i - 1)}} \sigma_i^j ; t_\beta = \arg \Phi^* \left( \frac{1 + \beta}{2} \right)$$

Le tableau  $X_i$  peut être transformé en un pavé de  $\mathbb{R}^d$ .

$$X_i \rightarrow R_i = \prod_{J=1}^d I_i^j(\beta)$$

où  $I_i^j(\beta)$  est un intervalle de  $\mathbb{R}$ .

C'est le cas de données avec erreurs de mesure,  $R_i$  est un pavé de dimension  $d$  centré à l'origine de  $\mathbb{R}^d$  [chapitre 5]

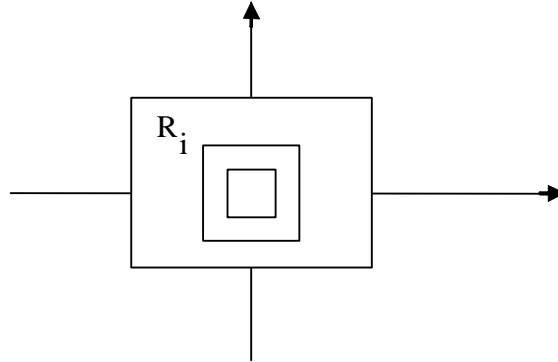


fig 10

P.Cazes, A. Chouakria, E.Diday, Y Scheketman (1997), [19] ont adapté une analyse en composantes principales pour étudier les éléments constitutifs d'une structure de ce type.

### 3 Indices de Distances

#### 3.1 Indice basé sur la distance de Kullback-Leibler

##### 3.1.1 Cas discret

Soient  $F_1, F_2$  2 densités multidimensionnelles empiriques dont la loi de probabilité est donnée par le tableau

$F_{(\cdot)} \rightarrow$

$$\left[ \begin{array}{cccc} (m_1^1, m_1^2, \dots, m_1^d) & (m_1^1, m_2^2, m_1^3, \dots, m_1^d) \dots & (m_{i_1}^1, m_{i_2}^2, \dots, m_{i_d}^d) \dots & (m_{r_1}^1, m_{r_2}^2, \dots, m_{r_d}^d) \\ P_{11..1}^{(\cdot)} & P_{12..1}^{(\cdot)} & P_{i_1 i_2 \dots i_d}^{(\cdot)} & P_{r_1 r_2 \dots r_d}^{(\cdot)} \end{array} \right]$$

avec  $\sum_{i_1=1}^{r_1} \sum_{i_2=1}^{r_2} \dots \sum_{i_d=1}^{r_d} (P_{i_1 \dots i_d}^{(\cdot)}) = 1$ ;  $m_l^j$  est la  $l^{\text{ième}}$  modalité ou classe de la variable  $V^j$ ,

$r_j$  est le nombre de modalités ou classes.

Ce tableau peut être obtenu à partir d'une structure de juxtaposition de données catégorielles (*TofcD*) (voir Burt, [12]).

##### Définition 2

La distance de Kullback-Leibler entre les 2 distributions est donnée par

$$\begin{aligned}
d_1(F_1, F_2) &= \sum_{i_1=1}^{r_1} \dots \sum_{i_d=1}^{r_d} \left[ \left( P_{i_1 \dots i_d}^{(1)} - P_{i_1 \dots i_d}^{(2)} \right) \log_2 \left( \frac{P_{i_1 \dots i_d}^{(1)}}{P_{i_1 \dots i_d}^{(2)}} \right) \right] \\
&= K_u(F_1, F_2) + K_u(F_2, F_1) - [K_u(F_1, F_2) + K_u(F_2, F_1)]
\end{aligned}$$

avec

$$K_u(F_1, F_2) = \sum_{i_1=1}^{r_1} \dots \sum_{i_d=1}^{r_d} \left[ \left( P_{i_1 i_2 \dots i_d}^{(1)} \right) \log_2 \left( P_{i_1 i_2 \dots i_d}^{(2)} \right) \right]$$

### 3.1.2 Cas continu

Soit  $f_1, f_2$  2 densités de probabilités.

#### Définition 3

La distance de Kullback-Leibler entre les 2 densités est donnée par

$$\begin{aligned}
d_1(f_1, f_2) &= \int_{\mathbb{R}^d} (f_1 - f_2) \log_2 \left( \frac{f_1}{f_2} \right) d\mu \\
&= K_u(f_1, f_2) + K_u(f_2, f_1) - [K_u(f_1, f_2) + K_u(f_2, f_1)],
\end{aligned}$$

avec

$$K_u(f_1, f_2) = \int_{\mathbb{R}^d} (f_1) \log_2(f_2) d\mu.$$

Après Shannon[90], Kullback[55] explique que la quantité évaluée l'information moyenne perdue si on utilise la distribution  $f_1$  alors que la vraie distribution est  $f_2$ .

Si  $f_1$  et  $f_2$  sont 2 densités de probabilité de la loi normale multidimensionnelle de dimension  $d$  de paramètres  $\{(\mu_i, \Sigma_i); i = 1, 2\}$ . Pour  $i = 1, 2$

$$\begin{aligned}
f_i &\longrightarrow N(\mu_i; \Sigma_i)' \\
\Rightarrow f_i(x) &= (2\Pi)^{-\frac{d}{2}} |\Sigma_i|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} \|x - \mu_i\|_{\Sigma_i^{-1}}^2 \right],
\end{aligned}$$

#### Proposition 6

On a

$$d_1(f_1, f_2) = \frac{1}{2 \log(2)} \left[ |\Sigma_1 \Sigma_2^{-1}| + |\Sigma_1^{-1} \Sigma_2| - \|\mu_1 - \mu_2\|_{\Sigma_1^{-1}}^2 - \|\mu_1 - \mu_2\|_{\Sigma_2^{-1}}^2 - 2 \right].$$

### Hypothèse

$$\text{Si } d_1(f_1, f_2) = 0 \Rightarrow f_1 = f_2 \Rightarrow \omega_1 = \omega_2.$$

Dans le cas où  $\mu_1 = \mu_2 \Rightarrow$

$$d_1(f_1, f_2) = \frac{1}{2 \log(2)} \left[ |\Sigma_1 \Sigma_2^{-1}| + |\Sigma_1^{-1} \Sigma_2| - 2 \right],$$

cet indice de distance peut être utilisé en associant un indice d'agrégation adéquat pour construire une hiérarchie indicée sur les éléments constitutifs de la structure de tableaux de données de mesure ou de données catégorielles.

Cependant, le choix de la distribution qui représentera une classe d'individus décrits par des distributions reste un problème posé, en conséquence un critère du type intra-classe basé sur cet indice de distance ne peut pas être défini.

## 3.2 Indice de distance basé sur des techniques factorielles

Cet indice est largement étudié dans Naab, (2001)[68] dans le cadre d'une thèse de Magister effectué sous ma direction. Soit

$$(X)' = [(X_1)', \dots, (X_n)']$$

une structure de juxtaposition de tableaux multiples.  $X_i$  est le tableau qui contient les observations de l'individu  $\omega_i \in \Omega$ . On peut utiliser les techniques factorielles pour réduire chacun des tableaux  $X_i$ . Chaque analyse du tableau  $X_i$  nous donne un système d'axes

$$\{E_i = \{\Delta_{u_r}; r \succeq 1\}; i = 1, n\}$$

Le problème consistera à rechercher un système commun. Bouroche[10] a proposé un certain nombre de critères qui permettent de choisir ce système d'axes commun à toutes les analyses, mais les calculs sont compliqués et parfois non justifiés. Nous retenons 3 critères pour la recherche du système d'axes commun à tous les systèmes obtenus par les différentes analyses factorielles. Ces nouveaux systèmes permettent de comparer les différents éléments de la structure mais beaucoup d'informations sera perdue par les procédures proposées. Ces approches ne sont donc pas très recommandées.

Soit  $X$  un tableau de données, on suppose que les objets fictifs sont classifiés en  $n$  classes



$$\{C_1, \dots, C_n\}$$

d'effectifs respectifs

$$L_1, \dots, L_n.$$

Chaque classe représente un individu.

L'analyse en composantes principales du tableau  $X$  conduit à la diagonalisation de la matrice

$$W = \sum_{i=1}^n (X_i)' D_i (X_i),$$

$D_i$  est une matrice diagonale contenant respectivement les poids associés aux observations de chacun des objets. S'il n'y a pas d'observations privilégiées alors

$$D_i = \frac{1}{L_i} \mathfrak{S}_{L_i},$$

$\mathfrak{S}_{L_i}$  est la matrice identité de dimension  $L_i$

Soit  $u_r$  le  $r$ ème axe factoriel. On a

$$\|u_r\|_{\mathfrak{S}_d} = 1 \text{ et } u_r \perp_{\mathfrak{S}_d} u_{r'} \text{ si } r \neq r'$$

Le représentant de l'individu  $\omega_i$  ( $C_i$ ) sur cet axe factoriel est donné par

$$(z_i)_r = \widetilde{X}_i \cdot u_r \in \mathbb{R}^{L_i},$$

$\widetilde{X}_i$  tableau centré

Le représentant de la classe  $C_i$  est

$$T_i = [(z_i)_1, \dots, (z_i)_{r \max}],$$

$T_i$  est une matrice de dimension  $L_i \times r \max$ ,  $r \max$  est le nombre maximum retenu pour les valeurs propres de la matrice  $W$ .

**Proposition 7**

La quantité

$$d_2(C_1, C_2) = \left( \sum_{r=1}^{r \max} \left( \|(z_1)_r - (z_2)_r\|_{\mathfrak{S}_{r \max}} \right)^2 \right)^{\frac{1}{2}},$$

est un indice de distance entre éléments représentatifs de la structure.

### Hypothèse

$$d_2(C_1, C_2) = 0 \Rightarrow C_1 = C_2 \Rightarrow \omega_1 = \omega_2.$$

La quantité d'information perdue par cette procédure est évaluée à

$$\epsilon = 100 \left( 1 - 100 \left( \frac{\sum_{\alpha=1}^{r \max} (\lambda_\alpha)}{d} \right) \right).$$

### 3.3 Indice basé sur le produit salaire de Hilbert-Schmidt

Soit  $W_i$  la matrice symétrique des produits scalaires du tableau  $X_i$  le  $i^{\text{ème}}$  élément de la structure, le teme général de cette matrice est

$$(W_i)_{lj} = \sum_{m=1}^{L_i} p^j x_{lm}^i (x_{mj}^i),$$

$p^j$  est le poids associé à la variable  $V^j$  avec  $\sum_{j=1}^d p^j = 1$

La matrice symétrique  $W_i$  d'ordre  $L_i$  représente l'individu  $\omega_i$  car elle regroupe les produits scalaires entre les observations des variables.

Si  $L_i = L$  et  $p^j = \frac{1}{d}$  # Cette hypothèse se traduit en pratique par le fait que les objets ou les individus sont observés le même nombre de fois et qu'il n'y a pas d'observations privilégiée #

$W_i$  peut être considérée comme un point de l'espace  $\mathbb{R}^{L^2}$  en empilant ses L colonnes

$$W_i \mapsto (\widetilde{W}_i)' = [(W_i^1), \dots, (W_i^L)],$$

$$\text{avec } \mapsto (W_i^l)' = (W_{i1}^l, \dots, W_{iL}^l)$$

$$\Rightarrow W_i^l \in \mathbb{R}^L \quad \forall l = 1, \dots, L \text{ et } \widetilde{W}_i \in \mathbb{R}^{L^2}$$

On défini ainsi un nuage de  $n$  vecteurs de  $\mathbb{R}^{L^2}$   $\{\widetilde{W}_i; i = 1, \dots, n\}$  représentant les  $n$  tables  $\{X_i; i = 1, \dots, n\}$ .

### 3.3.1 Le produit scalaire de Hilbert-Schmidt

#### Proposition 8

L'application

$$\delta : M_L \times M_L \rightarrow \mathbb{R}_+$$

$$(W_1, W_2) \rightarrow \delta(W_1, W_2) = \prec W_1, W_2 \succ_{HS} = tr((DW_1 \cdot DW_2))$$

est un produit scalaire,  $M_L$  est l'ensemble des matrices symétriques d'ordre  $L$ , la matrice diagonale  $D$  d'ordre  $L$  contient les poids affectés aux observations.

*Ce produit scalaire est le produit scalaire de Hilbert-Schmidt.*

Si  $D$  est la matrice identité d'ordre  $L$  alors

$$\prec W_1, W_2 \succ_{HS} = tr(W_1 \cdot W_2) = \prec \widetilde{W}_1, \widetilde{W}_2 \succ_{id},$$

$\widetilde{W}_1, \widetilde{W}_2 \in \mathbb{R}^{L^2}$ ,  $\prec$ ,  $\succ_{id}$  produit scalaire classique.

### 3.3.2 Norme et distance induite par le produit scalaire de Hilbert-Schmidt

La norme d'un élément représentatif  $W_1$  est donnée par

$$\prec W_1, W_1 \succ_{HS} = tr(W_1 \cdot W_1) = \prec \widetilde{W}_1, \widetilde{W}_1 \succ_{id} = \left\| \widetilde{W}_1 \right\|^2,$$

La distance entre deux objets  $\omega_1$  et  $\omega_2$  est donnée par

$$d_3(\omega_1, \omega_2) = \|W_1 - W_2\|_{HS} = \left\| \widetilde{W}_1 - \widetilde{W}_2 \right\|_{id}.$$

$$d_3(\omega_1, \omega_2) = \sqrt{\left\| \widetilde{W}_1 \right\|^2 + \left\| \widetilde{W}_2 \right\|^2 - 2 \left\| \widetilde{W}_1 \right\| \left\| \widetilde{W}_2 \right\| \cos(\widetilde{W}_1, \widetilde{W}_2)}$$

En pratique, les éléments représentatifs sont souvent normés pour éviter que les éléments qui ont une grande norme n'influent sur l'analyse. Ainsi,

Si  $\left\| \widetilde{W}_1 \right\| = \left\| \widetilde{W}_2 \right\| = 1$ , alors

$$d_3(\omega_1, \omega_2) = \sqrt{2 - R_v(1, 2)}$$

où

$$R_v(1, 2) = \cos(\widetilde{W}_1, \widetilde{W}_2) = \frac{\langle \widetilde{W}_1, \widetilde{W}_2 \rangle}{\|\widetilde{W}_1\| \|\widetilde{W}_2\|} \stackrel{id}{=} \frac{tr(W_1 \cdot W_2)}{\|W_1\|_{HS} \|W_2\|_{HS}}$$

$R_v$  est donc le coefficient d'association entre les individus décrits par leur élément représentatif respectif. C'est un coefficient de corrélation vectoriel entre éléments représentatifs.

## 4 Critère et problème d'optimisation

Nous désirons regrouper les  $n$  individus en  $k$  classes où chaque individu est décrit par son élément représentatif. L'homogénéité des classes est mesurée par

$$H_k = \sum_{\{i\} \in P_k} d_{HS}(W_i, l_k),$$

où  $l_k$  est l'élément représentatif du noyau de la classe  $P_k$  et  $d_{HS}$  est la distance induite par le produit scalaire de *Hilbert Schmidt*.  $H_k$  mesure la dispersion des individus à l'intérieur de la classe.  $H_k$  est aussi l'inertie totale des objets de la classe autour de son noyau.

Soit  $P_K$  l'ensemble des partitions à  $k$  classes

$$P \in P_K \Rightarrow P = \{P_1, \dots, P_K\}$$

est une partition à  $K$  classes,  $L_K$  l'ensemble des  $K$  représentations

$$(L \in L_K) \Rightarrow L = \{l_1, \dots, l_k\}; l_t \in M_L.$$

On définit l'application

$$C_r : P_K \times L_K \rightarrow \mathbb{R}_+$$

$$(P, L) \rightarrow C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} d_{HS}(W_i, l_k)$$

où

$$d_{HS}(W_i, l_k) = \left\| \widetilde{W}_i - \widetilde{L}_k \right\| = tr(W_i \times l_k)$$

$$\widetilde{W}_i; \widetilde{L}_k \in \mathbb{R}^{L^2}, W_i; l_k \in M_l$$

Le critère s'écrit

$$C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} \left( \sqrt{\|\widetilde{W}_i\|^2 + \|\widetilde{L}_k\|^2 - 2 \langle \widetilde{W}_i, \widetilde{L}_k \rangle} \right)^{1/2}$$

Ce critère exprime l'adéquation entre les éléments représentatifs des individus dans les classes avec les éléments représentatifs des noyaux des classes.  $C_r$  s'écrit

$$C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} \left( 2 \left( 1 - R_v(\widetilde{W}_i, \widetilde{L}_k) \right) \right)^{\frac{1}{2}}$$

Le coefficient  $R_v$  avec ses propriétés permet d'expliquer l'adéquation entre la partition et les représentants de ses classes. On note que plus  $R_v$  est proche de 1 plus petit est le critère et plus les individus sont en adéquation avec les noyaux des classes où ils sont affectés.

On recherche  $(P^*, L^*)$  qui réalise

$$\underset{\substack{P \in P_K \\ L \in L_K}}{\text{Min}} C_r(P, L)$$

Les algorithmes utilisés pour résoudre ce type de problème sont du type k-means. L'algorithme est basé sur la définition de la fonction de représentation  $g$  et de la fonction d'affectation  $f$ .

## 4.1 La fonction de représentation

$$g : P_k \rightarrow L_k$$

$$P = \{P_1, \dots, P_k\} \rightarrow g(P) = \{L_1, \dots, L_k\}$$

$g$  doit vérifier

$$\underset{L \in L_k}{\text{Min}} C_r(P, L) = C_r(P, g(P))$$

**Proposition 9**

Ce minimum est obtenu pour

$$g(P) = \widehat{L} = \{\widehat{L}_1, \dots, \widehat{L}_k\} = \widehat{G}; \widehat{L}_l = \widehat{G}_l \in M_L$$

est le centre de gravité de la classe  $P_l$ . Ce centre de gravité ou noyau de la classe est donné par

$$\widehat{G}_l = \left[ (\widehat{G}_1)', \dots, (\widehat{G}_L)' \right]$$

où pour  $l = 1, k$

$$(\widehat{G}_l)' = [G_{i1}^l, \dots, G_{iL}^l] \Rightarrow \widehat{G}_l \in \mathbb{R}^L,$$

avec

$$L_{it}^l = G_{it}^l = \frac{1}{N_l} \sum_{(i) \in P_l} (W_i)_{it}; N_l = |P_l|$$

**4.2 La fonction d'affectation**

$$f : L_k \rightarrow P_k$$

$$L = \{L_1, \dots, L_k\} \rightarrow f(L) = \{P_1, \dots, P_k\}$$

$f$  doit vérifier

$$\underset{P \in P_k}{MinCr}(P, L) = Cr(f(L), L)$$

**Proposition 10**

Ce minimum est obtenu pour

$$P_l = \{\omega_i \in \Omega / d_{HS}(\omega_i, L_l) \preceq d_{HS}(\omega_i, L_t); \forall t \neq l \text{ et } l \prec t \text{ en cas d'égalité}\}$$

## 5 L'algorithme

On choisit au hasard ou par d'autres artifices  $L^{(0)}$  et on utilise alternativement les fonctions  $f$  et  $g$  pour faire décroître le critère. L'algorithme se déroule comme suit

$$L^{(0)} \xrightarrow{f} P^{(1)} \xrightarrow{g} L^{(1)} \xrightarrow{f} P^{(2)} \dots L^{(*)} \xrightarrow{f} P^{(*)}$$

L'algorithme s'arrête dès que la partition obtenue ne change plus. Nous construisons ainsi 2 suites  $U_n$  et  $V_n$  telles que

$$\begin{cases} U_n = Cr(P^{(n)}, L^{(n)}) \\ V_n = (P^{(n)}, L^{(n)}) \end{cases}$$

### Proposition 11[31]

La suite  $U_n$  converge en décroissant et la suite  $V_n$  est stationnaire à partir d'un certain rang donc convergente

## 6 Structure générale de l'algorithme : Organigramme

### Déclaration des données

$T_i = [t_{hl}^i]$  tableau de dimension  $H \times J$

$i = 1, \dots, n$  ;  $\varepsilon =$  seuil d'arrêt

$n =$  nombre d'individus,  $J =$  nombre de jours

$H =$  nombre d'heures,  $K =$  nombre de classes

↓

### Détermination des éléments représentatifs

$i = 1, \dots, n$ ;  $W_i = T_i \times (T_i)'$

↓

Empilage des colonnes des éléments représentatifs

$$\widetilde{W}_i = \frac{\widetilde{W}_i}{\|\widetilde{W}_i\|_{HS}}$$

↓

*Construction de la partition initiale*

↓

### Partition $P^{(0)}$

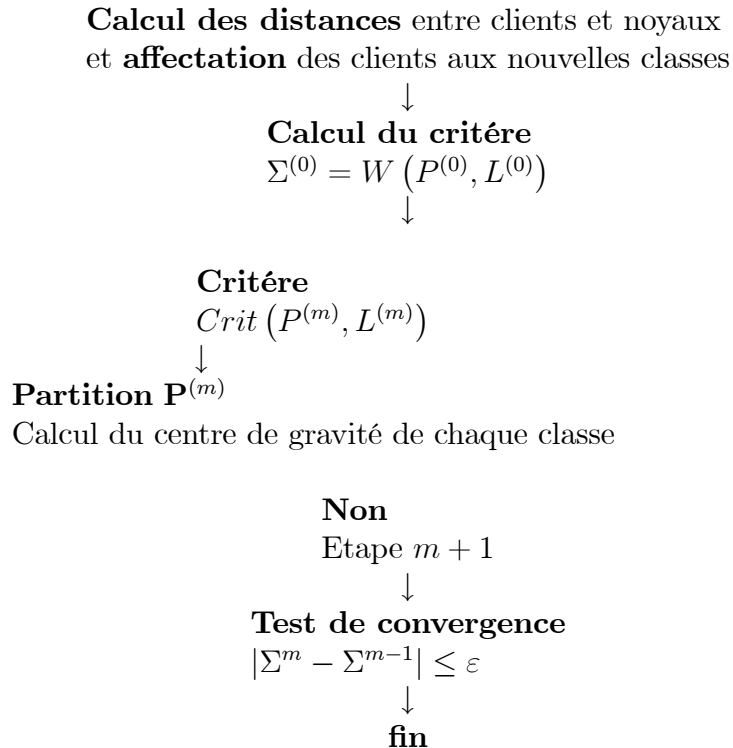
Construction de  $P^{(0)}$  à partir d'un tirage pseudo-aléatoire de  $K$  nombres choisis parmi les  $n$  numéros

↓

### Représentation $L^{(0)}$

Calcul des  $K$  noyaux

↓



## 7 Application

### 7.1 Données en entrée

Ces données ont fait l'objet d'une étude dans le cadre de plusieurs projets de fin d'études. Nous les utilisons pour montrer que l'analyse faite après une étape de réduction conduit à des résultats non interprétables par les paramètres qui caractérisent les clients et que si la variabilité de la consommation d'électricité est assez importante il faut éviter l'étape réduction.

L'algorithme est implémenté en MATLAB 5.3 software et est appliqué pour classifier 47 individus en 5 classes. Les 47 individus constituent un panel de clients domestiques basse tension choisis pour représenter les abonnés de la zone *II*. Ce choix a été fait par des techniques de sondage éprouvées. La stratification des données s'est faite sur la base du critère de consommation de l'électricité de manière à ce qu'il y ait égalité de la consommation annuelle par strates (10 strates). La répartition des individus par strates s'est faite en minimisant la variance de la consommation des échantillons. Ainsi, le tirage des échantillons s'est fait par zones (9 zones). Un pas de tirage de chaque strate et de chaque zone a été défini et on a tiré dans chaque zone de chaque strate un abonné de la liste principale (900000 abonnés).

Pour ces 47 clients, la société de distribution a installé des compteurs spéciaux qui mesurent l'appel de puissance toutes les 10 minutes pendant toute la durée de l'étude (les 6 premiers mois de l'année 2000). Devant la grande masse de données relevées, l'étude s'est restreinte aux relevés du mois d'avril qui marque la fin de l'hiver et pas encore l'été (le chauffage et la climatisation ne sont pas nécessaires). Les 3 jours éliminés du mois



sont les jours de fin de semaine car certains clients sont absents du logement tandis que pour d'autres, il y a une consommation maximale durant les 3 jours, due à leur présence continue.

L'entreprise a choisi d'observer les pics de consommations de chaque client durant les 28 jours ouvrables du mois d'avril.

## 7.2 Classification après réduction

Pour chaque client, on a construit sa courbe de charge ou d'appels de puissance regroupant en abscisses les 24 heures de la journée et en ordonnée ses pics de demande de consommation ou de charge. Pour cette approche, on a résumé la courbe des appels de puissance pour chaque client du panel par un vecteur regroupant les moyennes des appels de puissance pour les 28 jours pour chaque heure. Ainsi,

$$\widehat{T}_i \xrightarrow{(28,720)} C_i \rightarrow T_i \xrightarrow{((28 \times 24))} Y_i = [\overline{x}_i^1, \overline{x}_i^2, \dots, \overline{x}_i^d]; \quad d = 28$$

$$Y \xrightarrow{(47,28)} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_{47} \end{bmatrix}$$

**Table 4:** Meilleure partition, en 5 classes, obtenue après 10 itérations en déroulant l'algorithme k-means sur le tableau Y

classes	1	2	3	4	5
Individus dans les	1, 17, 18, 22	26	16	33	2, 3, 4, 5, 6, 7, 8, 0, 10, 11, 12, 13, 14
classes	28, 29, 35, 40	42	27	41	15, 19, 20, 21, 23, 24, 25, 26, 30, 31
	44, 46	43			32, 34, 36, 37.38.39.41.45.47
<b>effectif</b>	<b>10</b>	<b>3</b>	<b>2</b>	<b>2</b>	<b>30</b>

Pour chaque client, on a construit sa courbe de charge ou d'appels de puissance regroupant en abscisses les 24 heures de la journée et en ordonnées ses pics de demande de consommation ou de charge. Pour le mois d'étude, les observations de  $\omega_i$  sont donc regroupées dans des courbes de charges de ce type

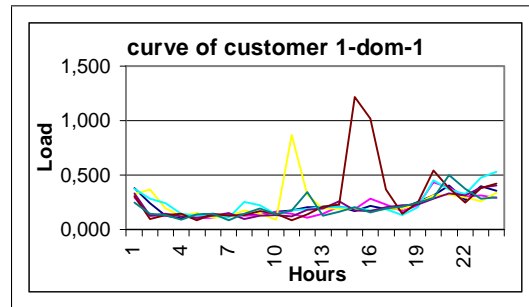


fig 11

La grande variabilité de la consommation durant les minutes et les heures de la journée (en période de pointe et en période creuse) fait que l'interprétation des classes à partir des caractéristiques de chaque client (type d'habitation, qualité des produits électriques utilisés, nombre de lampes, nombre de pièces, nombre de personnes constituant le ménage, ...) n'a pas donné de bons résultats.

Pour cet exercice, la condition de normalité n'est pas vérifiée. De plus, les courbes dépendent les unes des autres. La consommation moyenne pour chaque heure ne peut donc résumer les différentes observations.

### 7.3 Classification basé sur l'indice de Hilbert-Schmidt

$C_i$  est une courbe construite à partir d'un nombre fini de points. La structure de données en entrée peut être considérée comme une juxtaposition de tableaux de mesure.

$$\hat{T}_i \rightarrow C_i \rightarrow T_i = \begin{matrix} & & & & 1 & 2 & \dots & j & \dots & 24 \\ & & & & 1 & & & \dots & & \\ & & & & \vdots & & & & & \\ & & & & l & \left( \begin{matrix} \vdots & & & \dots & & \vdots \\ & t_{il}^j & & & & \\ \vdots & & & & & \\ & & & & & \vdots \\ 28 & & & & \dots & \end{matrix} \right) & & & & \end{matrix}$$

$t_{il}^j$  est la consommation de l'individu durant l'heure du jour  $j$ .  $T_i$  est un tableau de dimension  $l \times 24$ . Les clients sont identifiés par les nombres 1 à 47.

**Table 5:** partition initiale construite autour de 5 individus choisis au hasard parmi les 47 par la fonction random.

classes	1	2	3	4	5
Individus dans les	5, 6, 8, 9	1, 10	34	2, 3, 4, 7, 12, 13	11, 16, 17
classes	14, 15, 18	21, 33	39	20, 22, 23, 26, 27, 28	24, 25, 31
	30, 44, 47	37, 41	43	29, 32, 35, 36, 40, 42, 47	38, 45
<b>effectifs</b>	<b>10</b>	<b>6</b>	<b>3</b>	<b>20</b>	<b>8</b>

Valeur du critère pour cette partition

$$Cr = 9.9455$$

**Table 6:** Partition optimale obtenue après 8 itérations

classes	1	2	3	4	5
Individus dans les	5, 6, 8, 9	1, 4, 10, 12, 13, 19, 21	20, 23	2, 3, 7	16, 17
classes	14, 15, 18	23, 25, 26, 27, 30, 33	22, 32	11, 28	24, 31
	44, 47	35, 36, 37, 40, 41, 42, 46	34	29, 43	38, 45
<b>effectifs</b>	<b>9</b>	<b>20</b>	<b>5</b>	<b>7</b>	<b>6</b>

Valeur du critère pour cette partition

$$Cr = 4.0498$$

## 7.4 Interprétation de ces résultats à l'aide des courbes représentatives des centres ou noyaux des classes

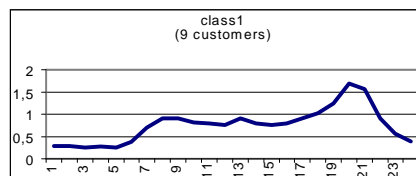


fig 12

La courbe représentant la **classe 1** regroupe 9 clients et montre un pic d'appels de puissance en période de pointe, un pic dominant en matinée et une demande uniforme en période creuse.

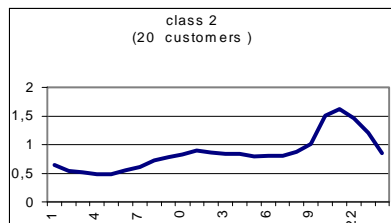


fig 13

La **classe 2** regroupe 20 clients et présente un pic d'appels de puissance identique à celui de la classe 1 mais les 2 courbes n'ont pas la même forme.

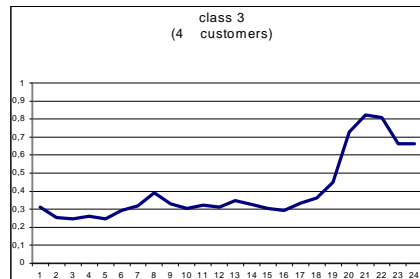


fig 14

La courbe de la **classe 3** regroupe 4 clients et montre un pic d'appel de puissance en heures de pointe équivalent à 0,8 kwh et des appels de puissance uniformes en périodes creuses.

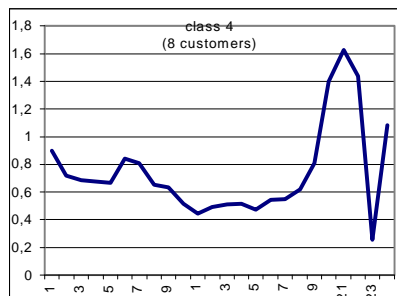


fig 15

On note que la **classe 4** regroupe 8 clients et est caractérisée par un pic d'appels de puissance dominant en soirée de 1.6 kwh et une décroissance dans les autres périodes de la journée.

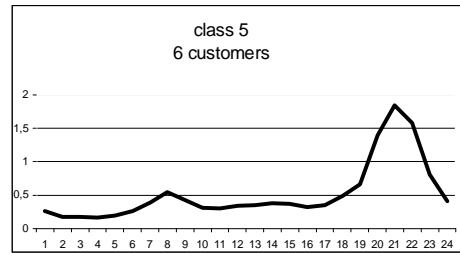


fig 16

La **Classe 5** regroupe 6 clients et est caractérisée par un pic dominant en soirée de 1.9 kw/h et une uniformité des appels de puissance dans les autres périodes de la journée avec un petit pic en matinée.

Cette application a montré que les courbes représentatives des noyaux finaux obtenues par l'algorithme présentent des formes semblables avec des niveaux différents. De plus, l'algorithme a donné des classes bien séparées. Pour ce type de données où le ratio entre les appels de puissance durant les heures de pointes et les heures creuses est supérieur à 3 , il n'est pas recommandé de réduire les données avant la classification . Enfin, nous avons pu expliquer les classes en fonction des autres caractéristiques des clients autres que la consommation d'électricité (type d'habitation, nombre de pièces, nombre de lampes utilisées, nombre de personnes ,...) ce que nous n'avons pas pu faire pour les classes obtenues après réduction.

## 8 Conclusion

La distance induite par le produit scalaire d'Hilbert Schmidt a permis de construire un algorithme de type k-means pour classifier les éléments constitutifs d'une structure de juxtaposition de tableaux de mesure de même dimension. Cependant, cet indice de distance ne s'étend pas au cas d'une structure de données catégorielles et il serait intéressant de trouver l'équivalent comme dans le cas de la distance euclidienne et la distance du khi-deux car on ne doit pas oublier que cette distance conduit à la distance euclidienne classique ,je pense que c'est une piste intéressante.

# Chapitre IV

## Classification des éléments constitutifs d'une structure de tableaux de données catégorielles

### 1 Introduction

Dans ce chapitre, on s'intéresse à la classification d'objets matriciels qui constituent les composants d'une structure des tableaux multiples et de manière particulière, à des matrices qui sont des tableaux du type disjonctifs complets donc à des matrices 0/1. Les indices de distance basées sur des trajectoires factorielles et sur le produit scalaire de Hilbert-Schmidt proposés dans le chapitre précédent ne s'adaptent ni à ce type de structure (données qualitatives) ni au cas où les objets matriciels sont de dimensions différentes *cf* : *application développée dans ce chapitre*.

Nous proposons un formalisme basé sur des notions simples de la statistique descriptive et de la théorie de l'information pour construire un indice de distance entre ces objets matriciels et qui s'étend au cas de données qualitatives de manière plus générale.

Nous considérons pour ce faire, qu'une observation est une réalisation d'une variable aléatoire. Dans le paragraphe 2 de ce chapitre, nous introduisons la notion de système physique comme modèle qui régit la description de chaque objet. Le nombre des états de chaque système est fini et la probabilité que le système se trouve dans un des états est mesurée par sa fréquence. Nous utilisons l'entropie pour mesurer l'incertitude des états des systèmes Shannon[90] et enfin, nous définissons dans le paragraphe 3 une distance entre les éléments constitutifs de la structure.

Cette démarche a permis de construire une hiérarchie indicée entre les éléments constitutifs de la structure de juxtaposition de tableaux multiples.

Un exemple numérique, l'algorithme, une application sur données réelles et quelques perspectives de développement possibles complètent ce chapitre.

Le programme de l'algorithme écrit en visual basic sous excell est placé en annexe. Ce programme est déroulé sur l'exemple numérique.

### 2 Formalisme adapté

Soit  $\Omega = \{\omega_1, \dots, \omega_n\}$  un ensemble fini de  $N$  objets élémentaires, et  $\{V^1, \dots, V^d\}$   $d$  variables discrètes définies sur  $\Omega$  et prenant un nombre fini de valeurs respectivement dans  $D_1, \dots, D_d$  où  $D_j = \{m_1^j, \dots, m_{r_j}^j\}$ ,  $m_l^j$  est la  $l$ ème modalité de la variable  $V^j$ .

$$E_i^j = \begin{pmatrix} v_1^j \\ \vdots \\ v_{L_i}^j \end{pmatrix} \text{ où } v_l^j \in [m_1^j, \dots, m_{r_j}^j], \forall l = 1, \dots, L_i ; \forall j = 1, \dots, d$$

$E_i^j$  est le vecteur avec  $L_i$  composantes correspondantes aux différentes observations de l'individu  $\omega_i \in \Omega$ .  $E_i^j$  contient par exemple, les observations répétées de l'objet  $\omega_i$  pour la variable discrète  $V^j$  qui le décrit.

La structure de juxtaposition de tableaux de données catégorielles (*TofcD*) est

$$\Delta = [\Delta_1, \dots, \Delta_N] \text{ avec } \Delta_i = [\Delta_i^1, \dots, \Delta_i^d],$$

$\Delta_i$  est une matrice de dimension  $L \times M$ ,  $M = \sum_{j=1}^d r_j$  et  $\Delta_i^j$  est un tableau de données catégorielles défini comme suit

$$\Delta_i^j = \begin{bmatrix} \overbrace{m_1^j m_2^j \dots m_{r_j}^j}^{V^j} \\ 00000100 \\ 00010000 \\ \dots z_{lt}^j \dots \\ \dots \\ 00001000 \end{bmatrix}$$

$$Z_{it}^j = \begin{cases} 1 & \text{si } \omega_i \text{ a choisi la modalité } m_t^j \\ 0 & \text{ailleurs} \end{cases}$$

Pour  $\omega \in \Omega$ , on définit l'application

$$V_\omega : \Theta \longrightarrow V(\Omega) = \{m_1, \dots, m_q\}$$

$$(\theta_\omega)_t \longmapsto V((\theta_\omega)_t) = m_j$$

$m_j$  est la  $t$ ème observation de l'individu  $\omega \in \Omega$  pour la variable  $V$  et  $\Theta$  est l'ensemble qui contient les  $L$  observations de  $\omega$ . On pose

$$P[V((\theta_\omega)_t) = m_j] = P[V_\omega = m_j] = f_{V_\omega}^j \text{ et } \sum_{j=1}^q f_{V_\omega}^j = 1$$

#### Définition 4

La paire  $(\Theta, P_{V_\omega}) = S_\omega = [V \leftrightarrow P_{V_\omega} = P_\omega]$  est appelé : *système physique aléatoire simple*.

$S = [V \leftrightarrow P = F]$  décrit  $\omega \in \Omega$  si et seulement si la distribution de la variable  $V$  pour les observations de l'individu  $\omega$  est  $F$ , le système physique  $S$  s'écrit

$$S = \bigwedge_{j=1}^q [V = m_j; \Pr [V = m_j] = p_j],$$

où  $\wedge$  est la conjonction entre événements aléatoires.

Dans le cas multidimensionnel,

**Définition 5**

le système physique

$$S_\omega = (\Theta, P_{V_\omega}) = [V \hookrightarrow P_{V_\omega} = F_d^\omega] = [V = (V^1, \dots, V^d) \hookrightarrow F_d^\omega],$$

où

$$P_{V_\omega}(l_1, \dots, l_d) = P[V_\omega = (m_{l_1}^j, \dots, m_{l_d}^j)] = f_d^\omega(l_1, \dots, l_d)$$

$$\text{et } \sum_{l_1=1}^{r_1} \sum_{l_2=1}^{r_2} \dots \sum_{l_d=1}^{r_d} [f_d^\omega(l_1, \dots, l_d)] = 1$$

est appelé: *système physique multiple aléatoire*.

Soit  $F_\omega^j$  la *jième* distribution marginale de  $F_d$ , le système physique aléatoire multiple associé aux distributions marginales est

$$\widehat{S}_\omega = \widetilde{\wedge}_{j \in Q_\omega} [S_\omega^j],$$

$\widetilde{\wedge}$  conjonction entre systèmes physiques.

Les systèmes physiques simples  $\{S_\omega^j; j = 1\}$  sont donnés par

$$\forall j \in Q_\omega; S_\omega^j = [V^j \hookrightarrow F_\omega^j].$$

### 3 Distance entre systèmes physiques aléatoires multiples

#### 3.1 Entropie comme mesure de l'incertitude des états des systèmes

Pour mesurer le degré d'incertitude des états des systèmes physiques ou de la variable aléatoire discrète, on utilise une caractéristique spéciale couramment utilisée en théorie de l'information appelée : *entropie*

Soit  $S$  système physique aléatoire



$$S = (\Theta, P_V) = [V \hookrightarrow P_V = P] = \widetilde{\bigwedge}_{i=1}^q [(S) \hookrightarrow v^j; P [(S) \hookrightarrow v^j] = p_j].$$

Le symbole  $(S) \hookrightarrow v^j$  traduit le fait que le système se trouve à l'état  $v^j$ .

### Formulation de Shannon [1948] de l'entropie

#### Definition 6

L'entropie du système est la quantité positive

$$H(S) = - \sum_{j=1}^q p_j \log_2(p_j),$$

La fonction  $H$  possède quelques propriétés élémentaires qui justifie son utilisation comme caractéristique pour mesurer le degré d'incertitude du système

#### Proposition 12

1. Si un des états est certain ( $\exists l \in \{1, 2, \dots, q\}$  tel que  $:P_l = \Pr [(S) \hookrightarrow v_l] = 1$ ) alors  $H(S) = 0$ . En effet, si  $p_l = 1 \Rightarrow p_l \log_2(p_l) = 0 \Rightarrow \lim_{p \rightarrow 0} (p_l \log_2(p_l)) = 0 \Rightarrow H(S) = 0$ .
  2. L'entropie d'un système physique aléatoire avec un nombre fini d'états  $(m_1, m_2, \dots, m_q)$  est maximale si tous les états sont équiprobables:  $\forall j \in \{1, 2, \dots, q\}; p_j = \Pr [S \leftrightarrow m_j] = \frac{1}{q}$ . On a :  $0 \leq H(S) \leq \log_2(q)$ .
- La caractéristique de la fonction d'entropie exprime le fait que la distribution de probabilité avec le maximum d'entropie est la plus biaisée et la plus consistante avec l'information spécifiée par les contraintes.
  - Cette distribution est aussi la plus dispersée si il n'y a pas de contraintes supplémentaires que celles spécifiées dans le problème d'optimisation[90].

## 3.2 Entropie d'un système physique multiple aléatoire

Dans le cas multidimensionnel, Soit  $F_d$  une distribution multidimensionnelle d'ordre be probability  $d$ .  $F_d$  est donnée par la table

$$F_d \rightarrow \begin{array}{cccc} (m_1^1, m_1^2, \dots, m_1^d) & (m_2^1, m_2^2, \dots, m_2^d) & \dots & (m_{r_1}^1, m_{r_2}^2, \dots, m_{r_d}^d) \\ p_{11\dots 1} & p_{21\dots 1} & \dots & p_{r_1 r_2 \dots r_d} \end{array}$$

Le système physique associé est

$$S = \bigwedge_{l_1=1}^{r_1} \bigwedge_{l_2=1}^{r_2} \dots \bigwedge_{l_d=1}^{r_d} [V = (m_{l_1}^1, \dots, m_{l_d}^d); \Pr [V = (m_{l_1}^1, \dots, m_{l_d}^d) = p_{(l_1, \dots, l_d)}]] ,$$

où  
 $r_j$  est le nombre de valeurs prises par la variable aléatoire  $V^j$  c-à-d: le nombre d'états  
 et

$$\sum_{l_1=1}^{r_1} \sum_{l_2=1}^{r_2} \cdots \sum_{l_d=1}^{r_d} [p_{(l_1, \dots, l_d)}] = 1.$$

$S$  est appelé système physique multiple

Le système physique multiple aléatoire associé aux distributions marginales est donné par

$$(\widehat{S}) = \widetilde{\Lambda}_{j=1}^d [S_j].$$

Le système physique simple et discret  $S_j$  est

$$S_j = \bigwedge_{l=1}^{r_j} [V^j = m_l^j; \text{Pr} [V^j = v_l^j] = p_l^j] = [\Theta^j, P_{V^j}],$$

où

$$p_l^j = P [V^j = v_l] = \sum_{l_1=1}^{r_1} \sum_{l_2=1}^{r_2} \cdots \sum_{l_{j-1}=1}^{r_{j-1}} \sum_{l_{j+1}=1}^{r_{j+1}} \cdots \sum_{l_d=1}^{r_d} [p_{l_1, l_2, \dots, l_d}];$$

$$\text{avec } \sum_{l=1}^{r_j} p_l^j = 1$$

### Proposition 13

Si les systèmes physiques simples ( $(S^j) ; j \in Q_j$ ) sont indépendants, alors

$$H(S) = \sum_{j \in Q_i} H(S^j)$$

### Définition 7

Le système physique conditionnel  $S^1 / [S^2 \hookrightarrow m_l^2]$  est donné par

$$S^1 / [S^2 \hookrightarrow m_l^2] = \bigwedge_{j=1}^{r_1} [(S) \hookrightarrow m_j^1 / S^2 \hookrightarrow m_l^2; P [(S) \hookrightarrow m_j^1 / S^2 \hookrightarrow m_l^2] = p_{j/l}].$$

L'entropie de ce système est

$$H(S^1 / [S^2 \hookrightarrow m_l^2]) = - \sum_{j=1}^{r_1} p_{j/l} \log_2 (p_{j/l}).$$

Le système physique multiple aléatoire  $(S^1/S^2)$  s'écrit

$$(S^1/S^2) = \hat{\bigwedge}_{l=1}^{r_2} \left[ \hat{\bigwedge}_{j=1}^{r_1} [(S) \hookrightarrow m_j^1 / S^2 \hookrightarrow m_l^2; P [(S) \hookrightarrow m_j^1 / S^2 \hookrightarrow m_l^2] = p_{j/l}] \right],$$

$$\implies H(S^1/S^2) = - \sum_{l=1}^{r_2} \left[ \sum_{j=1}^{r_1} p_{j/l} \log_2(p_{j/l}) \right],$$

**Proposition 14**

$$H(S) = H(S^1) + H(S^2/S^1) + H(S^3/S^1 \hat{\bigwedge} S^2) + \dots + H(S^d/S^1 \hat{\bigwedge} S^2 \hat{\bigwedge} \dots \hat{\bigwedge} S^{d-1}),$$

Pour la démonstration se référer à [100].

**Proposition 15**

Les quantités

$$K(P, Q) = - \sum_{i=1}^n P_i \log_2(Q_i/P_i) \quad ; \quad K(Q, P) = - \sum_{j=1}^n P_j \log_2(P_j/Q_j),$$

sont non négatives et on a

$$\sum_{j=1}^n P_j \log_2(P_j) \geq \sum_{j=1}^n P_j \log_2(Q_j); \quad K(P, Q) = K(Q, P) \Leftrightarrow P = Q \text{ presque sûrement}$$

$K(., .)$  n'est une fonction symétrique donc ce n'est pas une distance au sens classique du terme mais caractérise du point de vue statistique la déviation entre les distributions  $P$  et  $Q$ .

Kullback[55] explique que la quantité  $K(P, Q)$  évalue la quantité d'information moyenne perdue si l'on utilise la distribution  $P$  alors que la vraie distribution est  $Q$ .

Soit  $S_{\Pi_d}$  l'ensemble des systèmes physiques multiples aléatoires avec  $\prod_{j=1}^d r_j$  états

$$S \in S_{\Pi_d} \implies (S) = (\Theta, P) = \hat{\bigwedge}_{l_1=1}^{r_1} \hat{\bigwedge}_{l_2=1}^{r_2} \dots \hat{\bigwedge}_{l_d=1}^{r_d} \left[ [(S) \hookrightarrow (m_{\cdot l_1}, \dots, m_{l_d})]; p_{l_1 l_2 \dots l_d} \right].$$

Soit  $\delta$  l'application définie par

$$S_{\Pi_d} \times S_{\Pi_d} \xrightarrow{\delta} \mathfrak{R}_+$$

$$((S_1), (S_2)) \longrightarrow \delta(S_1, S_2) = H(S_1) + H(S_2) - [\Psi(S_1/S_2) + \Psi(S_2/S_1)]$$

Où

$$\Psi(S_1/S_2) = K_d(P^{(1)}, P^{(2)})$$

$H$  est la fonction d'entropie,  $P^{(1)}$  et  $P^{(2)}$  sont des distributions multivariées de dimension  $d$  qui définissent les systèmes physiques multiples aléatoires  $S_1, S_2$  et  $K_d$  est définie comme suit

$$K_d(P^1, P^2) = \sum_{I_1} \sum_{I_2} \dots \sum_{I_d} P_{I_1 \dots I_d}^{(1)} \log_2 \left( P_{I_1 \dots I_d}^{(2)} \right).$$

### Proposition 16

$\delta$  vérifie

$$1. \delta(S_1, S_2) \geq 0$$

$$2. \delta(S_1, S_2) = 0 \Leftrightarrow S_1 = S_2$$

$$3. \delta(S_1, S_2) = \delta(S_2, S_1) \quad (\text{symétrie})$$

On admet que si  $\delta(S_1, S_2) = 0 \Leftrightarrow S_1 = S_2 \Leftrightarrow \omega_1 = \omega_2$ .

$\delta$  mesure la dissimilarité entre les systèmes physiques. Plus la valeur de  $\delta$  est petite plus l'incertitude des systèmes est grande.  $\delta$  représente la quantité d'information moyenne perdue si on utilise la distribution  $P^1$  ( $P^2$ ) pour décrire le système alors que c'est l'autre distribution qui est vraie .

$\delta$  est une nouvelle version de la distance de Kullback-Leibler entre les distributions  $P^1$  et  $P^2$ . En effet, La distance de Kullback-Leibler entre  $P^1$  et  $P^2$  est donnée par

$$Ku(P^1, P^2) = \sum_{I_1} \sum_{I_2} \dots \sum_{I_d} (P_{I_1 \dots I_d}^1 - P_{I_1 \dots I_d}^2) \log_2 \left( P_{I_1 \dots I_d}^1 / P_{I_1 \dots I_d}^2 \right).$$

En développant cette expression, on retrouve  $\delta$ .

## 4 Application numérique

### 4.1 Procédure pour estimer la distribution conjointe

Dans le cas où toutes les variables qui interviennent dans la description des individus sont discrètes, on donne une procédure inspirée de l'analyse factorielle pour estimer la distribution conjointe et en déduire la valeur de l'entropie des systèmes multiples aléatoires correspondants.

Soit  $\Delta_i = [\Delta_i^1, \dots, \Delta_i^d]$  une juxtaposition de  $d$  tableaux de données catégorielles . Pour  $\omega_i \in \Omega$  fixé, on a

$$P_{l_1 l_2 \dots l_d}^{(i)} = \Pr [S \hookrightarrow (m_{l_1}^1, \dots, m_{l_d}^d)] = \frac{d}{L} N^{(i)}(l_1, \dots, l_d) = f_{l_1 \dots l_d}^{(i)}$$

$N^{(i)}(\cdot)$  est le nombre d'occurrences simultanées des modalités  $m_{l_1}^1, m_{l_2}^2, \dots, m_{l_d}^d$ .

$$\sum_{l_1=1}^{r_1} \sum_{l_2=1}^{r_2} \dots \sum_{l_d=1}^{r_d} f_{l_1 \dots l_d}^{(i)} = 1.$$

## 4.2 Exemple numérique

On considère 6 individus décrits par 2 variables qualitatives à 2 et 3 modalités. 10 observations pour chaque individu.

**Table1:** La juxtaposition des matrices 0/1 associées aux individus

$\omega_1$					$\omega_2$					$\omega_3$				
$V^1$		$V^2$			$V^1$		$V^2$			$V^1$		$V^2$		
$m_1^1$	$m_2^1$	$m_1^2$	$m_2^2$	$m_3^2$	$m_1^1$	$m_2^1$	$m_1^2$	$m_2^2$	$m_3^2$	$m_1^1$	$m_2^1$	$m_1^2$	$m_2^2$	$m_3^2$
1	0	1	0	0	1	0	0	1	0	1	0	1	0	0
0	1	0	0	1	0	1	0	0	1	0	1	0	0	1
1	0	0	1	0	0	1	1	0	0	1	0	0	1	0
0	1	1	0	0	0	1	1	0	0	0	1	0	1	0
0	1	1	0	0	1	0	0	1	0	0	1	1	0	0
1	0	0	0	1	0	1	1	0	0	1	0	0	0	1
0	1	0	1	0	0	1	0	0	1	1	0	1	0	0
1	0	0	1	0	1	0	0	1	0	0	1	0	1	0
1	0	0	1	0	0	1	0	1	0	0	1	1	0	0
1	0	1	0	0	1	0	0	0	1	1	0	0	1	0
$\omega_4$					$\omega_5$					$\omega_6$				
$V^1$		$V^2$			$V^1$		$V^2$			$V^1$		$V^2$		
$m_1^1$	$m_2^1$	$m_1^2$	$m_2^2$	$m_3^2$	$m_1^1$	$m_2^1$	$m_1^2$	$m_2^2$	$m_3^2$	$m_1^1$	$m_2^1$	$m_1^2$	$m_2^2$	$m_3^2$
1	0	1	0	0	1	0	0	1	0	1	0	1	0	0
0	1	0	1	0	0	1	1	0	0	1	0	0	0	1
1	0	0	0	1	1	0	0	0	1	0	1	0	1	0
0	1	1	0	0	1	0	0	1	0	1	0	1	0	0
0	1	0	1	0	0	1	0	1	0	0	1	0	0	1
1	0	0	1	0	1	0	0	1	0	0	1	0	0	1
1	0	0	1	0	0	1	0	0	1	0	1	1	0	0
1	0	0	0	1	1	0	0	1	0	1	0	1	0	0
0	1	0	0	1	1	0	0	1	0	1	0	0	1	0
1	0	1	0	0	1	0	1	0	0	1	0	1	0	0

### 4.2.1 Procédure pour construire une hiérarchie indicée sur ces objets

**Table 2:** Les distributions empiriques qui représentent les individus .

	$(m_1^1, m_1^2)$	$(m_1^1, m_2^2)$	$(m_1^1, m_3^2)$	$(m_2^1, m_1^2)$	$(m_2^1, m_2^2)$	$(m_2^1, m_3^2)$
$F_1$	0.2	0.3	0.1	0.2	0.1	0.1
$F_2$	0.1	0.2	0.1	0.3	0.1	0.1
$F_3$	0.2	0.2	0.1	0.2	0.2	0.1
$F_4$	0.2	0.2	0.2	0.1	0.2	0.1
$F_5$	0.1	0.4	0.2	0.1	0.1	0.1
$F_6$	0.4	0.1	0.1	0.1	0.1	0.1

Le programme écrit en visual basic déroulé sur cet exemple a donné les résultats suivants:

**Table 3:** Les entropies des systèmes physiques conditionnelles associés aux 6 objets

	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$	$S_6$
$S_1$	2, 4464	2, 6049	2, 5219	2, 6219	2, 6219	2, 8219
$S_2$	2, 6049	2, 4464	2, 6219	2, 8219	2, 8219	2, 9219
$S_3$	2, 6049	2, 7049	2, 5219	2, 6219	2, 8219	2, 8219
$S_4$	2, 7049	2, 8634	2, 6219	2, 5219	2, 7219	2, 8219
$S_5$	2, 4879	2, 6634	2, 6219	2, 5219	2, 3219	3, 0219
$S_6$	2, 6634	2, 8634	2, 6219	2, 6219	3, 0219	2, 3219

*Etape 1.* A partir de la matrice de similarité, on obtient

$$\min_{l \neq t; l, t=1, \dots, 5} \{dist(S_l, S_t)\} = dist(S_1, S_3) = 0, 1585.$$

Alors les individus  $\omega_1$  et  $\omega_3$  sont agrégés en un individu artificiel  $\omega_7$  qui sera placé dans la dernière ligne de la matrice la ligne et la colonne correspondantes aux objets  $\omega_1$  et  $\omega_3$  sont supprimés

*Etape 2.* De la nouvelle matrice de similarité, on obtient

$$\min_{l \neq t; l, t=2, 4, 5, 6} \{dist(S_l, S_t)\} = dist(S_2, S_7) = 0, 317.$$

Les individus  $\omega_2$  et  $\omega_7$  sont agrégés en un individu artificiel  $\omega_8$  .

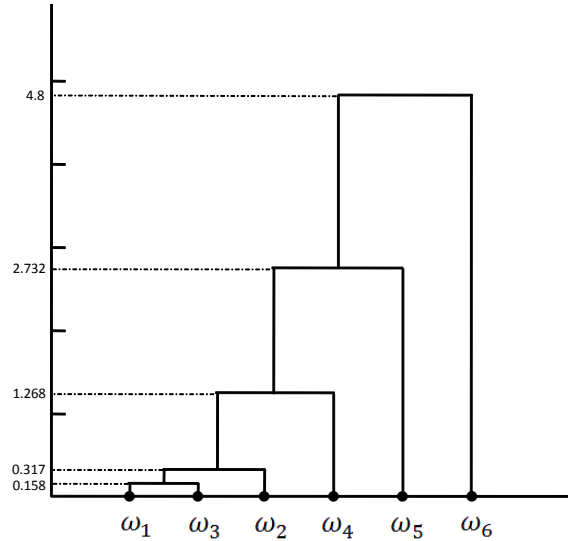
*Etape 3.*  $\min_{l \neq t; l, t=4, 5, 6} \{dist(S_l, S_t)\} = dist(S_4, S_8) = 1, 268.$

Les individus  $\omega_4$  et  $\omega_8$  sont agrégés en un individu artificiel  $\omega_9$  .

*Etape 4.*  $\min_{l \neq t; l, t=5, 6, 9} \{dist(S_l, S_t)\} = dist(S_5, S_9) = 2, 732.$

Les individus  $\omega_5$  et  $\omega_9$  sont agrégés en un individu artificiel  $\omega_{10}$ . Les individus  $\omega_6$  et  $\omega_{10}$  sont agrégés en un individu artificiel  $\omega_{11}$   $dist(S_6, S_{10}) = 4, 8.$

**La hiérarchie obtenue est**



- Dans cet exemple, on commence par regrouper les 2 objets les plus proches au sens de l'indice de distance entre les systèmes physiques correspondants.
- Plus on monte dans la construction de la hiérarchie plus les états des systèmes mélanges obtenus sont incertains.
- L'exemple montre que l'indice de distance de Kullback Leibler et l'indice d'aggregation du saut minimum conduisent à la construction de systèmes physiques avec le maximum d'entropie donc à un système physique où tous les états sont équiprobables.
- Si le nombre total des modalités des variables est assez grand comparativement au nombre d'observations, la fréquence de choisir telle ou telle modalité devient petite et plusieurs fréquences seront nulles: Les modalités de fréquences nulles seront élaguées des états du système et n'interviennent pas dans le calcul des distances. Cela peut conduire à l'impossibilité de comparer les systèmes donc les éléments de la structure.

<i>case number</i>	<i>cluster</i>	<i>distance</i>
$\omega_1$	1	4.201
$\omega_2$	2	8.040
$\omega_3$	1	3.598
$\omega_4$	2	8.040
$\omega_5$	1	3.324
$\omega_6$	1	3.4118

- Si l'on ne tient pas compte de la variabilité des variables, les individus  $\omega_2$  et  $\omega_4$  sont affectés à la première classe et sont équidistants du noyau de cette classe comme le montre la partition ci dessus obtenue dans le cas où chacun des 6 individus est décrit par sa valeur la plus fréquente pour chacune des variables qui le décrit .
- Si l'on utilise la procédure que nous avons développée et en coupant la hierarchie au niveau du 3<sup>ème</sup> palier, l'individu  $\omega_2$  et l'individu seront dans 2 classes différentes.

- Dans le cas où la variabilité des observations joue un rôle important particulièrement en médecine, dans la description des individus la classification faite sans tenir compte de cette variabilité peut conduire à des résultats erronés comparativement à la réalité des données.

### 4.3 Détails sur le déroulement de cet algorithme

Tableau contenant les entropies conditionnelles (T)  
 $T = [H(S_i/S_j)]_{\substack{i=1,6 \\ j=1,6}}$  et  $H(S_i/S_i) = H(S_i) \forall i = 1, 2, \dots, 6$

<b>2, 44643934</b>	2, 60493559	2, 52192809	2, 62192809	2, 62192809	2, 82192809
2, 60493559	<b>2, 44643934</b>	2, 62192809	2, 82192809	2, 82192809	2, 92192809
2, 60493559	2, 70493559	<b>2, 52192809</b>	2, 62192809	2, 82192809	2, 82192809
2, 70493559	2, 86343184	2, 62192809	<b>2, 52192809</b>	2, 72192809	2, 82192809
2, 48794309	2, 66343184	2, 62192809	2, 52192809	<b>2, 32192809</b>	3, 02192809
2, 66343184	2, 86343184	2, 62192809	2, 62192809	3, 02192809	<b>2, 32192809</b>

Matrice de dissimilarité à l'étape 0

$$\Delta^0 = \begin{matrix} & S_1 & S_2 & S_3 & S_4 & S_5 & S_6 \\ \begin{matrix} S_1 \\ S_2 \\ S_3 \\ S_4 \\ S_5 \\ S_6 \end{matrix} & \begin{pmatrix} 0 & & & & & & \\ 0,3169925 & 0 & & & & & \\ 0,15849625 & 0,35849625 & 0 & & & & \\ 0,35849625 & 0,7169925 & 0,2 & 0 & & & \\ 0,34150375 & 0,7169925 & 0,6 & 0,4 & 0 & & \\ 0,7169925 & 1,0169925 & 0,6 & 0,6 & 1,4 & 0 & \end{pmatrix} \end{matrix}$$

$\text{Min}_{l,t=1,6} \{\delta(S_l, S_t)\} = \delta(S_1, S_3) = 0,15849625$

Les objets 1 et 3 sont agrégés en un seul objet 7 et sont supprimés et remplacés par l'objet 7 que l'on a placé dans la dernière ligne.

Si  $A$  est une classe de systèmes physiques multiples aléatoires, alors

$$\Delta(S, A) = \text{Min} \{\delta(S, H) ; H \in A\}$$

$\Delta$  est l'indice d'agrégation entre classes d'objets de Sneath[?].

Matrice de dissimilarité à l'étape 1

$$\begin{matrix} & S_2 & S_4 & S_5 & S_6 & S_7 = S_1 \tilde{\wedge} S_3 \\ \begin{matrix} S_2 \\ S_4 \\ S_5 \\ S_6 \\ S_7 = S_1 \tilde{\wedge} S_3 \end{matrix} & \begin{pmatrix} 0 & & & & & \\ 0,7169925 & 0 & & & & \\ 0,7169925 & 0,4 & 0 & & & \\ 1,0169925 & 0,6 & 1,4 & 0 & & \\ 0,3169925 & 0,2 & 0,34150375 & 0,6 & 0 & \end{pmatrix} \end{matrix}$$

$\text{Min}_{l \neq t=2,4,5,6,7} \{\delta(S_l, S_t)\} = \delta(S_2, S_7) = 0,3169925$

Les objets 2 et 7 sont agrégés en un seul objet 8, ils sont supprimés et remplacés par l'objet 8 que l'on a placé dans la dernière ligne.



Matrice de dissimilarité à l'étape 2

$$\begin{array}{cccc}
 S_4 & S_5 & S_6 & S_8 = S_2 \tilde{\wedge} S_7 \\
 S_4 & 0 & & \\
 S_5 & 2,86797 & 0 & \\
 S_6 & 4,06797 & 5,6 & 0 \\
 S_8 & 1,26797 & 1,366015 & 2,4 & 0 \\
 \text{Min}_{l \neq t=4,5,6,8} \{ \delta(S_l, S_t) \} & = \delta(S_4, S_8) & = 1,26797
 \end{array}$$

Les objets 4 et 8 sont supprimés et remplacés par l'objet 9 placé dans la dernière ligne:

Matrice de dissimilarité à l'étape 3

$$\begin{array}{cccc}
 & S_5 & S_6 & S_9 = S_4 \tilde{\wedge} S_8 \\
 S_5 & 0 & & \\
 S_6 & 11,2 & 0 & \\
 S_9 = S_4 \tilde{\wedge} S_8 & 2,73203 & 4,8 & 0 \\
 \text{Min}_{l \neq t=5,6,8} \{ \delta(S_l, S_t) \} & = \delta(S_5, S_9) & = 2,73203
 \end{array}$$

Les objets 5 et 9 sont supprimés et remplacés par l'objet 10

L'objet 6 est donc agrégé à l'objet 10. et  $\delta(S_6, S_{10}) = 4.8$

#### 4.4 Application à des données réelles

Les données proviennent d'une étude du niveau de développement humain et économique de 48 départements

$$[R_1, \dots, R_{48}]$$

d'un pays donné. L'objectif est de construire une hiérarchie indicée sur les 48 départements selon leur niveau de développement économique et humain afin de faire des comparaisons entre départements et entre sous régions à l'intérieur des départements

Chaque département  $R_i$  regroupe  $L_i$  sous régions

$$[C_1^i, \dots, C_{L_i}^i] \left( \sum_{i=1}^n L_i = 1541 \right)$$

qui constitue un territoire.

$\forall i = 1, \dots, n$  et  $l = 1, \dots, L_i$ , on mesure l'indice de développement économique et l'indice de développement humain  $idE$  et  $idH$ . Ces 2 indices composites ont été développés par les experts du PNUD et dépendent de la situation géographique et de la spécificité de la sous région

$$0 \leq idE(C_i^i) \leq 1 \quad \text{et} \quad 0 \leq idH(C_i^i) \leq 1$$

Plus la valeur de l'indice est proche de 1 plus le développement économique ou humain de la sous régions est jugé satisfaisant. Ces indices ne sont pas calculés de la même façon selon que la sous région est classée rurale ou urbaine: Le classement fait à partir de ces valeurs n'a donc aucun sens.

Pour chaque  $i = 1, \dots, n$

$$R_i \xrightarrow{(L_i, 2)} \begin{bmatrix} & idE & idH \\ C_1^i & \cdot & \cdot \\ \vdots & \vdots & \vdots \\ C_{L_i}^i & \cdot & \cdot \end{bmatrix}$$

La structure de données n'est pas exploitable sous cette forme, il est nécessaire de la transformer. Les experts du PNUD découpent les observations en 5 classes ou quintiles et affectent chacune sous région au quintile qui correspond à sa valeur pour chaque des indices.

Nous déterminons ainsi, pour les 48 séries d'observations des 2 indices

$$idE \rightarrow q_{i1}^1, \dots, q_{i5}^1$$

$$idH \rightarrow q_{i1}^2, \dots, q_{i5}^2$$

$\Rightarrow$  Les intervalles interquintiles

$$\left\{ \begin{array}{l} \mathfrak{S}_{i1}^1, \dots, \mathfrak{S}_{i5}^1 \\ \mathfrak{S}_{i1}^2, \dots, \mathfrak{S}_{i5}^2 \end{array} \right.$$

Nous nous sommes pas intéressés par les valeurs des

$$\{q_{i1}^j, \dots, q_{i5}^j; j = 1, 2\}$$

mais aux appartenances des sous régions aux intervalles interquintiles.

Pour  $i = 1, \dots, n$ , le tableau  $R_i$  est transformé en une table de données binaires

$$T_i = \left( \begin{array}{c} \overbrace{idE} \\ \mathfrak{S}_1^1 \dots \mathfrak{S}_5^1 \\ 0 \ 0 \ 1 \ 0 \ 0 \\ 0 \ 0 \ 0 \ 0 \ 1 \\ \dots \\ 0 \ 1 \ 0 \ 0 \ 0 \\ 0 \ 0 \ 1 \ 0 \ 0 \end{array} \right) \left( \begin{array}{c} \overbrace{idH} \\ \mathfrak{S}_1^2 \dots \mathfrak{S}_5^2 \\ 0 \ 0 \ 0 \ 1 \ 0 \\ 0 \ 1 \ 0 \ 0 \ 0 \\ \dots \\ 0 \ 0 \ 0 \ 1 \ 0 \\ 0 \ 0 \ 1 \ 0 \ 0 \end{array} \right)$$

Exhiber les éventuelles disparités entre sous régions à l'intérieur des départements. Les observations sont regroupées dans les tableaux  $T_1, \dots, T_n$  qui constituent une structure de juxtaposition de données catégorielles.

## 5 Conclusion

- Dans ce chapitre, la définition de l'entropie est celle émise par Shannon Weiner[1936]. Cette définition est encore utilisée dans la théorie de l'information et du signal.
- Le formalisme proposé donne une nouvelle explication de l'aspect pratique de la distance de Kullback-Leibler type comme un indice de distance entre les éléments représentatifs d'une structure de tableaux de données catégorielles.
- Il est peut être possible d'étendre les résultats obtenus dans le cas d'une structure de juxtapositions de tableaux de mesure et dans le cas de données fonctionnelles.

# Chapitre V

## Classification des données avec erreurs de mesure

### 1 Introduction

- Lors des observations répétées, nous obtenons généralement des dispersions dans les résultats. Si les résultats sont aléatoires un traitement statistique permet de connaître la valeur la plus probable de la grandeur mesurée et de fixer les limites de l'incertitude.
- Ainsi, si l'on est conduit à effectuer pour chaque objet plusieurs mesures et qu'on obtienne à chaque fois un résultat différent ceci est certainement dû à des phénomènes perturbateurs: dispersion statistique. Cette dispersion statistique désigne le fait que les observations sont proches les unes des autres ou au contraire éparpillées.
- Une structure de juxtaposition de tableaux de mesure peut donc être transformée en une structure de tableau de données avec erreurs de mesure. Sous l'hypothèse raisonnable de normalité des observations. Nous apportons un autre éclairage en proposant 2 approches pour classifier des objets matriciels selon que les variables actives soient corrélées ou non. La première approche traite le cas où les variables sont non corrélées donc indépendantes tandis que la seconde traite le cas où les variables sont quelconques.

Les procédures proposées consistent à transformer dans les 2 cas, ces objets matriciels en objets intervalles. Ainsi, la valeur la plus probable de chaque individu  $\omega_i$  est donnée à l'intérieur d'un sous ensemble  $R_i$  de  $\mathbb{R}^d$ . Les sous ensembles  $R_i$  sont des pavés, ou des ellipsoïdes de  $\mathbb{R}^d$  estimés de manière classique sous la condition que  $\Pr[x_i \in R_i] \succeq 1 - \beta$  avec  $\beta$  seuil fixé.

Si les variables ayant participé à la description des individus, sont corrélées donc quelconques, ces ensembles sont des ellipsoïdes. Dans chaque cas, nous adaptons un indice de distance et justifions son utilisation.

- Dans le second cas, plus général, nous avons utilisé la distance pondérée de Mahalanobis et avons construit un algorithme de type k-means.
- Nous montrons qu'une fois les clusters obtenues, leur caractérisation n'est pas entachées d'erreurs et est donnée de manière exacte. Ce qui paraît tout à fait dans l'ordre des choses.

Le premier cas se ramène à la construction d'un algorithme de type k-means avec la distance classique euclidienne entre sommets des pavés qui décrivent les individus. Nous étudions plusieurs cas selon que les erreurs soient placées au niveau des observations et ou au niveau des noyaux des classes.

Habituellement, pour analyser les éléments constitutifs d'une structure de tableaux multiples, on résume chaque colonne de chaque tableau par un nombre. Les objets matriciels sont transformés en vecteurs et la structure devient un tableau classique définissant la description d'un nombre fini d'individus par un nombre fini de variables. Cependant, cette synthèse nécessite plusieurs hypothèses quant au choix de la valeur qui résume les différentes observations de l'individu pour la variable [voir chapitre II].

Le cas où chaque objet est décrit par une matrice est très fréquent en pratique, il génère des problèmes mathématiques intéressants et reste un vaste domaine de recherche et d'investigation.

Les structures de données de ce type sont très fréquentes en pratique. On peut citer

1. Le cas de données concernant la pollution engendrée par les automobilistes au niveau d'un important carrefour de trafic routier d'une grande métropole. Les taux des principaux polluants sont mesurés chaque heure de la journée sur toute l'année. Les principaux polluants peuvent être : le taux de poussière, d'ozone, de monoxyde de carbone, de monoxyde de zinc et de dioxyde de soude etc... Chaque jour est décrit par une matrice de données quantitatives de dimension 24 fois le nombre de polluants. On désire regrouper les jours selon les seuils de pollution dont les classes obtenues peuvent être expliquées par d'autres variables concernant par exemple la météorologie.
2. En médecine, si l'on dispose d'une population finie de patients atteints d'une longue maladie où chaque malade subit des contrôles réguliers sur l'état et de l'évolution de sa maladie sur une grande période. Il s'agit là aussi de regrouper les malades selon l'évolution de la maladie.

Pour ce type de données, la phase de regroupement des observations par le cumul journalier ou par la moyenne journalière pour chaque polluant dans le premier exemple peut conduire à des résultats erronnés. En effet, les pics de pollution dépendent des heures de pointe, des jours ouvrables et d'autres événements qui peuvent soit ralentir ou augmenter la densité de la circulation donc le seuil de pollution .

De même qu'en médecine, la moyenne journalière ou le cumul journalier de la tension artérielle n'ont pas de sens et ce sont ses variations qui peuvent dénoter des cardiopathies. Aussi, il est nécessaire de traiter l'information de la manière la plus proche de la réalité.

Cette structure a été présentée dans [1] et n'a pas eu l'attention à la mesure de sa fréquence en pratique. Dans Rebbouh[80], est proposé un algorithme du type k-means basé sur le produit scalaire de Hilbert-Schmidt pour classifier des d'objets matriciels de même dimension.

## 2 Estimation des sous ensembles $R_i$

Soit  $\Omega = \{\omega_1, \dots, \omega_N\}$  un ensemble fini de  $N$  individus à classifier. On suppose que pour tout  $i \in \{1, \dots, N\}$ ,  $\omega_i$  est décrit par  $d$  variables quantitatives dont les observations sont regroupées dans le tableau  $T_i$  de dimension  $N_i \times d$ .

$$\omega_i \rightarrow T_i = \begin{bmatrix} \theta_1^i \\ \vdots \\ \theta_l^i \\ \vdots \\ \theta_{N_i}^i \end{bmatrix} \begin{bmatrix} V^1 & \dots & V^d \\ x_{i1}^1 & \dots & x_{i1}^d \\ \vdots & \vdots & \vdots \\ x_{il}^1 & \dots & x_{il}^d \\ \vdots & \vdots & \vdots \\ x_{iN_i}^1 & \dots & x_{iN_i}^d \end{bmatrix},$$

$x_{il}^j$  est la  $l^{\text{ième}}$  observation de l'individu  $\omega_i$  pour la variable  $V^j$

$$T_i = [ X_1^i \quad \dots \quad X_{N_i}^i ]; X_l^i \in \mathbb{R}^d,$$

$T_i$  est un objet matriciel qui regroupe les observations répétées de l'individu  $\omega_i$  pour les  $d$  variables quantitatives qui le décrivent. Ces objets matriciels peuvent être de dimensions différentes. La structure

$$\Delta = [T_1, \dots, T_n],$$

est une structure de juxtaposition de tableaux de mesure (*TofM*).

## 2.1 Les $d$ variables sont non corrélées

### 2.1.1 Les intervalles de confiance associés

Pour un seuil de signification  $\beta$  fixé, on construit les intervalles de confiance associés aux observations variées de chaque individu pour les variables qui le décrivent même si le nombre des observations diffère d'un individu à l'autre et que ces nombres soient relativement petits.

### 2.1.2 Premier cas $\sigma_i^j$ sont connus pour $i = 1, \dots, N$ et $j = 1, \dots, d$

L'individu  $\omega_i$  est observé à l'intérieur de l'intervalle  $I_i^j$  de longueur  $2\varepsilon_i^j(\beta)$ , pour  $j = 1, \dots, d$ ,

#### Proposition 17

Cet intervalle de confiance au niveau de signification  $\beta$  fixé est donné par

$$I_i^j(\beta) = ]\mu_i^j - \varepsilon_i^j(\beta) \ ; \ \mu_i^j + \varepsilon_i^j(\beta)[,$$

où

$$\varepsilon_i^j(\beta) = \sqrt{\frac{\sigma_i^j}{N_i - 1}} \arg \Phi^* \left( \frac{1 + \beta}{2} \right),$$

$\Phi^*(y)$  est la valeur pour laquelle la fonction cumulative *cdf* de la distribution normale est  $y$  et

$$\left\{ \begin{array}{l} \mu_i^j = \frac{1}{N_i} \sum_{l=1}^{L_i} x_{il}^j, \\ \sigma_i^j = \left( \frac{1}{N_i - 1} \sum_{l=1}^{N_i} [x_{il}^j - \mu_i^j]^2 \right)^{1/2}. \end{array} \right.$$

### 2.1.3 Deuxième cas $\{\mu_i^j ; j \succeq 1\}$ sont connus et $\sigma_i^j \neq 0 ; \forall j \succeq 1$

De la même manière, nous obtenons les intervalles de confiance estimés des variances au seuil de signification  $\beta$  fixé sont

$$I_i^j(\beta) = ]-t_\beta \cdot D_i^j ; +t_\beta \cdot D_i^j[$$

où

$$D_i^j = \sqrt{\frac{0.8 \cdot N_i + 1/2}{N_i(N_i - 1)} \sigma_i^j} \quad \text{et} \quad t_\beta = \arg \Phi^* \left( \frac{1 + \beta}{2} \right).$$

Il est à noter que ces estimateurs dépendent du nombre d'observations. Cependant, si ces matrices n'ont pas la même dimension, on peut envisager une étape de complétion. En effet, si les tableaux

$$\{T_i ; i = 1, \dots, N\}$$

n'ont pas la même dimension  $\Leftrightarrow [\exists i \neq i' \text{ tel que } N_i \neq N_{i'}]$ . Soit  $L$  le plus petit commun multiple des

$$\{N_i ; i = 1, \dots, N\}$$

$$P = PPCM(N_i, i = 1, N),$$

Il existe donc  $r_i$  tel que

$$p = N_i \times r_i.$$

Maintenant, on duplique  $r_i$  fois chaque tableau  $T_i$ , on obtient ainsi, un nouveau tableau  $\widehat{T}_i$  de dimension  $P \times d$ . Si les  $N_i$  sont des grands nombres, le plus petit commun multiple devient très grand et la procédure conduira nécessairement à une structure de données de grands tableaux. De plus, cette technique de complétion détruit complètement la chronologie des observations. Ce n'est donc pas une technique recommandée. Il est plus raisonnable de ne pas procéder à une étape de complétion avant d'entamer la phase de classification.

L'objet  $\omega_i$  est donc décrit par le pavé  $R_i$  de  $\mathbb{R}^d$  qui s'écrit

$$R_i = \prod_{j=1}^d I_i^j(\beta) ; I_i^j(\beta) \subset \mathbb{R}(I_i^j(\beta) \text{ est un intervalle})$$

## 2.2 Les $d$ variables sont quelconques

Deux approches peuvent être envisagées:

1. La première approche consiste à utiliser une procédure heuristique basée sur les techniques factorielles de réduction. Cette procédure permet de construire de nouvelles variables non corrélées à partir des variables initiales pour chacun des objets matriciels en présence. A ce stade, deux problèmes importants sont posés. Le premier problème concerne le choix de la dimension ou du nombre de nouvelles variables à retenir. Le second problème concerne le choix du sous espace compromis ou commun qui permet de comparer les différents éléments de la structure et donc de procéder à leur classification.
2. La seconde approche, plus appropriée au sens où elle évite les 2 problèmes cités précédemment et qui demeurent non totalement résolus, consiste à estimer directement le couple  $(\mu_i, \Sigma_i)$  qui résumera la description de l'individu  $\omega_i$ .

### 2.2.1 Première approche: Choix du sous espace compromis et estimation des intervalles de confiance

- L'analyse en composantes principales déroulée sur chaque tableau  $T_i$  conduit à la construction de  $r_i$  axes factoriels orthogonaux sur lesquels les projections des  $N_i$  observations de  $\omega_i$  donnent de nouvelles variables non corrélées qui sont les composantes principales. Nous obtenons ainsi, en déroulant  $N$  analyses,  $N$  systèmes d'axes

$$\left\{ \left\{ \Delta_{u_1^{(i)}}, \dots, \Delta_{u_{r_i}^{(i)}} \right\}; i = 1, \dots, N \right\}$$

- Dans le but de comparer les objets et de pouvoir les regrouper en classes, on doit les représenter dans un même référentiel. Ce problème fondamental de la recherche de système compromis ou commun des  $N$  systèmes, fait encore l'objet de recherche en mathématiques et intéresse particulièrement la géométrie différentielle, on cite les travaux de Orlick.P and Tera.H[74][1992]. Dans[10] sont proposés des critères pour le choix de ce référentiel mais les procédures proposées sont difficilement justifiables et ne constituent pas l'objet de notre étude.
- La procédure heuristique que nous proposons s'inspire des techniques utilisées dans le cas de la représentation simultanée des individus et des variables dans un même référentiel. Il s'agit de construire  $d^*$  axes pour l'analyse du tableau  $T_1$ . Ces  $d^*$  axes, dont le nombre est choisi par une procédure simple voir (\*), constitueront le repère de référence et on plongera les autres axes dans ce système. On effectuera donc  $N - 1$  rotations de centre  $0_{\mathbb{R}^d}$  et d'angles

$$\left\{ \theta_l^{(1)}; l = 1, \dots, N - 1 \right\}$$

La projection de la  $i$ ème observation de l'individu  $\omega_i$  sur le  $r$ ème axe factoriel est donnée par

$$W_i^r = T_i \cdot u_r,$$



où  $u_r$  est le vecteur propre normé associé à la rième plus grande valeur propre  $\lambda_r$  de la matrice d'inertie  $\Sigma_i$ . Nous supposons que le tableau  $T_i$  est centré

$$\Sigma_i = \frac{1}{N_i} (T_i)' (T_i),$$

$$d^* = \text{Min} \{r_i, i = 1, \dots, n\} (*)$$

où  $r_i$  est le plus petit nombre positif tel que

$$\frac{\lambda_1^i + \lambda_2^i + \dots + \lambda_{r_i}^i}{\text{tr}(\Sigma_i)} \geq a,$$

$a$  est fixé ( $a = 0.7$ ),  $\lambda_l^i$  est la lième plus grande valeur propre de la matrice d'inertie  $\Sigma_i$ . Les  $d^*$  nouvelles variables sont centrées et données par

$$\begin{cases} W_{(i)}^l = T_i \cdot u_l & \forall l = 1, \dots, d^* \\ W_i^l = 0 \end{cases}$$

Les intervalles de confiance  $I_i^j$  qui décrivent l'individu  $\omega_i$  sont estimés à partir des observations contenues dans la matrice

$$Y^{(i)} = \begin{bmatrix} W^1 & \dots & W^j & \dots & W^{n^*} \\ y_{i1}^1 & \dots & y_{i1}^j & \dots & y_{i1}^{n^*} \\ \vdots & \vdots & \vdots & & \vdots \\ y_{il}^1 & \dots & y_{il}^j & \dots & y_{il}^{n^*} \\ \vdots & \vdots & \vdots & & \vdots \\ y_{iN_i}^1 & \dots & y_{iN_i}^j & \dots & y_{iN_i}^{n^*} \end{bmatrix}$$

$$I_i^j(\beta) = ] -t_\beta \cdot D_i^j ; + t_\beta \cdot D_i^j [$$

où

$$D_i^j = \sqrt{\frac{0.8 \cdot N_i + 1/2}{N_i (N_i - 1)}} \sigma_i^j; \quad t_\beta = \arg \Phi^* \left( \frac{1 + \beta}{2} \right).$$

où

$$\sigma_i^l = \left( \frac{1}{N_i - 1} \sum_{t=1}^{L_i} [(T_i) \cdot (u_l)_t]^2 \right)^{1/2} = \left( \frac{1}{N_i - 1} \sum_{t=1}^{L_i} (y_{it}^j)^2 \right)^{\frac{1}{2}}$$

### 2.2.2 Seconde approche: Estimation des couples $(\mu_i, \Sigma_i)$

Cette approche fera l'objet d'une étude détaillée car elle nous paraît plus générale de par ses conditions d'application. Les espérances mathématiques  $\{\mu_i^j; j = 1, \dots, d\}$  et les composantes des matrices de variances covariances estimées  $\Sigma_i$  sont données par

$$\mu_i = \begin{pmatrix} \mu_i^1 \\ \vdots \\ \mu_i^d \end{pmatrix} \text{ avec } \mu_i^j = \frac{1}{L_i} \sum_{l=1}^{L_i} x_{il}^j$$

$$(\Sigma_i)_{jk} = Cov(V_{(i)}^j, V_{(i)}^k) = \Sigma_{ijk}^i = \frac{1}{N_i - 1} \sum_{l=1}^{L_i} (x_{il}^j - \mu_i^j) (x_{il}^k - \mu_i^k),$$

Nous assurons dans tous les cas que la condition de normalité des données est vérifiée et que

$$P[x_i \in R_i] \succeq 1 - \beta$$

De plus, ces estimateurs sont sans biais, convergents et consistants même si le nombre  $N_i$  d'observations est assez petit. Il faut que  $N_i \succ d; \forall i = 1, N$  pour que les matrices de variance covariances estimées soient régulières. De plus, pour que la normalité des observations soit envisageable, il faut que ces nombres dépassent 30.

## 3 Choix de l'indice de distance entre éléments représentatifs

### 3.1 Les individus sont décrits par des pavés de $\mathbb{R}^d$

Soit  $\mathfrak{R}_{V_d}$  l'ensemble des pavés droits de  $\mathbb{R}^d$

$$R \in \mathfrak{R}_{V_d} \iff R = \prod_{j=1}^d \mathfrak{S}^j \text{ où } \mathfrak{S}^j \subset \mathbb{R}; \mathfrak{S}^j = [x_{j-s_j}; x_{j+s_j}]$$

#### Proposition 18

Le nombre de sommets présents dans un pavé droit de  $\mathbb{R}^d$  qui représente l'ordre de ce graphe vaut  $2^d$ .

A chaque

$$R = \prod_{j=1}^d [[x_{j-\sigma_j}; x_{j+\sigma_j}]]$$

On lui associe les sommets  $S_t^{(\cdot)} \in \mathbb{R}^d$

$$\left(S_t^{(\cdot)}\right)' = [a_1^t, \dots, a_d^t] \in \mathbb{R}^d \text{ avec } a_l^t = \begin{cases} x_{l-s_l} \\ \text{ou} \\ x_{l+s_l} \end{cases}$$

### 3.1.1 Indice de distance entre objets intervalles

#### Proposition 19

L'application  $\delta$  définie par

$$\delta : \mathfrak{R}_{V_d} \times \mathfrak{R}_{V_d} \longmapsto \mathbb{R}^+$$

$$(R_1, R_2) \longmapsto \delta (R_1, R_2) = \frac{1}{2^d} \sum_{t=1}^{2^d} \left[ \left( \left\| S_t^{(1)} - S_t^{(2)} \right\| \right) \right],$$

où  $S_t^{(1)}, S_t^{(2)} \in \mathfrak{R}^d$  et  $\|\cdot\|$  est la norme euclidienne,  $\delta$  est une distance entre objets intervalles.

On suppose que

$$\delta (R_1, R_2) = 0 \Leftrightarrow \omega_1 = \omega_2$$

### 3.1.2 Indice d'aggrégation entre objets et pavés de $\mathbb{R}^d$

L'application  $D$

$$D : \mathbb{R}^d \times \mathfrak{R}_{V_d} \longmapsto \mathbb{R}^* :$$

$$(X_i, R) \longmapsto D (X_i, R) = \frac{1}{2^d} \sum_{t=1}^{2^d} \left\| X_i - S_t^{(\cdot)} \right\| = \frac{1}{2^d} \sum_{t=1}^{2^d} \left( \sum_{j=1}^d \left( (X_i)_j - (S_t^{(\cdot)})_j \right)^2 \right)^{\frac{1}{2}}$$

Où  $S_t^{(\cdot)} \in \mathbb{R}^d$  est le  $t^{\text{ième}}$  sommet du pavé  $R$ ,  $(X_i)_j$  et  $(S_t^{(\cdot)})_j$  sont respectivement la  $j^{\text{ième}}$  coordonnée de  $X_i$  et  $S_t^{(\cdot)}$ .  $D$  mesure la distance moyenne entre le pavé  $R$  et le vecteur  $X_i \in \mathbb{R}^d$ .

## 3.2 Les individus sont décrits par des ellipsoïdes de $\mathbb{R}^d$

### 3.2.1 Distance de Mahannalobis pondérée.

#### Définition 8

Soit  $M_d(\mathbb{R})$  l'ensemble des matrices réelles définies positives, l'application  $d_{ma}$  définie par

$$(\mathbb{N} \times \mathbb{R}^d \times M_d(\mathbb{R})) \times (\mathbb{N} \times \mathbb{R}^d \times M_d(\mathbb{R})) \rightarrow \mathbb{R}$$

$$(T_i, T_l) \rightarrow d_{ma}(T_i, T_l) = \frac{N_i \|\mu_i - \mu_l\|_{\Sigma_i^{-1}} + N_l \|\mu_i - \mu_l\|_{\Sigma_l^{-1}}}{N_i + N_l}$$

est la distance pondérée de Mahannolobis .

- Cette distance tient compte non seulement de la moyenne des observations et de leur nombre mais aussi de leur variabilité.
- La distance de Mahanalobis (nommée d'après **Prasantra Chandra Mahanalo-bis**[1936] scientifique indien, chercheur en statistique) permet de comparer deux échantillons . En effet, si leur précision est grande, alors une différence entre leurs moyennes est significative, et on peut les considérer comme différents. En revanche, si la précision est faible, donc la dispersion est grande, alors une différence entre leurs moyennes n'est pas significative, elle est peut-être uniquement dûe au manque de précision.
- Comme toujours en statistique, on peut parfois affirmer que deux quantités diffèrent mais jamais qu'elles sont vraiment égales : on peut parfois dire non, jamais oui.
- Mathématiquement, cette distance quadratique est la différence entre les deux moyennes au carré, pondérée par leur précision (plus précisément, divisée par la somme de leur variance).

On suppose de plus que

$$1. d(T_i, T_l) = 0 \Leftrightarrow \omega_i = \omega_l$$

$$2. N_i + N_l \succ d.$$

Cette seconde hypothèse assure que les matrices  $\Sigma_i$  et  $\Sigma_l$  soient singulières.

## 4 Classification d'objets décrits avec erreurs de mesure

### 4.1 Objets décrits par des pavés de $\mathbb{R}^d$

Le modèle des données est de la forme

$$\Delta = [R_1, \dots, R_N],$$

ou  $R_i$  est un pavé de  $\mathbb{R}^d$ .

Les valeurs prises par l'individu  $\omega_i$  pour les  $d$  variables quantitatives sont dans un pavé  $R_i$  de  $\mathbb{R}^d$ .

Si les variables sont non corrélées, on peut utiliser des techniques de l'analyse factorielle pour réduire chaque tableau  $T_i$  et construire ainsi de nouvelles variables non corrélées, on retrouve ainsi le précédent cas. En général, les techniques de l'analyse factorielle nécessitent des critères pour le choix d'un espace commun qui permettent de comparer les éléments réduits pour chacun des tableaux. Dans notre cas, ces réductions conduisent à des pavés centrés à l'origine de  $\mathbb{R}^d$  ce qui les rend comparables et cela nous permet de comparer les éléments de la structure des données.

#### 4.1.1 Procédure

La procédure suivante montre comment construire ces nouvelles variables non corrélées

Soit  $T_i$  le tableau contenant les observations de l'individu  $\omega_i$

$$\omega_i \longrightarrow T_i = \begin{bmatrix} V^1 & \cdot & V^j & V^d \\ x_{i1}^1 & \cdot & x_{i1}^j & \dots x_{i1}^d \\ \vdots & \dots & \vdots & \dots \vdots \\ x_{iL_i}^1 & \dots & x_{iL_i}^j & \dots x_{iL_i}^d \end{bmatrix} \longrightarrow T_i = \begin{bmatrix} X_i^1 \\ \vdots \\ X_i^{L_i} \end{bmatrix}; (X_i^l)' \in (\mathfrak{R}^d, Jd)$$

#### Proposition 20

La projection de la  $lième$  observation de l'individu  $\omega_i$  sur le  $rième$  axe factoriel est donnée par

$$W_i^r = T_i \cdot u_r,$$

où  $u_r$  est le vecteur propre normalisé associé à la  $rième$  plus grande valeur propre  $\lambda_r$  de la matrice d'inertie  $\Sigma_i$ ,  $T_i$  est un tableau centré et

$$\Sigma_i = \frac{1}{L_i} (T_i)' (T_i),$$

Soit

$$n^* = \text{Min} \{ r_i, i = 1, n \}$$

où  $r_i$  est le plus petit nombre positif tel que

$$\frac{\lambda_1^i + \lambda_2^i + \dots + \lambda_{r_i}^i}{\text{tr}(\Sigma_i)} \geq a,$$

$\lambda_l^i$  est la  $lième$  plus grande valeur propre de la matrice d'inertie  $\Sigma_i$ . Les  $n^*$  nouvelles variables sont centrées et données par

$$\forall l = 1, n^*$$

$$\begin{cases} W_{(i)}^l = T_i \cdot u_l \\ \overline{W}_i^l = 0 \end{cases} .$$

Cependant, l'intervalle  $\mathfrak{S}_i^l$  qui décrit l'individu  $\omega_i$  pour les nouvelles variables  $W_{(i)}^l$  est estimé par

$$\mathfrak{S}_i^l = \left[ -\frac{1}{2}\sigma_i^l; +\frac{1}{2}\sigma_i^l \right],$$

où

$$\sigma_i^l = \left( \frac{1}{L_i - 1} \sum_{t=1}^{L_i} [(T_i) \cdot (u_l)_t]^2 \right)^{1/2} .$$

#### 4.1.2 Problème d'optimisation

On désire regrouper les  $N$  individus en  $K$  classes homogènes. Les valeurs prises par les individus sont soit à l'intérieur d'un pavé de  $\mathbb{R}^d$  soit des vecteurs de  $\mathbb{R}^d$ . L'homogénéité des classes se mesure à l'aide d'un critère de la somme des inertias intra-classes. Ce critère s'exprime

$$C_r(P, L) = \sum_{k=1}^K H_k = \sum_{k=1}^K \left[ \sum_{\{l\} \in P_k} \delta(x_l, l_k) \right],$$

$l_k$  est le noyau de la classe  $P_k$ ,  $x_l$  est l'observation de l'individu  $\omega_l$  et  $\delta$  est un indice de distance entre objets et éléments représentatifs des classes.

Soit  $IP_k, IL_k$  respectivement l'ensemble des  $k$ -partitions, l'ensemble des  $k$ -représentatifs éléments des classes

$$P \in IP_k \Rightarrow P = \{P_1, \dots, P_k\}$$

et

$$L \in IL_k \Rightarrow L = (l_1, \dots, l_k)$$

Le problème de classification se réduit à la recherche d'une partition à  $k$  classes et de ses  $k$  noyaux qui minimisent le critère  $C_r$ . Il s'agit de résoudre le problème d'optimisation suivant

$$\underset{P \in IP_k, L \in IL_k}{Min} C_r(P, L).$$

Les algorithmes habituellement utilisés pour résoudre ce type de problème sont du type  $k$ -means. Ces algorithmes consistent à définir une fonction d'affectation  $f$  et une fonction de représentation  $g$ . Les fonctions  $f$  et  $g$  sont utilisées alternativement pour faire décroître le critère.  $f$  et  $g$  doivent vérifier les conditions suivantes

$$f : IL_k \mapsto IP_k$$

$$L \in IL_k \mapsto f(L) = P = \{P_1, \dots, P_k\}$$

$$\underset{P \in IP_k}{Min} [C_r(P, L)] = C_r(f(L), L)$$

$$g : IP_k \mapsto IL_k$$

$$P \mapsto g(P) = L = (l_1, \dots, l_k)$$

$$\underset{L \in IL_k}{Min} [C_r(P, L)] = C_r(P, g(P))$$

Dans ce chapitre, nous présentons une approche pour classifier des objets dont les descriptions sont entachées d'erreurs de mesure et nous étudions trois cas.

### 4.1.3 Algorithme de classification

On choisit  $L^{(0)}$  et on utilise alternativement les fonctions  $f$  et  $g$ . L'algorithme se déroule comme suit

$$L^{(0)} \xrightarrow{f} P^{(1)} \xrightarrow{g} L^{(1)} \xrightarrow{f} P^{(2)} \xrightarrow{g} \dots \xrightarrow{f} P^{(n)} \xrightarrow{g} L^{(n)} \mapsto \dots P^{(*)} \xrightarrow{g} L^{(*)}$$

L'algorithme s'arrête dès la partition générée ne change plus. On construit ainsi 2 suites  $V_n$  et  $U_n$ .

**Proposition 21**

La suite

$$U_n = C_r(P^{(n)}, L^{(n)})$$

converge en décroissant et la suite

$$V_n = (P^{(n)}, L^{(n)})$$

est stationnaire à partir d'un certain rang

Preuve:[32]

## 4.2 Etude de cas

### 4.2.1 Alternative 1: Les erreurs sont cummulées au niveau des noyaux des classes

Soit  $\mathfrak{R}_{V_d}$  l'ensemble des pavés de  $\mathbb{R}^d$

$$R \in \mathfrak{R}_{V_d} \iff R = \prod_{j=1}^d \mathfrak{S}^j \text{ où } \mathfrak{S}^j \subset \mathbb{R}.$$

Les objets sont décrits par des vecteurs de  $\mathbb{R}^d$ . Les noyaux des classes sont observés à l'intérieur de pavés de  $\mathbb{R}^d$ . Pour chaque pavé  $R \in \mathfrak{R}_{V_d}$ , on associe l'ensemble

$$S_d = \{ \text{som}(\cdot, t) = s_t \in \mathbb{R}^d, t = 1, \dots, 2^d \},$$

#### Proposition 22

$$|S_d| = 2^d$$

Pour  $d = 2 \iff |S_2| = 2^2 = 4 \Rightarrow S_2 = \{ \text{som}(\cdot, t); t = 1, 2, 3, 4 \}$

Les sommets associés pour ces 2 intervalles  $\{ \mathfrak{S}_i = ]m_i; M_i[; i = 1, 2 \}$  sont

$$\left\{ \begin{array}{l} \text{som}(\cdot, (1)) = s_1 = \begin{pmatrix} m_1 \\ m_2 \end{pmatrix} \\ \text{som}(\cdot, (2)) = s_2 = \begin{pmatrix} m_1 \\ M_2 \end{pmatrix} \\ \text{som}(\cdot, (3)) = s_3 = \begin{pmatrix} M_1 \\ m_2 \end{pmatrix} \\ \text{som}(\cdot, (4)) = s_4 = \begin{pmatrix} M_1 \\ M_2 \end{pmatrix} \end{array} \right.$$

**Indice de distance entre objets** On définit l'application  $\delta$  par

$$\delta : \Omega \times \Omega \longmapsto \mathfrak{R}^+$$

$$(\omega_i, \omega_l) \longmapsto \delta(\omega_i, \omega_l) = \|x_i - x_l\|,$$

où  $x_i, x_l \in \mathfrak{R}^d$  et  $\|\cdot\|$  est la norme euclidienne classique,  $\delta$  est une distance entre objets.

**Indice d'agrégation entre objets et groupe d'objets. Définition 9**

L'application  $D$  définie par

$$D : \Omega \times \mathfrak{R}_{V_d} \longmapsto \mathbb{R}$$



$$(\omega_i, R) \longmapsto D(\omega_i, R) = \frac{1}{2^d} \sum_{t=1}^{2^d} \|x_i - l_t\|$$

où  $l_t$  est le  $t$ ème sommet du pavé  $R$  qui représente la classe  $C$ ,  $l_t \in \mathbb{R}^d$  et  $\|\cdot\|$  est la norme euclidienne classique.  $D$  est un indice d'aggrégation entre objets et groupe d'objets et mesure la distance moyenne entre le pavé  $R$  et les vecteurs  $\{x_i \in \mathbb{R}^d; \omega_i \in C\}$ .

$$D(\omega_i, R) = \frac{1}{2^d} \sum_{t=1}^{2^d} \left( \sum_{j=1}^d (x_{ij} - (l_t)_j)^2 \right)^{\frac{1}{2}},$$

$x_{ij}$  et  $(l_t)_j$  sont respectivement la  $j$ ème coordonnée de  $x_i$  et  $l_t$

**Critère de classification** L'application  $C_r$

$$C_r : IP_k \times IL_k \longmapsto \mathbb{R}_+$$

$$(P, L) \longmapsto C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} \left[ \frac{1}{2^d} \sum_{t=1}^{2^d} \|x_i - l_k\| \right]$$

est le critère d'adéquation  $C_r$  s'écrit

$$C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} \left[ \frac{1}{2^d} \sum_{t=1}^{2^d} \left[ \left( \sum_{j=1}^d (x_{ij} - l_{tj}^k)^2 \right)^{1/2} \right] \right]$$

Ce critère exprime l'adéquation entre les individus et les noyaux des classes. Dans le cas où les erreurs sont cumulées au niveau des noyaux. Les observations des individus sont données sans erreurs.

La fonction d'affectation  $f$  et la fonction de représentation  $g$  sont données par

$$f : IL_k \longmapsto IP_k$$

$$L \longmapsto f(L) = P$$

$$P = \{P_1, \dots, P_k\} \in IP_k$$

et

$$P_t = \{\omega \in \Omega / D(\omega, l_t) \leq D(\omega, l_m) \text{ et } t \prec m \text{ s'il y a égalité}\}$$

**Proposition 23**

$f$  vérifie

$$\underset{P \in IP_k}{Min} C_r(P, L) = C_r(f(L), L)$$

$$g : IP_k \mapsto IL_k$$

$$P \mapsto g(P) = L^*$$

$$L^* = \{l_1, \dots, l_k\} \in IL_k$$

avec

$$l_t = \frac{1}{|P_t|} \left( \frac{1}{2^d} \right) \sum_{\{i\} \in P_k} (x_i),$$

**Proposition 24**

$g$  vérifie

$$\underset{L \in IL_K}{Min} [C_r(P, L)] = C_r(P, g(P) = L^*)$$

En effet, sans restreindre la généralité, étudions le cas où  $d = 2$

On a : Pour tout  $i = 1, \dots, n$

$$\omega_i \rightarrow R_i = \mathfrak{S}_i^1 \times \mathfrak{S}_i^2 \quad \text{où } \mathfrak{S}_i^1 = [\mu_{i1} - \frac{\sigma_{i1}}{2}; \mu_{i1} + \frac{\sigma_{i1}}{2}]; \mathfrak{S}_i^2 = [\mu_{i2} - \frac{\sigma_{i2}}{2}; \mu_{i2} + \frac{\sigma_{i2}}{2}]$$

Posant

$$\begin{cases} \mu_{i1} = x_i; \sigma_{i1} = a_i \\ \mu_{i2} = y_i; \sigma_{i2} = b_i \end{cases}$$

Les 4 sommets correspondants sont

$$S(i, 1) = S_{i1} = \begin{pmatrix} x_i - a_i \\ y_i - b_i \end{pmatrix}; S(i, 2) = S_{i2} = \begin{pmatrix} x_i - a_i \\ y_i + b_i \end{pmatrix}$$

$$S(i, 3) = S_{i3} = \begin{pmatrix} x_i + a_i \\ y_i + b_i \end{pmatrix}; S(i, 4) = S_{i4} = \begin{pmatrix} x_i + a_i \\ y_i - b_i \end{pmatrix}$$

$$\text{Pour } k = 1, K; L_k \in L \Rightarrow L_k \in \mathbb{R}^2 \Rightarrow L_k = \begin{pmatrix} x^{(k)} \\ y^{(k)} \end{pmatrix}$$

Le critère à optimiser s'écrit

$$Cr(P, L) = \sum_{k=1}^K \left[ \sum_{\omega_i \in P_k} \left[ \frac{1}{2^d} \sum_{t=1}^{2^d} \|S(i, t) - L_k\| \right] \right]$$

On suppose que les objets de la classe  $P_k$  sont numérotés par  $\{1, \dots, n_k\}$

$$\text{et } \begin{cases} x^{(k)} = x \\ y^{(k)} = y \end{cases}$$

On a

$$\|S(i, 1) - L_k\|^2 = ((x_i - a_i) - x)^2 + ((y_i - b_i) - y)^2$$

$$\|S(i, 2) - L_k\|^2 = ((x_i - a_i) - x)^2 + ((y_i + b_i) - y)^2$$

$$\|S(i, 3) - L_k\|^2 = ((x_i + a_i) - x)^2 + ((y_i + b_i) - y)^2$$

$$\|S(i, 4) - L_k\|^2 = ((x_i + a_i) - x)^2 + ((y_i - b_i) - y)^2$$

$$\begin{aligned} \varphi(x, y) &= 2 \sum_{i=1}^{n_l} [((x_i - a_i) - x)^2 + ((x_i + a_i) - x)^2] + \\ &\quad 2 \sum_{i=1}^{n_l} [((y_i - b_i) - y)^2 + ((y_i + b_i) - y)^2] \\ \Rightarrow \begin{cases} \frac{\partial \varphi(x, y)}{\partial x} = 0 \\ \frac{\partial \varphi(x, y)}{\partial y} = 0 \end{cases} &\Rightarrow x = \frac{1}{n_k} \sum_{i=1}^{n_l} (x_{i1}) \quad ; \quad y = \frac{1}{n_k} \sum_{i=1}^{n_l} (x_{i2}) \end{aligned}$$

#### Remarque 4

Les noyaux sont des vecteurs de  $\mathbb{R}^d$ .

#### 4.2.2 Alternative 2: Les erreurs sont cumulées au niveau de la description des individus et les noyaux sont donnés sans erreurs

Dans ce cas  $\omega_i$  est décrit par le pavé  $R_i = \prod_{j=1}^d \mathfrak{S}_j^i$  avec

$$\mathfrak{S}_i^l \subset \mathbb{R}$$

**Critère d'optimisation** Le critère  $C_r$  s'écrit

$$C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} D(R_i, l_k)$$

où

$$L = (l_1, \dots, l_k) \in IL_k, l_k \in \mathbb{R}^d$$

et  $D$  est définie par

$$D : \Omega \times IL \mapsto \mathbb{R}^*$$

$$(\omega_i, L) \mapsto D(\omega_i, L) = \frac{1}{2^d} \left( \sum_{t=1}^{2^d} \|S_i(t) - L\| \right)$$

$\|\cdot\|$  est la norme euclidienne,  $L, S_i(t) \in \mathbb{R}^d$ .  $S_i(t)$  est le sommet  $t$  du pavé  $R_i$ .  
Le critère s'écrit

$$C_r(P, L) = \sum_{k=1}^K \sum_{\omega_i \in P_k} \left[ \frac{1}{2^d} \sum_{l=1}^{2^d} \left( \sum_{j=1}^d (S_i(t, j) - (l_k)_j)^2 \right)^{\frac{1}{2}} \right]$$

La fonction représentative  $g$  est donnée par

$$g(P) = L = (l_1, \dots, l_k)$$

avec

$$(l_k)_j = \frac{1}{n_k} \sum_{i=1}^{n_k} (x_{ij})$$

$x_{ij}$  est la moyenne prise par l'individu  $\omega_i$  pour la variable  $j$ ,  $n_k = |P_k|$ .

**Proposition 25**

$g$  vérifie

$$\underset{L \in IL_K}{Min} C_r(P, L) = C_r(P, g(P))$$

La fonction d'affectation  $f$  est

$$f : IL_k \longmapsto IP_k$$

$$L = (l_1, \dots, l_k) \longmapsto f(L) = P = \{P_1, \dots, P_k\}$$

$$P_t = \{\omega \in \Omega / D(\omega, l_t) \leq D(\omega, l_m); \text{ et } t \prec m, \text{ si on a égalité}\}$$

**Proposition 26**

$f$  vérifie

$$\underset{P \in IP_k}{\text{Min}} C_r(P, L) = C_r(f(L), L)$$

On montre que les noyaux ou représentants des classes, dans ce cas, sont des vecteurs de  $\mathbb{R}^d$ . Il n'y a donc pas  $2^d$  sommets représentants chaque classe mais chaque noyau est représenté par un vecteur de  $\mathbb{R}^d$  et

$$D(\omega_i, l) = \|X_i - l\|; l \in \mathbb{R}^d$$

**4.2.3 Alternative 3: Les erreurs sont cumulées dans la description des individus et dans les noyaux des classes**

**Proposition 27**

L'application  $\Delta$  est définie par

$$\Delta : \mathfrak{R}_{V_d} \times \mathfrak{R}_{V_d} \longmapsto \mathbb{R}_+$$

$$(R_1, R_2) \longmapsto \Delta(R_1, R_2)$$

$$\begin{aligned} \Delta(R_1, R_2) &= \frac{1}{2^d} \sum_{l=1}^{2^d} D(S(1, l), R_2) = \frac{1}{2^d} \sum_{l=1}^{2^d} \sum_{t=1}^{2^d} \delta(S(1, l), S(2, t)) \\ &= \frac{1}{2^d} \sum_{l=1}^{2^d} \sum_{t=1}^{2^d} \left[ \left( \sum_{j=1}^d (S^j(1, l) - S^j(2, t))^2 \right)^{1/2} \right] \end{aligned}$$

$$S(1, l), S(2, t) \in \mathbb{R}^d,$$

$\delta$  est la distance euclidienne.  $\Delta$  est un indice de distance entre pavés de  $\mathbb{R}^d$ .

### Critère à optimiser

$$C_r(P, L) = \sum_{k=1}^K \sum_{\omega_i \in P_k} \left[ \frac{1}{2^d} \sum_{l=1}^{2^d} \left( \sum_{j=1}^d (S_i(t, j) - N_k(t, j))^2 \right)^{\frac{1}{2}} \right]$$

La fonction de représentation  $g$  est

$$g : IP_k \longmapsto IL_k$$

$$P = (P_1, \dots, P_k) \longmapsto g(P) = L = (l_1, \dots, l_k)$$

Si

$$\omega_i \rightarrow \left[ (x_{ij})_{j=1, d} ; (\varepsilon_{ij})_{j=1, d} \right],$$

on montre que le noyau  $l_k$  de la classe  $P_k$  est donné par

$$l_k \rightarrow \left[ (\mu_{kj})_{j=1, d} ; (\sigma_{kj})_{j=1, d} \right]$$

avec

$$\mu_{kj} = (\mu_k)_j, \quad \sigma_{kj} = (\sigma_k)_j$$

et

$$\left\{ \begin{array}{l} (\mu_k)_j = \frac{1}{n_k} \sum_{i=1}^{n_l} (x_i)_j \\ (\sigma_k)_j = \frac{1}{n_k} \sum_{i=1}^{n_l} (\varepsilon_{ij}) \end{array} \right.$$

En effet, sans restreindre la généralité, étudions le cas où  $d = 2$

On a

Pour tout  $i = 1, \dots, n$

$$\omega_i \rightarrow R_i = \mathfrak{S}_i^1 \times \mathfrak{S}_i^2 \quad \text{où } \mathfrak{S}_i^1 = \left[ \mu_{i1} - \frac{\sigma_{i1}}{2}; \mu_{i1} + \frac{\sigma_{i1}}{2} \right]; \mathfrak{S}_i^2 = \left[ \mu_{i2} - \frac{\sigma_{i2}}{2}; \mu_{i2} + \frac{\sigma_{i2}}{2} \right]$$

Pour  $k = 1, K; L_k \in L \Rightarrow L_k \in \mathfrak{R}_{V_d} \Rightarrow L_k = I_k^1 \times I_k^2$  avec

$$\left\{ I_k^j = \left[ x_j^{(k)} - a_s^{(k)}; x_j^{(k)} + a_s^{(k)} \right]; j = 1, 2 \right\}$$

Les 4 sommets associés à chacun des pavés dans l'ordre sont:

$$\begin{array}{ccc}
i \downarrow & & L_k \downarrow \\
S(i, 1) = S_{i1} = \begin{pmatrix} x_i - a_i \\ y_i - b_i \end{pmatrix} & ; & N(k, 1) = \begin{pmatrix} x - a \\ y - b \end{pmatrix} \\
S(i, 2) = S_{i2} = \begin{pmatrix} x_i - a_i \\ y_i + b_i \end{pmatrix} & ; & N(k, 2) = \begin{pmatrix} x - a \\ y + b \end{pmatrix} \\
S(i, 3) = S_{i3} = \begin{pmatrix} x_i + a_i \\ y_i + b_i \end{pmatrix} & ; & N(k, 3) = \begin{pmatrix} x + a \\ y + b \end{pmatrix} \\
S(i, 4) = S_{i4} = \begin{pmatrix} x_i + a_i \\ y_i - b_i \end{pmatrix} & ; & N(k, 4) = \begin{pmatrix} x + a \\ y - b \end{pmatrix}
\end{array}$$

Le critère à optimiser s'écrit

$$Cr(P, L) = \sum_{k=1}^K \left[ \sum_{\omega_i \in P_k} \left[ \frac{1}{2^d} \sum_{t=1}^{2^d} \|S(i, t) - N(k, t)\| \right] \right]$$

On a

$$\|S(i, 1) - N(k, 1)\|^2 = ((x_i - a_i) - (x - a))^2 + ((y_i - b_i) - (y - b))^2$$

$$\|S(i, 2) - N(k, 2)\|^2 = ((x_i - a_i) - (x - a))^2 + ((y_i + b_i) - (y + b))^2$$

$$\|S(i, 3) - N(k, 3)\|^2 = ((x_i + a_i) - (x + a))^2 + ((y_i + b_i) - (y + b))^2$$

$$\|S(i, 4) - N(k, 4)\|^2 = ((x_i + a_i) - (x + a))^2 + ((y_i - b_i) - (y - b))^2$$

$$\begin{aligned}
\varphi(a, b, x, y) &= 2 \sum_{i=1}^{n_l} [[(x_i - a_i) - (x - a)]^2 + [(x_i + a_i) - (x + a)]^2] + \\
&\quad 2 \sum_{i=1}^{n_l} [[(y_i - a_i) - (y - a)]^2 + [(y_i + a_i) - (y + a)]^2] \Rightarrow \\
&\quad \begin{cases} \frac{\partial \varphi(a, b, x, y)}{\partial x} = 0 \\ \frac{\partial \varphi(a, b, x, y)}{\partial y} = 0 \end{cases} \Rightarrow x = \frac{1}{n_k} \sum_{i=1}^{n_l} (x_i) \quad ; \quad y = \frac{1}{n_k} \sum_{i=1}^{n_l} (y_i) \\
&\quad \begin{cases} \frac{\partial \varphi(a, b, x, y)}{\partial a} = 0 \\ \frac{\partial \varphi(a, b, x, y)}{\partial b} = 0 \end{cases} \Rightarrow a = \frac{1}{n_k} \sum_{i=1}^{n_l} (a_i) \quad ; \quad b = \frac{1}{n_k} \sum_{i=1}^{n_l} (b_i)
\end{aligned}$$

La fonction d'affectation  $f$  est

$$f : IL_k \mapsto IP_k$$

$$L \mapsto f(L) = P$$

$$P = \{P_1, \dots, P_k\} \in IP_k$$

est

$$P_t = \{\omega \in \Omega / D(\omega, l_t) \leq D(\omega, l_m); \text{ et } t \prec m \text{ si on a égalité} \}$$

### 4.3 Objets décrits par des ellipsoïdes de $\mathbb{R}^d$

#### 4.3.1 Critère et problème d'optimisation

Nous désirons regrouper les  $n$  individus en  $k$  classes où chaque individu sera décrite par son élément représentatif. L'homogénéité des classes est mesurée par

$$H_k = \sum_{\{i\} \in P_k} d_{ma}(W_i, L_k),$$

où  $L_k$  est l'élément représentatif du noyau de la classe  $C_k$ ,  $W_i \in \mathbb{N} \times \mathbb{R}^d \times M_d(\mathbb{R})$  et  $d_{ma}$  est la distance de Mahannalobis ponderée.  $H_k$  mesure la dispersion des individus à l'intérieur de la classe.  $H_k$  est aussi l'inertie des objets de la classe autour de son noyau.

Soit  $P_K$  l'ensemble des partitions à  $k$  classes

$$P \in P_K \Rightarrow P = \{P_1, \dots, P_K\}$$

est une partition à  $K$  classes,  $L_K$  l'ensemble des  $K$  représentations

$$(L \in L_K) \Rightarrow L = \{L_1, \dots, L_k\}$$

où

$$L_t \in \mathbb{N} \times \mathbb{R}^d \times M_d(\mathbb{R}).$$

On défini l'application

$$C_r : P_K \times L_K \rightarrow \mathbb{R}_+$$

$$(P, L) \rightarrow C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} d_{ma}(W_i, L_k)$$



$$d_{ma}(W_i, L_k) = \frac{N_i \|\mu_i - \tilde{\mu}_k\|_{\Sigma_i^{-1}} + \tilde{N}_k \|\mu_i - \mu_l\|_{\tilde{\Sigma}_k^{-1}}}{N_i + \tilde{N}_k}$$

$$W_i; c_k \in \mathbb{N} \times \mathbb{R}^d \times M_d(\mathbb{R})$$

Le critère s'écrit

$$C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} d_{ma}(W_i, L_k)$$

Ce critère exprime l'adéquation entre les éléments représentatifs des individus dans les classes avec les éléments représentatifs des noyaux des classes.  $C_r$  s'écrit

$$C_r(P, L) = \sum_{k=1}^K \sum_{\{i\} \in P_k} d_{ma}(W_i, L_k)$$

On recherche le couple  $(P^*, L^*)$  qui réalise

$$\underset{\substack{P \in P_K \\ L \in L_K}}{\text{Min}} C_r(P, L)$$

Les algorithmes utilisés pour résoudre ce type de problème sont du type  $k - means$ . Ces algorithmes sont basés sur la définition de la fonction de représentation  $g$  et de la fonction d'affectation  $f$  qui seront utilisées simultanément pour faire décroître le critère.

### 4.3.2 Caractérisation d'une classe d'objets(résultat très important)

Il s'agit de rechercher le noyau représentant chacune des classes générée par cet algorithme

Si l'on note par  $i = 1, \dots, k$  les  $k$  objets dans la classe  $C$ .

soit  $\varphi$  l'application définie par

$$\mathbb{R}^d \times M_{(d)} \times \mathbb{N} \xrightarrow{\varphi} \mathbb{R}_+$$

$$(\mu, S, N) \longrightarrow \varphi(\mu, S, N) = \sum_{i=1}^K \left[ \frac{N_i}{N_i + N} \|\mu_i - x\|_{S_i^{-1}}^2 + \frac{N}{N_i + N} \|\mu_i - x\|_{S^{-1}}^2 \right]$$

Il s'agit de chercher  $\mu^*, S^*, N^*$  qui réalise le minimum de  $\varphi$ . Résoudre le problème d'optimisation

$$\begin{aligned} \text{Min } & \varphi(\mu, S, N) \\ (x, S, L) & \in \mathbb{R}^d \times M_{(d)} \times \mathbb{N} \end{aligned}$$

Sans restreindre la généralité, on s'intéresse au cas où  $p = N_i \quad \forall i$ . On cherche à résoudre le problème d'optimisation suivant

$$\begin{aligned} \text{Min } & \varphi(x, S) \\ (x, S, p) & \in \mathbb{R}^d \times M_{(d)} \end{aligned}$$

on pose  $\mu = x : \mu_i = x_i$

Condition nécessaire

$$\begin{cases} \frac{\partial \varphi}{\partial x} = 0 & (1) \\ \frac{\partial \varphi}{\partial S} = 0 & (2) \end{cases}$$

$$\frac{\partial \varphi}{\partial x} = 0 (1) \iff \sum_{i=1}^K [(x_i - x)' (S_i^{-1} + S^{-1})] = 0 (1) \iff \sum_{i=1}^K [(x_i - x)' S_i^{-1}] + \sum_{i=1}^K [(x_i - x)' S^{-1}] = 0$$

En ce qui concerne la seconde condition, on remarque que la première expression de  $\varphi$  ne dépend pas de  $S$ . Soit

$$\Phi(S) = \sum_{i=1}^K \|(x_i - x)\|_{S^{-1}}^2$$

On a

$$\Phi(S + H) - \Phi(S) = \sum_{i=1}^K [(x_i - x)' (S + H)^{-1} (x_i - x)] - \sum_{i=1}^K [(x_i - x)' (S)^{-1} (x_i - x)]$$

or

$$(S + H)^{-1} = [S(I + S^{-1}H)]^{-1} = (I + S^{-1}H)^{-1} S^{-1}$$

Le développement de  $(I + S^{-1}H)^{-1}$  donne

$$(I + S^{-1}H)^{-1} = I - S^{-1}H + \dots \Rightarrow (S + H)^{-1} = [I - S^{-1}H + \dots] S^{-1} \Rightarrow$$

$$\begin{aligned}
\Phi(S+H) - \Phi(S) &= \sum_{i=1}^K [(x_i - x)' (S^{-1}.H.S^{-1}) (x_i - x)] + \dots \\
&= \sum_{i=1}^K \left[ [S^{-1}(x_i - x)]' (H) [S(x_i - x)] \right] + \dots \\
&= \sum_{i=1}^K \|S^{-1}(x_i - x)\|_H^2 + \dots \\
\Rightarrow \frac{\partial \varphi}{\partial S} = 0 \quad (2) &\Rightarrow \sum_{i=1}^K \|S^{-1}(x_i - x)\|_H^2 = 0 \Rightarrow S^{-1}(x_i - x) = 0 \forall i \Rightarrow S^{-1} = 0
\end{aligned}$$

car  $x_i \neq x$

$$(1) \iff \sum_{i=1}^K [(x_i - x)' (S_i^{-1})] = 0 \Rightarrow (x)' \sum_{i=1}^K (S_i^{-1}) = \left( \sum_{i=1}^K (x_i)' (S_i^{-1}) \right) \Rightarrow$$

Si

$$\begin{aligned}
\Gamma_k &= \sum_{i=1}^K (S_i^{-1}) \\
\Rightarrow x &= (\Gamma_k)^{-1} \sum_{i=1}^K (S_i^{-1}) (x_i) = \tilde{\mu}_k
\end{aligned}$$

Ce qui montre qu'une fois les clusters obtenues, la caractérisation n'est pas entachées d'erreurs et est donnée de manière exacte Ce qui parait tout à fait dans l'ordre des choses.

La distance entre un individu  $\omega_i$  donné décrit par le couple  $(\nu_i, \Sigma_i)$  et la classe  $C_k$  caractérisée par son noyau  $L_t$  est donnée par

$$d_{ma}(\omega_i, L_k) = \|\mu_i - \tilde{\mu}_k\|_{\Sigma_i^{-1}}$$

### 4.3.3 La fonction de représentation

$$g : P_k \rightarrow L_k$$

$$P = \{P_1, \dots, P_k\} \rightarrow g(P) = \{L_1, \dots, L_k\}$$

$g$  doit vérifier

$$\underset{L \in L_k}{MinCr}(P, L) = Cr(P, g(P)) \quad (*)$$

**Proposition 28**

La fonction  $g$  définie par

$$P = \{P_1, \dots, P_k\} \rightarrow g(P) = \{L_1, \dots, L_k\}$$

$$\text{où } L_k = (\Gamma_k)^{-1} \sum_{i=1}^K (S_i^{-1})(x_i) \quad \text{avec } \Gamma_k = \sum_{i=1}^K (S_i^{-1})$$

verifie(\*)

#### 4.3.4 La fonction d'affectation

$$f : L_k \rightarrow P_k$$

$$L = \{L_1, \dots, L_k\} \rightarrow f(L) = \{P_1, \dots, P_k\}$$

$f$  doit vérifier

$$\underset{P \in P_k}{MinCr}(P, L) = Cr(f(L), L)$$

**Proposition 29**

Ce minimum est obtenu pour

$$P_l = \{\omega_i \in \Omega \quad \text{tels que } d_{ma}(\omega_i, L_l) \leq d_{ma}(\omega_i, L_t); \forall t \neq l \text{ et } l \prec t \text{ en cas d'égalité}\}$$

#### 4.3.5 L'algorithme

On choisi au hasard ou par d'autres artifices  $L^{(0)}$  et on utilise alternativement les fonctions  $f$  et  $g$  pour faire décroître le critère L'algorithme de déroule comme suit

$$L^{(0)} \xrightarrow{f} P^{(1)} \xrightarrow{g} L^{(1)} \xrightarrow{f} P^{(2)} \dots L^{(*)} \xrightarrow{f} P^{(*)}$$

L'algorithme s'arrête dès que la partition obtenue ne change plus. Nous construisons ainsi 2 suites  $U_n$  et  $V_n$  telles que

$$\begin{cases} U_n = Cr(P^{(n)}, L^{(n)}) \\ V_n = Cr(P^{(n)}, L^{(n)}) \end{cases}$$

**Proposition 30**

La suite  $U_n$  converge en décroissant et la suite  $V_n$  est stationnaire à partir d'un certain rang donc convergente.

### 4.3.6 Exemples

On désire classifier les 6 objets en 2 classes. Chaque objet est décrit par 3 variables quantitatives  $V_1, V_2$  et  $V_3$ . Les données en entrée se présentent comme suit

$$\omega_1 \mapsto T_1 = \begin{bmatrix} 3 & 5 & 6 \\ 7 & 9 & 10 \\ -5 & 2 & 0 \\ 3 & 1 & -1 \\ 5 & 4 & 2 \end{bmatrix}; T_2 = \begin{bmatrix} 2 & -1 & 0 \\ 3 & 5 & 11 \\ -2 & 3 & 4 \\ 6 & 7 & 8 \end{bmatrix}; T_3 = \begin{bmatrix} 4.1 & 2.6 & -1 \\ 2.4 & 0.5 & 4 \\ 5 & 6 & 7 \end{bmatrix}$$

$$T_4 = \begin{bmatrix} 3 & 3.6 & 4 \\ -1 & 17 & 0 \\ 10 & 3 & 12 \\ 13 & 0 & 11 \end{bmatrix}; T_5 = \begin{bmatrix} 4 & -1 & 2 \\ 3 & 2 & 6 \\ 2 & 3 & 7 \\ 3 & 5 & -11 \\ 1 & 0.2 & 3 \\ 0.4 & 7 & 2 \end{bmatrix}; T_6 = \begin{bmatrix} 2 & 7 & 2.5 \\ -1 & 3 & 1.5 \\ 2 & 5 & 0 \\ 4 & 0.5 & 14 \\ 6 & 0.7 & 3 \end{bmatrix}$$

- On remarque que les tables  $T_1, T_2, \dots, T_6$  n'ont pas la même dimension
- Les moyennes, les écart-types et les intervalles de confiance pour chaque variable et pour chaque individu sont donnés dans les tableaux suivants

$$\omega_1 \rightarrow \begin{bmatrix} \text{mean}(\mu) & \text{std.deviation}(\sigma) \\ 2.6 & 4.56070 \\ 4.2 & 3.11448 \\ 3.4 & 4.56070 \end{bmatrix}; \omega_1 \rightarrow \begin{bmatrix} \mu - \frac{\sigma}{2} & \mu + \frac{\sigma}{2} \\ 0.31965 & 4.88034 \\ 2.64275 & 5.75724 \\ 1.11965 & 5.68035 \end{bmatrix}$$

$$\omega_2 \rightarrow \begin{bmatrix} \text{mean}(\mu) & \text{std.deviation}(\sigma) \\ 2.25 & 10.90489 \\ 3.50 & 10.08298 \\ 5.75 & 4.78714 \end{bmatrix}; \omega_2 \rightarrow \begin{bmatrix} \mu - \frac{\sigma}{2} & \mu + \frac{\sigma}{2} \\ -3.20245 & 7.70244 \\ -1.54149 & 8.54149 \\ 3.35643 & 8.14357 \end{bmatrix}$$

$$\omega_3 \rightarrow \begin{bmatrix} \text{mean}(\mu) & \text{std.deviation}(\sigma) \\ 3.8333 & 1.32035 \\ 3.0333 & 2.77549 \\ 3.3333 & 4.04145 \end{bmatrix}; \omega_3 \rightarrow \begin{bmatrix} \mu - \frac{\sigma}{2} & \mu + \frac{\sigma}{2} \\ 0.317283 & 4.49318 \\ 1.64526 & 4.42075 \\ 1.31228 & 5.35373 \end{bmatrix}$$

$$\omega_4 \rightarrow \begin{bmatrix} \text{mean}(\mu) & \text{std.deviation}(\sigma) \\ 6.25 & 6.39661 \\ 5.9 & 7.56571 \\ 6.75 & 5.73730 \end{bmatrix}; \omega_4 \mapsto \begin{bmatrix} \mu - \frac{\sigma}{2} & \mu + \frac{\sigma}{2} \\ 3.05170 & 9.44831 \\ 2.11715 & 9.68286 \\ 3.88135 & 9.61865 \end{bmatrix}$$

$$\omega_5 \mapsto \begin{bmatrix} \text{mean}(\mu) & \text{std.deviation}(\sigma) \\ 2.2333 & 1.35892 \\ 2.7 & 2.97658 \\ 1.5 & 6.47302 \end{bmatrix}; \omega_5 \mapsto \begin{bmatrix} \mu - \frac{\sigma}{2} & \mu + \frac{\sigma}{2} \\ 1.55384 & 2.91276 \\ 1.21171 & 4.18829 \\ 1.73651 & 4.73651 \end{bmatrix}$$

$$\omega_6 \mapsto \begin{bmatrix} \text{mean}(\mu) & \text{std.deviation}(\sigma) \\ 2.6 & 2.60768 \\ 3.24 & 2.79517 \\ 4.2 & 5.59687 \end{bmatrix}; \omega_6 \mapsto \begin{bmatrix} \mu - \frac{\sigma}{2} & \mu + \frac{\sigma}{2} \\ 1.29616 & 3.90384 \\ 1.84242 & 4.63759 \\ 1.40157 & 6.99844 \end{bmatrix}$$

- Les sommets associés pour chaque individu s'écrivent

Pour  $i = 1, n; S(i, l) \in \mathbb{R}^d$

$$[S(i, l)]_j = \mu_i^j + \varepsilon \cdot \sigma$$

; avec  $\varepsilon = \pm 1$

$$\begin{array}{cccccccc} \overbrace{S(i, 1)}^{\mu_i^1 - \sigma_i^1} & \overbrace{S(i, 2)}^{\mu_i^1 + \sigma_i^1} & \overbrace{S(i, 3)}^{\mu_i^1 - \sigma_i^1} & \overbrace{S(i, 4)}^{\mu_i^1 + \sigma_i^1} & \overbrace{S(i, 5)}^{\mu_i^1 - \sigma_i^1} & \overbrace{S(i, 6)}^{\mu_i^1 + \sigma_i^1} & \overbrace{S(i, 7)}^{\mu_i^1 - \sigma_i^1} & \overbrace{S(i, 8)}^{\mu_i^1 + \sigma_i^1} \\ \mu_i^2 - \sigma_i^2 & \mu_i^2 - \sigma_i^2 & \mu_i^2 + \sigma_i^2 & \mu_i^2 + \sigma_i^2 & \mu_i^2 - \sigma_i^2 & \mu_i^2 - \sigma_i^2 & \mu_i^2 + \sigma_i^2 & \mu_i^2 + \sigma_i^2 \\ \mu_i^3 - \sigma_i^3 & \mu_i^3 - \sigma_i^3 & \mu_i^3 - \sigma_i^3 & \mu_i^3 - \sigma_i^3 & \mu_i^3 + \sigma_i^3 & \mu_i^3 + \sigma_i^3 & \mu_i^3 + \sigma_i^3 & \mu_i^3 + \sigma_i^3 \end{array}$$

- A ce niveau, 2 formes de la description des objets peuvent être considérées

-Chaque objet est décrit par un vecteur de  $\mathbb{R}^d$  qui contient les valeurs moyennes prises par les variables qui les décrivent .

$$\omega_i \mapsto (\mu_i)' = [\mu_i^1, \dots, \mu_i^d]$$

où

$$\mu_i^j = \frac{1}{L_i} \sum_{l=1}^{L_i} x_{il}^j.$$

Les données en entrée sont regroupées dans le tableau  $X$ .

$$X = \begin{bmatrix} \mu_i^1 & \mu_i^2 & \mu_i^3 \\ 2.6 & 4.2 & 3.4 \\ 2.25 & 3.5 & 5.75 \\ 3.83 & 3.03 & 3.33 \\ 6.25 & 5.9 & 6.75 \\ 2,23 & 2.7 & 1.5 \\ 2.6 & 3.24 & 4.2 \end{bmatrix}$$

L'algorithme des  $k$  - means déroulé pour classifier les 6 individus en 2 classes donne

$$\begin{array}{c} \text{Centres finaux des classes} \\ \begin{bmatrix} \text{class1} & \text{class2} \\ v11 & 2.7 & 6.25 \\ v22 & 3.3 & 5.90 \\ v33 & 3,64 & 6.75 \end{bmatrix} \\ \begin{bmatrix} \text{case number} & \text{cluster} & \text{distance} \\ \omega_1 & 1 & ,90305 \\ \omega_2 & 1 & 2,16773 \\ \omega_3 & 1 & 1,20818 \\ \omega_4 & 2 & ,00000 \\ \omega_5 & 1 & 2,27794 \\ \omega_6 & 1 & ,58051 \end{bmatrix} \end{array}$$

- Chaque objet est décrit par les  $2^d$  ( $d = 3$ ) sommets associés.

Les coordonnées des sommets associés aux 6 individus sont

Pour  $\omega_1 \downarrow$

,31965	2,64276	1,11965	4,88035	2,64276	1,11965	,31965
5,75724	1,11965	4,88035	5,75724	1,11965	4,88035	2,64276
5,68035	,31965	5,75724	5,68035	,31965	2,64276	5,68035
	4,88035	5,75724	5,68035	<b>1</b>	<b>4,20132</b>	

Pour  $\omega_2 \downarrow$

-3,20245	-1,54149	3,35643	7,70245	-1,54149	3,35643	-3,20245
8,54149	3,35643	7,70245	8,54149	3,35643	7,70245	-1,54149
8,14357	-3,20245	7,70245	8,14357	-3,20245	-1,54149	8,14357
	7,70245	5,20783	8,14357	<b>2</b>	<b>8,03977</b>	

Pour  $\omega_3 \downarrow$

3,17283	1,64526	1,31228	4,49318	1,64526	1,31228	3,17283
4,42075	1,31228	4,49318	4,42075	1,31228	4,49318	1,64526
5,35373	3,17283	4,42075	5,35373	3,17283	1,64526	5,35373
	4,49318	4,42075	5,35373	<b>1</b>	<b>3,59757</b>	

Pour  $\omega_4 \downarrow$

3,05170	2,11715	3,88135	9,44831	2,11715	3,88135	3,05170
9,68286	3,88135	9,44831	9,68286	3,88135	9,44831	2,11715
9,61865	3,05170	9,68286	9,61865	3,05170	2,11715	9,61865
	9,44831	9,68286	9,61865	<b>2</b>	<b>8,03977</b>	

Pour  $\omega_5 \downarrow$

1,55384	1,21171	1,73651	2,91276	1,21171	1,73651	1,55384
4,18829	1,73651	2,91276	4,18829	1,73651	2,91276	1,21171
4,73651	1,55384	4,18829	4,73651	1,55384	1,21171	4,73651
	2,91276	4,18829	4,73651	<b>1</b>	<b>3,32424</b>	

Pour  $\omega_6 \downarrow$

1,29616	1,84242	1,40157	3,90384	1,84242	1,40157	1,29616	4,63759	1,40157
1,29616	4,63759	1,40157	3,90384	1,84242	6,99844	1,29616	4,63759	6,99844
1,29616								
	1,84242	6,99844	3,90384	4,63759	6,99844	<b>1</b>	<b>3,41777</b>	

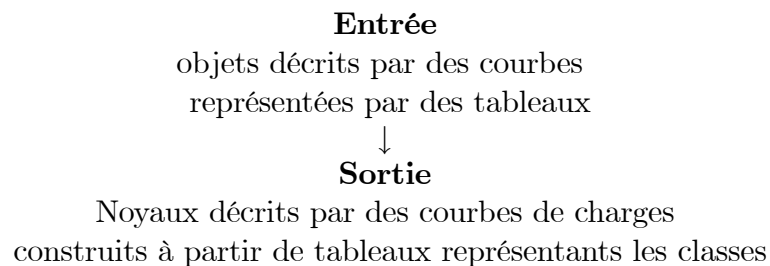
## Chapitre VI

# Conclusion générale

- Dans cette thèse, on s'est intéressé à l'analyse des éléments constitutifs d'une structure de juxtaposition de tableaux multiples qui dans certains cas des objets matriciels.
- Nous avons développé différentes approches adaptées aux types de structure que nous avons liées à des formalismes spécifiques qui permettent d'appréhender les éléments de la structure des données.
- Nous avons montré comment et dans quelles conditions les techniques factorielles peuvent y être adaptées quand il s'agit de décrire les éléments constitutifs de la structure. Nous avons étudié le cas d'une structure de juxtaposition de tableaux de mesure et le cas où chaque élément de la structure est décrit par un tableau de données quantitatives.

La démarche consiste à résumer chaque élément de la structure par sa représentation sur des sous-espaces de faible dimension. Ces sous-espaces sont obtenus à partir des déroulements d'ACP.

- Nous avons discuté les avantages et les inconvénients. Cette approche qui a fait l'objet des travaux de Bouroche[10] pose un problème fondamental qui est: *# La recherche d'un sous espace commun ou compromis à tous les sous espaces construits#*. Il propose que ces choix se fassent par des méthodes heuristiques suivant des critères de minimisation d'inertie mais le problème de fixer la dimension des sous-espaces n'a pas été abordé. Ce sous-espace affine ou vectoriel doit permettre d'avoir une vision globale des éléments constitutifs pour une étude plus fine.
- Ce problème intéresse d'autres disciplines de la mathématique et en particulier la géométrie différentielle, nous citons les travaux de P.Orlick et H Terao[74][1992]. Ce problème reste à mon sens assez ouvert.
- Nous avons contourné ce problème en proposant les 2 approches qui constituent l'ossature de cette thèse et qui conduisent à un niveau de connaissance en sortie équivalent à celui introduit en entrée (condition nécessaire d'une bonne analyse, Diday [1998]).
- Cette remarque se justifie aisément dans l'étude réalisée dans le chapitre 3 comme le montre le tableau suivant





- De plus, dans les 2 approches, nous tenons compte de la variabilité des observations de chaque individu pour les différentes variables qui les décrivent.
- La seconde approche nous paraît plus générale tant au point de vue méthodologique qu'au point de vue du formalisme utilisé. Nous avons manipulé des variables aléatoires qui sont des concepts sous tendus par des notions de variabilité et d'extension. L'approche peut s'adapter facilement à des données mesurables ou qualitatives car elle s'appuie essentiellement sur des calculs de fréquences relatives ou des probabilités d'apparition de telle ou telle modalité. Elle permet de s'étendre à tous les types de structure de données
- . Le fait de prendre les données sans manipulation préalable et mise en forme rend les résultats obtenus plus proches de la réalité que l'on désire appréhender donc conduit à des résultats plus crédibles et facilement interprétables.

Cette approche est basée sur la définition de Shannon Weiner du concept de système physique. La notion de système physique a longtemps fait l'objet de discussions entre scientifiques. Depuis que Carnot a ébauché les concepts fondamentaux de la thermodynamique et défini la caractéristique physique qui mesure le degré du désordre et qu'il a baptisé entropie. Ces discussions continuent avec l'apparition de la notion d'objets probabilistes [32](Diday,...) et de sous ensembles flous (Kauffman,Bezdek,...)

- Nous pensons que le choix de la définition de Shannon du système physique a permis de manière naturelle d'utiliser et d'expliquer l'aspect pratique de la distance de Kullback-Leibler comme un indice de distance entre les éléments constitutifs de structures complexes.

Ce travail ouvre la voie à plusieurs domaines d'applications pratiques et particulièrement en médecine.

Enfin, ce travail ouvre le chemin à plusieurs axes de recherche et d'investigations et les extensions sont nombreuses. Parmi les éclairages possibles, on peut citer

- Analyse d'objets décrits par des distributions paramétrées connues (normales, multinomiales,...)
- Etudes algébriques et propriétés d'objets aléatoires(complétude,...)

# Chapitre VII

## Bibliographie

### Bibliographie

- [1] T. Adamson (1996) , *Data structure and algorithms, a first course*. Springer Verlag.
- [2] A. Anderberg (1973) , *Cluster analysis for applications* . Academic Press  
L'un des premiers ouvrages spécialisé en langue anglaise qui contient l'historique des développements qu'a connu la classification automatique.
- [3] T.W. Anderson (1974) , *An introduction to multivariate analysis* John Wiley and sons, New York.  
Le premier livre qui préconise une modélisation des données en entrée avant d'entamer l'analyse. C'est un livre de référence qui fait le lien entre l'approche classification de l'école "française " et l'approche anglo-saxone dite "clustering
- [4] Arabie and Hubert (1992) , *Combinatorial data analysis*. Annual review of psychology. p: 168 – 203.
- [5] G.H. Ball,D.J Hall (1967) , *ISODATA,a novel method of data analysis and pattern classification*, Stanford research Inst. California,
- [6] J.D.Banfield and B.Basford (1989) , *Model based gaussian and non gaussian clustering*. Technical report n° **186** Dept. of Statistical, GN 22, University of Washington.
- [7] JP.Benzecri (1973) , *L'analyse des données* tome 1 & 2 .Dunod, Paris.  
Ouvrage spécialisé en langue française contenant d'importantes considérations historiques du sujet et de rigoureux développements formels sur la notion de classification automatique.
- [8] J.C. Bezdek (1975) , *Mathematical models for systematic and taxonomy* .Freidman,San Francisco,USA.
- [9] HH. Bock (1975) , *On significiance tests in cluster analysis*. Journal of Classification, p: 77–108, Springer Verlag, N.Y.
- [10] JM. Bouroche (1975) , *Analyse de données ternaires; Double analyse en composantes principales(DACP)*. Thèse de l'Université de Paris VI
- [11] X.Bry (1997) , *Analyses factorielles multiples*. Economica, Paris.
- [12] C.Burt (1950) , *The factorial analysis of qualitative data*. British J of statistical psychology **3, 3**, p: 166–185.  
Burt a été incontestablement l'inventeur du point de vue méthodologique de l'analyse factorielle multiple. Dans cet article ,il préconise le calcul du tableau qui

regroupe les croisements deux à deux de toutes les variables associées qui interviennent. Ce tableau qui porte son nom est une juxtaposition de tables de contingence. Après une normalisation de celui-ci, il devient un tableau disjonctif complet sur lequel opère l'analyse factorielle multiple (AFM)

- [13] T.Calcom (1981), *Analyse conjointe de matrices de données*.Thèse de l'Université de Grenoble, France
- [14] A.Carlier (1985), *Application de l'analyse des évolutions et de l'analyse intra-période*. Revue de Statistique et Analyse des Données **10-1**, p: 27–53.
- [15] A. Carlier, C. Lavit, M. Pages, M.O. Pernin, J.C.Turlot, (1988), *A comparative review of methods which handles a set of indexed data tables*, Multiway data analysis, p 85 -102
- [16] J.D.Carroll and J.J.Chang (1970), *Analysis of individual differences in multidimensional scaling via  $N$  way generalisation of "Eckart-Young" decomposition*. Psychometrika, **49- 3**, p: 403–423.
- [17] P.Cazes (1980), *Analyse de certains tableaux rectangulaires décomposés en blocs*. Les Cahiers de l'Analyse des Données , **5**, p: 387– 403.
- [18] P.Cazes (1986), *Une généralisation des correspondances multiples et des correspondances hiérarchiques*. Cahiers du BURO, **46-47**, Université de Pierre et Marie Curie, Paris , France
- [19] P.Cazes, A. Chouakria, E.Diday, Y Scheketman (1997), *Analyse en composantes principales d' objets décrits par des intervalles*. Revue de Statistiques appliquées. Vol **45**, n°3.
- [20] G.Celeux (1988), *Classification et modèles: RSA XXXVI*, **4**, p: 43–58.
- [21] J.L. Chandon and S. Pison (1981), *Analyse typologique: théorie et applications*. Masson, Paris.
- [22] L.Csizard (1967), *Information type distance measures and indirect observation*. Stud. Scien Math., Hungar, Vol. **2**, p :299–318.
- [23] R.M. Cormack and C Hubert (2000), *Hierarchical dependency models for multivariate survival data with censoring* . Life Data, **6**, pp: 299–320.
- [24] R.M. Cormack (1971), *A review of classification*.J.of Royal statist. Society, serie A, **134**, part.3, pp: 321– 367.  
L'un des premiers articles faisant état de la synthèse des techniques de classification.
- [25] A.D. Cordon (1987), *A review of hierarchical classification*. J.of Royal Statist. Society, serie A, **150**, part.2 , p: 119–137.  
Complément de l'article de Cormack par l'ajout des travaux les plus récents sur la classification hiérarchique

- [26] L.Cuttman (1941), *The quantification of a class of attributes*, p:319-348, Soc. Sc. Res. Council, NY
- [27] N.Day (1969), *Estimating the components of mixture of normal distribution*. *Biometrika*, **56**, p: 463-474.
- [28] F.Dazy and J.F Le Barzik (1997), *L'analyse des données évolutives : méthodes et applications*. Edition technip, Paris.
- C'est un ouvrage consacré aux techniques factorielles adaptées à l'analyse d'une juxtaposition de tableaux indexés par le temps. Les applications proposées sont très intéressantes et touchent le domaine de la vie sociale française '(étude de l'évolution de la délinquance et de l'évolution des votes en France). Les méthodes statis et dacp sont largement développées
- [29] F.A.T. De Carvalho (1994), *Proximity coefficients between boolean objects. New approaches in classication and data analysis*, Diday (Eds). Springer Verlag, p: 387-394.
- [30] E.Diday (1971), *La méthode des nuées dynamiques*. *RSA* **19**, N° 19, p: 19-34.
- La méthode des nuées dynamiques généralise la technique proposée par Ball & Hall[65] et permet de modifier le choix initial des centres de classes. Le problème de classification se pose en terme d'optimisation d'un critère dont l'algorithme des nuées dynamiques permet de donner une solution locale dépendante du choix initial. Diday a mis au point une procédure de groupements stables (formes fortes) qui permet de remédier au moins partiellement au principal inconvénient des algorithmes proposés à savoir la dépendance de la solution obtenue en fonction du choix initial.
- [31] E.Diday and all (1980), *Optimisation en classification automatique* tome 1 & 2. INRIA, Paris.
- Rassemble plusieurs travaux sur l'application de l'algorithme des nuées dynamiques et des développements de la classification par les techniques des centres mobiles. Les travaux ont fait l'objet de thèses et diverses applications sont exposées.
- [32] E Diday (1995), *Probabilist, possibilist and belief objects for knowledje analysis* *Annals of operationnel research*, **55**, p: 227-276.
- [33] P.Ducemetière (1970), *Les méthodes numériques de classification automatique*. *RSA*, Vol **48** n° 4.
- [34] R.O.Duda and P.E Hart (1973), *Pattern classification and science analysis*. J.Wiley, N.Y.
- [35] C.Eckart, G Young (1936) : The approximation of one matrix by another of lower rank, *Psychometrika*, **1**, p 211- 218
- [36] B Escofier and J. Pages (1988), *Analyses factorielles simples et multiples : objectifs, méthodes et interprétation*. Dunod, Paris.

- [37] B.Escofier (1989) , *Quelques indices pour comparer des tableaux de contingence* Re-  
vue de Statistique et d'Analyse des Données, **1**, p: 39–51.
- [38] B.Everit and D Hand (1981) , *Finite mixture distributions*. Chapman and Hall,  
U.S.A.
- [39] RA. Fisher (1936) ,*The use of multiple measurements in taxonomic problems*. Ann  
Of eugenics, **7**, p: 179 –188.
- Dans cet article Fisher a montré que dans le cas où les individus ne seraient pas  
décrits par un seul paramètre, le calcul d'une partition exacte est possible car il existe  
une relation d'ordre entre les individus ce qui limite considérablement l'éventail des  
partitions
- [40] E.W Forgy,(1965) ,Cluster analysis of multivariate data,Biometrics, **21**, 3, p 728
- [41] T.Foucart (1983) , *Une nouvelle approche de la méthode STATIS* . RSA, XX1, **2**, p:  
61–75.
- [42] T.Foucart (1984) , *Analyse factorielle de tableaux multiples*.Masson,France
- [43] C.Gowda & E.Diday (1991) , *Symbolic clustering using a new dissimilarity mea-  
sure*.Pattern Recognition Vol **24**, N° 6
- [44] M.Jambu and M.O.Lebeaux (1978) , *Classification automatique pour l'analyse des  
données* Dunod,Paris
- [45] R. A. Harshman,(1970) ,*Foundation of the parafac procedure*, Phonetics,**16**, UCLA,  
Los Angeles
- [46] J A. Hartigan (1975) .*Clustering algoritms* J.Wiley, New York.
- [47] L. Kauffman and P.J.Rousseeuw (1990) , *Finding groups in data*. J.Wiley, New York
- [48] C.Hayashi (1956) ,Theory and examples of quantification II, *procedure .of the in-  
stitutue StaT .Math*, **4**(2),P19-30
- [49] P. Horst (1965) ,*Factor analysis of data matrices*, Holt. Rinehart, Winston, NY
- [50] R. J. Ketteling (1971) , *Canonical analysis of several sets of variables*, Biometrika,  
**58**(3), p 433-450
- [51] WP Krijner (1993) , *The analysis of three-way arrays constrained parafac methods*.  
DSWO Press.
- [52] R. Kodell,and J.Chen (2001) , *Infering effects on tumor frequencies and times to  
observation in the analysis of tumor multiplicity data*. Biometrical Journal, **43**, **4**,  
p:447–460.
- [53] P. R Kroomenberg (1983) ,*Three mode principal component analysis*. DSWO Press.

- [54] S. Kullback (1959), *Information theory and statistics*. J.Wiley, New York
- [55] S. Kullback and R.A. Leibler (1967), *Information and sufficiency*. Ann. Math Statist., **32**, p: 79–86
- [56] C. Lavit (1988), *Analyse conjointe de tableaux quantitatifs*. Masson, France.
- [57] L.Lebart et N.Tabard (1973), *Recherches sur la description automatique des données*, Credoc,cr 13
- [58] L.Lebart, A.Morineau et M. Piron (1996), *Statistique exploratoire multidimensionnelle*. Dunod, France

Cet ouvrage est une synthèse des techniques factorielles utilisées en analyse de tableaux multiples( DACP,Statis,analyse procusteenne,regression multiple,analyse canonique,AFM analyse discriminate,...) . Une grande partie de cet ouvrage est consacré aux techniques de classification automatique. Une introduction à la validation de ces techniques constitue le dernier chapitre du livre (Tests ,simulation par le bootstrap ....). La présentation de ces différentes approches est faite avec un aperçu historique clair et sans ambiguïté.

- [59] A. Leclerc (1975), *L'analyse des correspondances sur juxtaposition de tableaux de contingence* RSA, **23**, p: 221–248.

Dans cet article sont étudiées les différentes propriétés concernant les sous tableaux de Burt.

- [60] I. C. Lerman (1970), *Les bases de la classification automatique*. Gauthiers Villard, Paris.
- [61] I .C.Lerman (1980), *Classification et analyse ordinaire des données*. Dunod, Paris.
- [62] I. C. Lerman (1980), *Analyse globale des éléments constitutifs d'une juxtaposition de tableaux de contingence* RSA XXVIII, **2**.
- [63] H. L'Hermier Des Plantes (1976), *Structuration des Tableaux A Trois Indices de la Statistique (STATIS)*. Thèse de l'Université de Montpellier. France.

- [64] G. Mac Lachlan (1982), *he classification and mixture maximum likelihood approach to cluster. analysis*.Hand Book of statistical ,Vol **2**, Krishnaiak. P.R and Kanak eds, pp: 199–208, Amesterdam, Holland.

- [65] J. B. Macqueen (1967), *Some methods for classification and analysis multivariate observations*. Proc. Sym . Math Statist and Probab. (5th), Berkely, **1**, pp: 281–297.

Dans cet article est introduite la méthode de classification par des partitions dite : méthode des K-means.

- [66] A. Morinneau, A.E Sammartino, M.Getter-Summa, C.Pardoux (1994), *Analyse des données et modélisations des séries chronologiques*. RSA, **42**, (4), pp: 61–81.

- [67] F. Murtagh (1985), *Multidimensionnal clustering algorithms* COMPSTAT lectures **4**, Physica Verlag, Vienna.
- [68] N. Naab(2001), *Analyse factorielle pour analyser les éléments constitutifs d'une structure de tableaux multiples*, thèse de magister, USTHB
- [69] J.P. Nakache (1973), *Influence du codage des données en analyse factorielle des correspondances*. RSA, **21**, N°2.
- [70] A. Nishisato (1980), *Analysis of categorial data, dual scaling and its applications*. Université de Toronto, Press.
- C'est un ouvrage où est développé l'analyse des correspondances multiples sous le nom de dual scaling. D'autres auteurs se sont penchés sur les techniques d'analyse de tableaux multiples et ont donné des extensions à l'AFM sous des noms tels que: Homogeneity analysis (J De Leuvv,...).
- [71] C. Perruchet (1983), *Une analyse bibliographique des épreuves de classificabilité en analyse des données*. RSA Vol. **8**, p: 18–41.
- [72] B.C. Petters and H.F Waker (1978), *An iterative procedure for obtaining maximum estimate of the parameters for a mixture of normal distributions*. SIAM, JAM, **35**, p : 362–378.
- [73] B.C. Petters and H.F Waker (1978), *The numerical evaluation of the maximum likelihood mixture proportions*. SIAM, JAM, **35**, p: 447–462.
- [74] P. Orlick et H Terao (1992) : *Arrangements of hyperplanes*, Springer Verlag
- [75] JO. Ramsay (2000), *Functional compment of variation in hardwritng*. JASA, serie B, **95**, p: 9–15.
- [76] A. Rebbouh (1983), *Classification automatique des données avec erreurs de mesure et application aux données médicales*. Thèse de Magister Université des sciences d'Alger.
- [77] A. Rebbouh (1993), *Classification des données avec erreurs de mesure: modèle gaussien*. IFCS, Paris
- [78] A. Rebbouh et S.Djemai (1996), *Classification d'objets aléatoires*. Prépublication, Institut de mathématiques université des sciences d'Alger, N° 183/36/13.
- [79] A. Rebbouh (2005), *Clustering the constitutive element of a structure of disjonctif table data*. 37èmes journées de statistique de la SFDS, juin, Pau, France.
- [80] A. Rebbouh (2006), *Clustering the constitutive element of a structure of measuring table data*. Communications in statistics simulation and computation, **35**, 2, p: 751–764.
- [81] A. Rebbouh (2007), *Classification d'objets aléatoires, deuxièmes journées de Statistiques appliquées et théoriques de Biskra, avril 2007*.

- [82] A. Rebbouh (2008) , *Clustering the constitutive element of a structure of binary data tables*  
à paraître dans JAMDS # Journal of Applied mathematics and Decision Science#
- [83] R.A. Renner and H.F. Walker(1984) , *Mixture density maximum likelihood and the em algorithm*, Dekker, N.Y.
- [84] M. Roux (1985) , *Algorithmes de classification*. Masson, Paris.
- [85] F. H. Ruspini (1969) , *A new approach to clustering*. Infor and control, **15**, p: 22–32.
- [86] G. Saporta (1981) , *Méthodes exploratoires d'analyse de données temporelles*. Thèse de l'Université de Paris VI.
- [87] G. Saporta (1990) , *Probabilités, analyse des données et statistique*. Edition technip-Paris.
- [88] A. Schroeder (1976) , *Analyse d'un mélange de distributions de probabilités de même type*. RSA, Vol **24** ,n° 1.
- [89] A.J. Scott and M.J.Symons (1971) , *Clustering methods based on likelihood ratio*. Biometrika, **27**, p: 387–397.
- [90] C E Shannon (1948) , *A mathematical theory of communication* Bell System technical journal, **27**(3) p: 349–423 & 623–659.
- [91] P H. Sneath & RR Sokal (1973) , *Numerical taxonomy*. Freeman, NY  
Livre de base et de référence sur le clustering.
- [92] M. Symons (1981) , *Clustering criteria and multivariate normal mixture*. Biometrika, **37**, p: 37–43.
- [93] PH Tassi et S Legait (1990) , *Théorie des probabilités, en vue des applications statistiques*. Edition technip, Paris.
- [94] N.Tennenhaus et Prieuret (1974) , *Analyse des séries chronologiques multidimensionnelles*. Revue RAIRO, **2**, p: 5–16.
- [95] N.Tennenhaus et F.W Young (1985) , *An analysis and synthesis of multiple correspondance analysis, optimal scaling, dual scaling homogeneity analysis and other methods for quantifying categorical multivariate data* Psychometrika, **50**, p: 91–119.
- [96] R.L. Thorndike (1953) , *Who belongs in the family*. Psychometrika, **18**, p: 267–276.  
Thorndike a introduit la notion de seuils ou protection destinée à modifier éventuellement le nombre de classes. Ses travaux ont été le prélude à la mise au point de plusieurs algorithmes de classification.



- [97] L.R.Tucker (1964), The extension of factor analysis to three -dimensional matrices, Harris.(ed), Univ. of Wisconsin, Press p:109-127
- [98] L.R. Tucker (1966), *Some mathematical notes on three mode factor analysis* Psychometrika, **31**, p: 231–279.  
Cet article est un exposé synthétique des diverses approches de l'AFM.
- [99] B.Van Cutsem (eds) (1989), *Classification and dissimilarity analysis*. Springer - Verlag, N.Y.
- [100] Ventsel (1973), *Théorie des probabilités*. Edition MIR, Moscou.
- [101] PGN.Von der Heidjer (1987), *Correspondance analysis longitudinal categorial data*. DSWO Press
- [102] J.H.Wolfe (1970), *Pattren clustering by multivariate mixture analysis Multivarariate*. Behavioral research, **5**, pp: 329–350.

# 1 Appendix

**Program of the algorithm that computes the ascending hierarchical of elements of a categorical data structure written under excel VBA.**

```

Sub test()
Dim matin(6, 6) As Variant
Dim matrixM(6, 6) As Variant
Dim matrixT(6, 6) As Variant
Dim vect(6) As Variant
Dim delta(6, 6) As Variant
Dim vectmin(6) As Variant
Dim delta2(6, 6)
size_matin = 6
Sheets("F1").Select
' vectmin contains the minimum value and the corresponding two objects
' input data matrix should be entered into the sheet F1
For i = 1 To size_matin
For j = 1 To size_matin
matin(i, j) = Cells(i, j)
Next
Next
' computation of matrix M
nl = 0
For l2 = 1 To size_matin
For l1 = 1 To size_matin
matrixM(l2, l1) = -Log(matin(l2, l1)) / Log(2)
Next
Next
' all the intermediate outputs will be put in the sheet F2
Sheets("F2").Select
' transposition of matrix M
For i = 1 To size_matin
For j = i + 1 To size_matin
X = matrixM(i, j)
matrixM(i, j) = matrixM(j, i)
matrixM(j, i) = X
Next
Next
' product of the transpose of matrix M by the input data matrix
For i = 1 To size_matin
For j = 1 To size_matin
s = 0
For k = 1 To size_matin
s = s + matin(i, k) * matrixM(k, j)
Next

```

```

matrixT(i, j) = s
Next
Next
nligne = size_matin
Do While nligne >= 3
mindelta = 100
For l1 = 1 To nligne
For l2 = l1 + 1 To nligne
x1 = matrixT(l1, l2)
x2 = matrixT(l2, l1)
x3 = matrixT(l1, l1)
x4 = matrixT(l2, l2)
delta(l2, l1) = x1 + x2 - (x3 + x4)
If mindelta > delta(l2, l1) Then
mindelta = delta(l2, l1)
vectmin(2) = l1
vectmin(3) = l2
End If
delta(l1, l2) = delta(l2, l1)
Next
delta(l1, l1) = 0
Next
vectmin(1) = mindelta
nl = nl + 1
' reconstruction of the matrix delta
lig = 0
For i = 1 To nligne
If i <> vectmin(2) And i <> vectmin(3) Then
lig = lig + 1
col = 0
For j = 1 To nligne
If j <> vectmin(2) And j <> vectmin(3) Then
col = col + 1
delta2(lig, col) = delta(i, j)
End If
Next
End If
Next
lig = nligne - 1
col = 0
For i = 1 To nligne
If i <> vectmin(2) And i <> vectmin(3) Then
If delta(vectmin(2), i) < delta(vectmin(3), i) Then
min1 = delta(vectmin(2), i)
Else

```

```
min1 = delta(vectmin(3), i)
End If
col = col + 1
delta2(lig, col) = min1
delta2(col, lig) = min1
End If
Next
delta2(lig, lig) = 0
nligne = nligne - 1
For i = 1 To nligne
For j = 1 To nligne
matrixT(i, j) = 0
delta(i, j) = 0
If i <= nligne And j <= nligne Then
matrixT(i, j) = delta2(i, j)
End If
Next
Next
Loop
Sheets("F1").Select
End Sub
```